

# Author's Accepted Manuscript

Robust visual tracking using structural region hierarchy and graph matching

Yi-Zhe Song, Chuan Li, Liang Wang, Peter Hall, Peiyi Shen

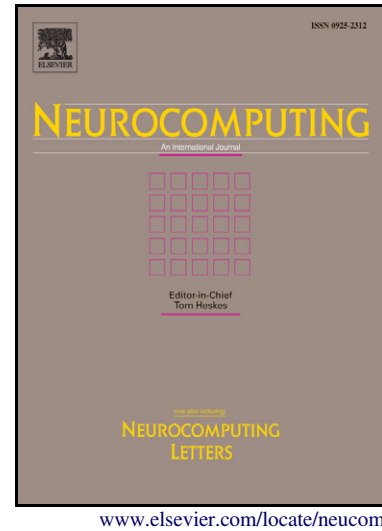
PII: S0925-2312(12)00195-6  
DOI: doi:10.1016/j.neucom.2011.11.030  
Reference: NEUCOM12640

To appear in: *Neurocomputing*

Received date: 23 June 2011  
Revised date: 15 November 2011  
Accepted date: 19 November 2011

Cite this article as: Yi-Zhe Song, Chuan Li, Liang Wang, Peter Hall, Peiyi Shen, Robust visual tracking using structural region hierarchy and graph matching, *Neurocomputing*, doi:10.1016/j.neucom.2011.11.030

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting galley proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



# Robust Visual Tracking using Structural Region Hierarchy and Graph Matching

Yi-Zhe Song<sup>a,b</sup>, Chuan Li<sup>b</sup>, Liang Wang<sup>c</sup>, Peter Hall<sup>b</sup>, Peiyi Shen<sup>d</sup>

<sup>a</sup>*School of Electronic Engineering and Computer Science, Queen Mary, University of London, London E1 4NS*

<sup>b</sup>*Department of Computer Science, University of Bath, Bath, BA2 7AY, UK*

<sup>c</sup>*National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences*

<sup>d</sup>*XiDian University, ShaanXi, P.R.China, 710071*

---

## Abstract

Visual tracking aims to match objects of interest in consecutive video frames. This paper proposes a novel and robust algorithm to address the problem of object tracking. To this end, we investigate the fusion of state-of-the-art image segmentation hierarchies and graph matching. More specifically: (i) we represent the object to be tracked using a hierarchy of regions, each of which is described with a combined feature set of SIFT descriptors and color histograms; (ii) we formulate the tracking process as a graph matching problem, which is solved by minimizing an energy function incorporating appearance and geometry contexts; and (iii) more importantly, an effective graph updating mechanism is proposed to adapt to the object changes over time for ensuring the tracking robustness. Experiments are carried out on several challenging sequences and results show that our method performs well in terms of object tracking, even in the presence of variations of scale and illumination, moving camera, occlusion, and background clutter.

**Keywords:** Tracking, Image Hierarchy, Graph Matching

---

## 1. Introduction

Fast and reliable object tracking is a prerequisite for a variety of vision applications such as monitoring and surveillance, robotic navigation, human computer interaction, etc. Although it has been one of the most active research topics over the past three decades in the field of computer vision, visual tracking remains a difficult problem due to a number of challenging factors, e.g., noise, occlusion, varying viewpoints, background clutter, illumination changes, and object appearance changes due to motion and articulation. There are numerous frameworks for tracking in the recent literature. For a comprehensive survey, see [1, 2].

In this paper, we are interested in investigating the effective fusion of structural region hierarchy and graph matching for the object representation and tracking task. Our work is motivated by recent developments in hierarchical image representations [3, 4, 5, 6] and their success in computer vision. In particular, a very latest study on filtering hierarchical image description [7] produces hierarchies that are not only accurate in terms of structural decompositions, but also manageable in terms of computational complexity. In summary, we treat the object to be tracked as a structural hierarchy of regions, and represent it as a relational graph. Each node (region) in the graph is augmented with an associated feature set consisting of local invariant features and color histograms, which is expected to be stable to cope with significant amount of image variability (such as appearance changes due to viewpoint and lighting). The tracking problem can be accordingly cast as a graph correspondence process. An adaptive graph updating scheme

is also proposed to account for object appearance variation and deformation. Experimental results on several challenging sequences demonstrate the effectiveness of our method.

Major contributions of this paper can be summarized as follows: i) We exploit the use of structural region hierarchy for the tracking problem. Although the representation of hierarchical regions has been widely studied and used in image segmentation, to the best of our knowledge, it seems that object tracking using structural region hierarchy has not yet been explored. ii) We propose an effective combination of image hierarchy description and graph matching, which leads to a robust object tracking algorithm. We show how this general principle can be effectively implemented in a graph matching and updating framework. iii) We perform a series of experiments on real-world videos, from which it is demonstrated that our method obtains both robustness to handle different challenging factors and flexibility to track different types of objects.

The remainder of this paper is organized as follows. Section 2 briefly reviews related work. The proposed tracking method is described in Section 3. Experiments and results are provided and analyzed in Section 4, prior to discussion and conclusion in Section 5.

## 2. Related Work

In this section, we restrict to three categories of recent methods related to our research, which are briefly reviewed as follows, respectively.

**Object Representation and Tracking:** Visual tracking is to continuously determine the position of the object in an image sequence against

dynamic scenes [1]. For example, tracking is cast as a sparse approximation problem in a particle filter framework, through a set of trivial templates [8]. In [9], a probabilistic framework is proposed for object tracking using a bag-of-pixels representation and a combination of a rigid registration between frames, a segmentation and online appearance learning.

There are two general ways to represent the objects: appearance-based (e.g., color histogram [10]) and shape-based (e.g., contour [11]). In [10], object appearance is modeled with a mixture of a fixed number of color-spatial Gaussians. In contour tracking, the object is represented as a contour, whose shape and location's evolution over time can be recovered by level sets [11] for example. Appearance-based representations often ignore structural information of the object, while shape-based ones neglect the object appearance. Neither of them is consistently satisfactory in the presence of dramatic intensity or color changes, object deformation and background clutter.

**Feature Detectors and Local Descriptors:** Detection of interest points and local invariant features constitute the basis for some vision tasks such as object recognition [12] and stereo matching [13]. Most algorithms start with interest point detection, followed by the local descriptor computation.

There have been several methods for interest point detection, e.g., difference-of-Gaussian [12], maximally stable extremal regions (MSER) [13], and Harris-Affine and Hessian-Affine corners [14]. The descriptors are usually designed to be invariant to lighting, scale and rotation changes, e.g., very popular SIFT [12]. See [15] for a comparative study on different descriptors. Recently, Grabner et al. [16] propose to speed up the computation of SIFT

using integral images. As a variant of SIFT, SURF [17] shares its distinctiveness and robustness, but has much faster computation speed.

**Image Representation via Structure Hierarchy:** Hierarchical image descriptions have been recognized as being valuable. As a well established example, the work of [4] uses morphological operators to generate a tree rooted around gray level extrema. Paris and Durand [6] use the notion of stability in feature space to produce a hierarchical description. The connected segmentation tree (CST) [5] takes into account the photometric properties, spatial organization and object structure to yield semantically meaningful hierarchies. The relatively successful boundary detector to date is the probability of boundary (Pb) maps [18]. The latter, i.e., global-Pb, leads to the state-of-the-art region hierarchies [3].

Ideally, hierarchies<sup>1</sup> reflect assemblies that comprise real objects, but in practice they can often be very large and complex. Some studies are carried out to simplify hierarchies, e.g., MSER simplifies a hierarchy in which thresholds make levels [13]. Recently, Song et al. [7] find semantic structures in image hierarchies using Laplacian graph energy as the complexity measure, which not only retains the semantics of the hierarchies but reduces their complexities by an order of magnitude.

**Tracking using Local Feature Descriptors:** Detection of interest points and local invariant features constitute the basis for some vision tasks

---

<sup>1</sup>Please note that the structure hierarchy herein means a representation based on hierarchical tree of image regions, which should be differentiated from the hierarchical representation of human body structure (e.g., body parts) in [19] and the hierarchy of multi-scale or multi-stage data representation and analysis strategy [20].

95 such as object recognition [12] and stereo matching [13]. The descriptors  
 96 are usually designed to be invariant to lighting, scale and rotation changes,  
 97 e.g., very popular SIFT [12]. See [15] for a comparative study on different  
 98 descriptors. Recently, several papers exploit the use of local invariant features  
 99 in the tracking problem. Ta et al. [21] propose an algorithm for continuous  
 100 image recognition and feature descriptor tracking. An integrated SIFT-based  
 101 mean shift algorithm is presented in [22]. In [23], tracking is based on motion  
 102 analysis of regional affine invariant features, and the object occupancy map  
 103 is updated according to the pixel motion consistency.

104 In addition to the use of local features, several algorithms explore their  
 105 relationships for the tracking task. Tang and Tao [24] present an attributed  
 106 relational graph (ARG) that incorporates distinctive SIFT features and their  
 107 relations for object tracking. In [25], a generative model consisting of con-  
 108 sistent and random components is used to depict the relationship between  
 109 local feature motions and object global motion. The latter is estimated in  
 110 term of maximum likelihood of local SURF feature observations.

111 Our work also benefits from modeling objects in terms of attributed rela-  
 112 tional graphs. However, our graphs are built directly from structural region  
 113 hierarchies other than that made from local features themselves [24, 25], and  
 114 hence are able to capture regional photometric and geometric properties to fa-  
 115 cilitate robust tracking via a novel graph matching and updating mechanism.  
 116 More specifically, 1) Our object representation is a hierarchy of regions, each  
 117 of which is denoted as a combined feature set of SIFT descriptors and color  
 118 histograms, and our relational graph is constructed directly from such hierar-  
 119 chy. However, in [24], observations in each frame are only the extracted SIFT

120 features, which are used as vertices to generate attributed relational graph.  
 121 In addition to different graph formation and updating manners, correlation-  
 122 based ARG matching in [24] could be far from robust compared with our  
 123 graph correspondence method. 2) The component tree is used in [26], but  
 124 it is much different from our hierarchical tree of regions, in terms of both  
 125 construction manners and semantics of leaves. Our tracking is implemented  
 126 via graph matching of hierarchy-derived graphs. In [26], to identify the best  
 127 fit to the input MSER, feature vectors (that are built for each of the ex-  
 128 tremal regions of the component tree) are compared to choose the region  
 129 with the smallest weighted Euclidean distance as the tracked MSER. 3) For  
 130 [25], individual SURF features based graph matching is used to estimate  
 131 affine motion for constraining the search region. However, ours is based on  
 132 feature set correspondences which results directly in tracking.

### 133 3. The Method

134 Our tracking method consists of two major components: building hier-  
 135 archical object descriptions and tracking via graph matching and updating.  
 136 Figure 2 offers a diagram of the proposed tracker, whose sub-components will  
 137 be discussed in detail in Sections 3.1 and 3.2.

#### 138 3.1. Object Representation via Structural Region Hierarchy

139 **From Contours to Structural Hierarchies:** At first, the method de-  
 140 scribed in [3] is chosen to construct a hierarchy of regions which works on  
 141 the output of any contour detector. In this paper, we use the state-of-the-  
 142 art ‘global-Pb’ contour detector [18] as its input. Briefly, Arbelaez et al.  
 143 [3] used oriented watershed transform (OWT) on the topographic surface



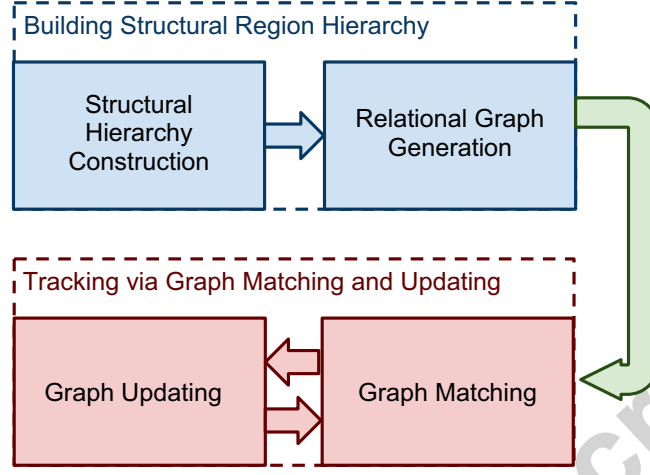


Figure 1: A system diagram of the proposed tracking framework

144 defined by  $O(x, y) = \max_{\theta} O(x, y, \theta)$ , where  $O(x, y, \theta)$  is the ‘global-Pb’ con-  
 145 tour response. The catchment basins of the minima  $\mathcal{P}_0$  provide the region of  
 146 the finest partition and corresponding watershed arcs  $\mathcal{K}_0$ . The strength of  
 147 the boundaries  $O(x, y, \theta)$  is then transferred to the locations  $\mathcal{K}_0$  by approx-  
 148 imating the watershed arcs with line segments and weighting each point in  
 149  $\mathcal{K}_0$  by the value of  $O(x, y, \theta)$  in direction  $\theta$  determined by the orientation of  
 150 the corresponding line segments. Afterwards, the region hierarchy is built  
 151 by a greedy graph-based region merging algorithm which merges the most  
 152 similar regions at one time. Figure 2(c) shows an example of the hierarchy,  
 153 which is visualized as an unrametric contour map (UCM) image obtained  
 154 by weighting each boundary between two regions by its scale of disappear-  
 155 ance. Locally, regional features describe the object parts’ details; globally,  
 156 the relations between regions encode the object structure. It is anticipated  
 157 that such a hierarchical representation allows parts of the object to have  
 158 different motions, hence is more flexible to represent objects and to handle

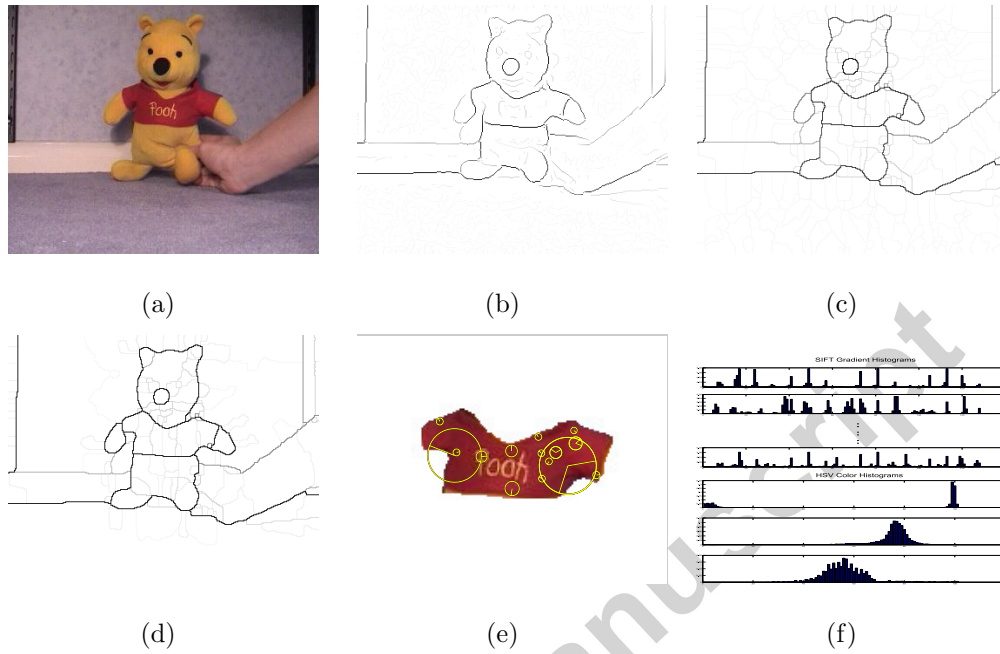


Figure 2: An example illustrating the process of ‘from image to hierarchy of regions’: a) Original image, b) Pb contours, c) original hierarchy of regions, d) filtered hierarchy of regions, e) an example region imposed by the detected SIFT points, and f) its SIFT and HSV histogram features.

159 object deformation and appearance changes.

160 Unfortunately, the hierarchy produced by [3] contains far more data than  
 161 is required for an efficient description, which makes it impractical in real  
 162 applications such as object tracking. There are significant practical advan-  
 163 tages to simplify hierarchical descriptions, e.g., gains in memory efficiency  
 164 and computational speed, which are all important factors for fast tracking.  
 165 This work uses the method described in [7] to filter complex hierarchies into  
 166 simpler ones. The reduced hierarchies preserve the semantic interpretation  
 167 in terms of objects and object parts; the number of levels of the reduces  
 168 hierarchies are typically an order of magnitude less than the original.

169 The principle in solving the problem of filtering hierarchy is to choose  
 170 those levels that are lower in complexity than their neighbors. The Laplacian  
 171 graph energy [27] is used to measure the complexity. Let  $G$  be a graph of  
 172  $n$  vertices and  $m$  edges, i.e., a  $(n, m)$ -graph. Let  $A$  and  $D$  be its adjacency  
 173 matrix and corresponding degree matrix. Then  $L = D - A$  is the graph  
 174 Laplacian. The Laplacian graph energy is defined by

$$\mathcal{LE}(G) = \sum_{j=1}^n |\mu_j - 2m/n|$$

175 where  $\mu_j$  is the eigenvalues of  $L$  and  $2m/n$  is the average vertex degree. In [7],  
 176 the affinity matrix is specially defined as  $A = \{a_{ij} | a_{ij} = \exp(-w_{ij}/w_{max})\}$   
 177 where  $w_{ij}$  is the average boundary strength between regions  $i$  and  $j$ , and  
 178  $w_{max}$  is a decay factor set to the maximum over all  $w_{ij}$ . Another extension,  
 179 called the *component-wise* Laplacian graph energy (cLGE) is introduced in  
 180 [7]. For a graph with  $k$  disconnected components, the cLGE is defined as

$$c\mathcal{LE}(G) = K \sum_{i=1}^k \frac{\mathcal{LE}(G_i)}{|n_i|}$$

181 in which  $G_i$  is  $i$ th connected component (or sub-graph) of  $|n_i|$  nodes, and  $K$   
 182 is the number of nodes in the whole graph. The cLGE at every level in the  
 183 hierarchy is computed independently using graphs built from the primitives  
 184 at the lowest level. At the bottom level of the hierarchy, each primitive is an  
 185 1-node sub-graph on its own, whereas the top level forms a single connected  
 186 graph. At intermediate levels, as segmentations become coarser, subgraphs  
 187 are merged to create larger ones, and so the number of disconnected compo-  
 188 nents will fall. cLGE for the level as a whole can rise or fall, depending on  
 189 the way these primitives are connected. So only those levels, at which cLGE

is locally minimal, are kept in the filtered hierarchy [7]. See Figure 2(d) for an example.

**From Structural Hierarchy to Relational Graph:** Graphs are a versatile and flexible representation formalism useful for a range of problems in information processing. We finish building structural regional hierarchies by turning them into attributed relational graphs  $G'$ . We first remove possibly repetitive nodes residing on different levels from the filtered hierarchy via active search, and the resulting tree is still a rooted connected tree, whose nodes represent connected regions within the input and edges define an inclusion relationship between connected regions. Then, for the links between a parent and multiple children, we assign the weights with the corresponding size ratios between the parent and children nodes. It is also assumed that only children of a common parent can be adjacent, similar to the CST [5]. Accordingly, we assign the weight to each link between intra-level nodes as their similarity.

For each node (region) of the graph, we augment it with a feature vector containing its local geometric and photometric information. In particular, we describe each region with a SIFT descriptor [12, 17, 15] which achieved great success in recent applications due to their appealing characteristics. In addition, a color histogram is also used to encode regional photometric content, a notion that has been widely used in object representation [10]. Each SIFT [12] feature is described as  $\mathbf{f}_s = \{p, s, o, \mathbf{h}_g\}$ , where  $p$  is the 2D position of the key point,  $s$  is the scale,  $o$  is the orientation of the main gradient within the local region, and  $\mathbf{h}_g$  is the gradient orientation distribution quantized into 128 bins. For color histograms, we use the HSV color

space and quantize each channel into 100 bins, leading to a concatenated  
 300-dimensional feature vector  $\mathbf{f}_c = \{\mathbf{h}_h, \mathbf{h}_s, \mathbf{h}_v\}$ . For the sake of efficiency,  
 all SIFT features and color histograms are computed incrementally during  
 creation of the hierarchy.

### 3.2. Tracking via Graph Matching and Updating

**Graph Matching:** We formulate the tracking process into a graph  
 matching and updating scheme. Given the graph representation  $G'$  of the  
 object to be tracked in current frame  $t$  and the (potentially larger) candidate  
 graph representation  $G''$  in next frame  $t + 1$ , the task of object tracking is  
 naturally cast as a graph matching problem, i.e.,  $\mathcal{M} : G' \rightarrow G''$ . In this  
 paper, we borrow the idea from the method in [28]. For description con-  
 venience, let  $P'$  and  $P''$  be the sets of features to be matched (in our case,  
 ‘features’ mean regional nodes in  $G'$  and  $G''$ ). Let  $R \subseteq P' \times P''$  the set of  
 potential assignments. A matching configuration between the two sets can  
 be denoted as a binary vector  $\mathbf{x} \in \{0, 1\}^R$ . A correspondence  $a \in R$  indexes  
 an entry  $x_a$  in  $\mathbf{x}$ , and it is active if  $x_a = 1$  and inactive otherwise. An energy  
 function  $E(\mathbf{x})$  is defined to formulate the matching task as minimization of  
 $E(\mathbf{x})$ . The uniqueness constraint is enforced via the set  $M$

$$M = \{\mathbf{x} \in \{0, 1\}^R \mid \sum_{a \in R(p)} x_a \leq 1, \forall p \in P\} \quad (1)$$

where  $P = P' \cup P''$  and  $R(p)$  is the set of correspondences involving feature  
 $p$ . The problem is to find a configuration  $\mathbf{x}$  that can minimize  $E(\mathbf{x})$ . In  
 our work, we define our energy as a weighted sum of only two energy terms  
 corresponding to appearance and geometry properties:

$$E(\mathbf{x}) = \lambda^{app} E^{app}(\mathbf{x}) + \lambda^{geo} E^{geo}(\mathbf{x}). \quad (2)$$

The term  $E^{app}(\mathbf{x})$  favors correspondences between features having similar appearance, which is defined as the sum of unary terms:  $E^{app}(\mathbf{x}) = \sum_{a \in R} \theta_a^{app} x_a$ . For an assignment  $a = (p', p'') \in R$ ,  $\theta_a^{app}$  is the distance between appearance descriptors computed at  $p'$  and  $p''$  in the respective sets. In this work, Chamfer distance [29] is used to compute the dissimilarity between SIFT-features sets. Assume that two regional feature sets are  $\mathcal{U} = \{\mathbf{u}_i\}_{i=1}^{N_f}$  and  $\mathcal{V} = \{\mathbf{v}_j\}_{j=1}^{M_f}$ , where  $N_f$  and  $M_f$  denote the number of SIFT features in the two associated regions, the symmetric Chamfer distance is defined as

$$d_{cham}(\mathcal{U}, \mathcal{V}) = \frac{1}{N_f} \sum_{\mathbf{u}_i \in \mathcal{U}} \min_{\mathbf{v}_j \in \mathcal{V}} \|\mathbf{u}_i - \mathbf{v}_j\| + \frac{1}{M_f} \sum_{\mathbf{v}_j \in \mathcal{V}} \min_{\mathbf{u}_i \in \mathcal{U}} \|\mathbf{v}_j - \mathbf{u}_i\|.$$

237 For color histogram features of the two regions, we compute their dissimilarity  
238 using the  $\chi^2$  distance, i.e.,

$$d_{\chi^2}(\mathbf{h}_1, \mathbf{h}_2) = \frac{1}{2} \sum_i \frac{(\mathbf{h}_1(i) - \mathbf{h}_2(i))^2}{\mathbf{h}_1(i) + \mathbf{h}_2(i) + \epsilon} \quad (3)$$

239 where the introduction of a non-zero constant  $\epsilon$  is just to avoid “division by  
240 zero” in practice. A weighted sum of normalized  $d_{cham}$  and  $d_{\chi^2}$  is used as the  
241 appearance energy function.

242 The term  $E^{geo}(\mathbf{x})$  measures geometric compatibility of correspondences  
243 only for neighboring features, which is defined as  $E^{geo}(\mathbf{x}) = \sum_{(a,b) \in N} \theta_{ab}^{geo} x_a x_b$ ,  
244 where the constraint set  $N$  is a neighbor system consisting of all correspon-  
245 dence pairs defined over neighboring features. This geometry energy term is  
246 computed using the same method described [28]. Feature correspondence is

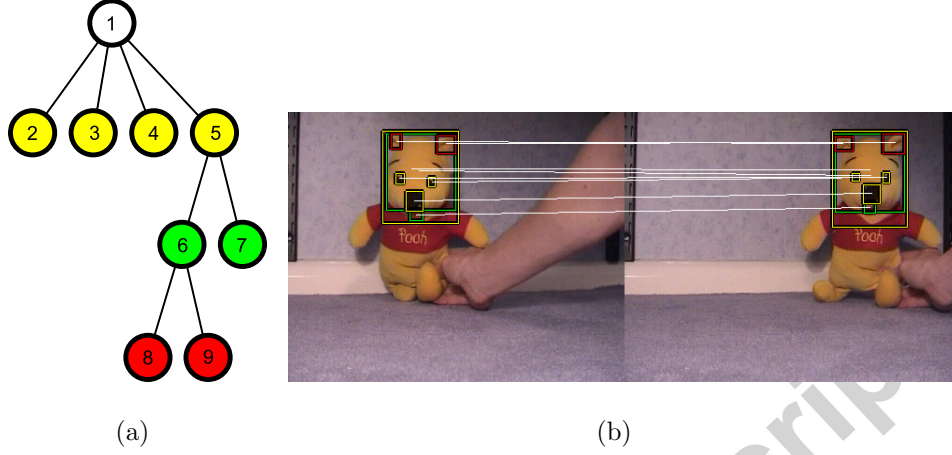


Figure 3: An example of graph matching: a) illustration to the hierarchy tree of regions, and b) the matched regions between two frames, where the colored nodes in (a) correspond to the rectangle regions with the same color.

247 finally written as

$$\min_{\mathbf{x}} E(\mathbf{x}|\bar{\theta}) = \sum_{a \in R} \bar{\theta}_a x_a + \sum_{(a,b) \in N} \bar{\theta}_{ab} x_a x_b \quad (4)$$

248 which is naturally referred to as graph matching if features  $P'$  and  $P''$  are  
 249 viewed as vertices of the two graphs. Pairwise term  $\bar{\theta}_{ab} x_a x_b$  with  $a = (p', p'')$ ,  
 250  $b = (q', q'')$  encodes compatibility between edges  $(p', q')$ ,  $(p'', q'')$  of the first  
 251 and second graph, respectively, while unary term  $\bar{\theta}_a x_a$  measures similarity  
 252 between vertices  $p', p''$ . We adopt the dual decomposition approach for this  
 253 optimization problem. See [28] for more details. Figure 3 shows an example of  
 254 graph matching of the hierarchy tree of regions, where each region is plotted  
 255 as its corresponding bounding box for simplicity.

256 **Graph Updating Mechanism:** One of the key problems in tracking is  
 257 how to dynamically update the object model to accommodate the object's  
 258 appearance and structure changes over time. We update the object graph

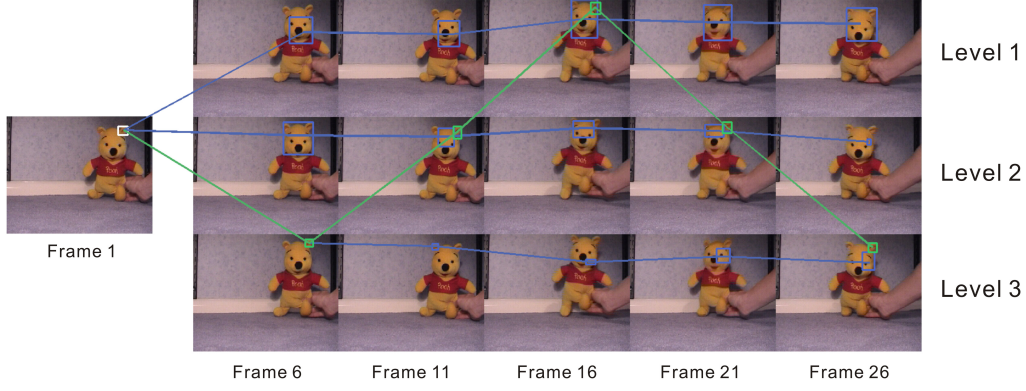


Figure 4: Tracking independent levels in the structural hierarchy (blue line) produces inaccurate results; robustness of tracking is improved by tracking structural region hierarchy instead (green line).

model using a collection of hierarchy of regions and their regional features. Regions with incompatible motion are discarded, while newly appearing object regions are incorporated into the graph representation. In particular, forward and backward matching are combined to improve the tracking robustness.

To track a video, we first build a graph  $G$  for each frame. The target object,  $g_1$ , is then initialized as a subgraph of  $G_1$  in the first frame. We set ‘key frames’ for every  $k$  frame, and name the other frames the “intermediate history frames”. Given  $g_t$  for key frame  $t$ , its updated graph in the next key frame contains two parts:  $g_{t+k}^f$ , a forward graph,  $g_{t+k}$ , computed from matching  $g_t$  to  $G_{t+k}$ ; and  $g_{t+k}^b$ , a backward graph computed from a sequence of matchings from  $G_{t+k}$  to the intermediate history frames  $[G_{t+1}, G_{t+2}, \dots, G_{t+k-1}]$ .

In details, we first perform a forward matching from  $g_t$  to  $G_{t+k}$  and use



273 RANSAC to reject outliers. The resulting forward graph  $g_{t+k}^f$  only contains  
 274 region features that already exist in frame  $t$ . To incorporate potentially new  
 275 object regions, we first define a candidate graph  $g_{t+k}^n$  that is connected by  
 276 regions within a short distance from  $g_{t+k}^f$ :

$$g_{t+k}^n = \forall r \in G_{t+k} : \min(|r - g_{t+k}^f|) < \gamma \quad (5)$$

277 A sequence of backward matchings is then performed between  $g_{t+k}^n$  and  
 278  $[G_{t+1}, G_{t+2}, \dots, G_{t+k-1}]$ . The confidence of a candidate region to be a true new  
 279 object region in  $g_{t+k}^n$  is evaluated by the number of times it is successfully  
 280 matched. Regions of high matching score are incorporated into the final  
 281 backward graph  $g_{t+k}^b$ :

$$g_{t+k}^b = \forall r \in g_{t+k}^n : \left( \sum_{i=1}^{k-1} \text{Score}(r, G_{t+i}) \right) > \lambda \quad (6)$$

282 The final updated graph  $g_{t+k}$  is the union of  $g_{t+k}^f$  and  $g_{t+k}^b$ . In practice,  $\gamma$   
 283 and  $\lambda$  are parameters that control the generosity of introducing new object  
 284 regions into the graph. We set  $k = 5$  to select the key frames. This value  
 285 works well for most videos and can be tuned for particularly slowly or fast  
 286 moving objects. In the updated graph, an edge is only kept if both the  
 287 parent and the child nodes are successfully matched. This means that  $g_{t+k}$   
 288 is not necessarily restricted to a complete subgraph of the original filtered  
 289 hierarchy and allows the existence of “isolated” nodes. In order to handle  
 290 occlusion, we also temporarily cease the updating procedure when forward  
 291 mapping rate falls below a threshold. Our global graph matching paradigm  
 292 naturally discovers the object when it reappears.

293 We use default parameters as in [28], except for the  $\lambda$  (Equation 6),  
 294 which is a new parameter introduced here and was set to be 2 in our ex-  
 295 periment, which means a new feature can only be added into the template  
 296 if it is matched to at least two of the five intermediate historic frames. For  
 297 RANSAC, we use the affine model as it provides efficient versatility in trans-  
 298 formation without introducing unnecessary complexity.

#### 299 4. Experiments and Results

300 In this section, we first provide a simple example illustrating the benefits  
 301 and effects of tracking using structural region hierarchies. We then concen-  
 302 trate on evaluating our tracker against general challenging tracking problems  
 303 such as illumination change, occlusion, moving background and so on — a  
 304 setup that was also used in the two state-of-the-art trackers we compare with.

305 In Figure 4, we illustrate different tracking results of the highlighted re-  
 306 gion (left ear of the bear) in frame 1: when it is tracked as a node in three  
 307 separate planar relational graphs (blue trajectories, one for each level) and  
 308 when it is tracked within a structural region hierarchy as proposed. It shows  
 309 that tracking independent levels of the structural hierarchy yields incorrect  
 310 results. First, higher level graphs, such as level 1 and 2, are usually com-  
 311 posed of large regions. The consequence is small objects, such as the ear of  
 312 the bear, will be lost of track because of not being detected in the first place.  
 313 Second, even on the lower level graphs (level 3) tracking can still fail because  
 314 the structural information is not sufficiently used to distinguish regions of  
 315 similar appearance. For example, on level 3 the tracker lost track from frame  
 316 11 onwards, where interestingly the left ear was matched to the right one

317 instead. However, the structural hierarchy encodes richer information and it  
 318 is able to pick out the correct matches across different levels.

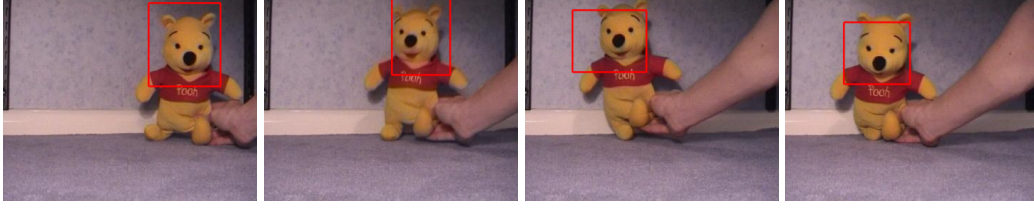
319 In order to evaluate the general tracking performance of our method, sev-  
 320 eral challenging sequences were chosen to carry out our experiments. These  
 321 sequences encapsulate moving objects of different types such as rigid cars and  
 322 articulated humans. In particular, they were recorded in indoor and outdoor  
 323 environments, including significant variations with respect to scale, illumi-  
 324 nation, moving camera, occlusion, pose and cluttered background. For each  
 325 sequence, we manually specified the tracking area, and plotted the global ob-  
 326 ject tracking results and local region tracking results using yellow and white  
 327 bounding boxes respectively.

328 We also compared our method with two state-of-the-art algorithms de-  
 329 scribed in [30] and [8]. The former [30] models the constantly changing fore-  
 330 ground shape as a small number of rectangular blocks with histogram-based  
 331 appearance representation, whose positions within the tracking window are  
 332 adaptively determined. The  $l_1$  tracker [8] describes each target candidate as  
 333 sparse representation in the space spanned by target and trivial templates,  
 334 and tracking is continued using a Bayesian state inference framework. Apart  
 335 from the first (the example Winne video used throughout), all sequences are  
 336 taken from VISOR<sup>2</sup>. For both of the two algorithms, we directly used their  
 337 published codes with default parameter settings.

338 The first sequence includes a hand-held moving Winne toy. It exhibits  
 339 variations in terms of scale, rotation and perspective effects. The tracking

---

<sup>2</sup>[http://imagelab.ing.unimore.it/visor/video\\_categories.asp](http://imagelab.ing.unimore.it/visor/video_categories.asp)



(a) Nejhum et al. 2008 [30]



(b) Mei and Ling 2009 [8]

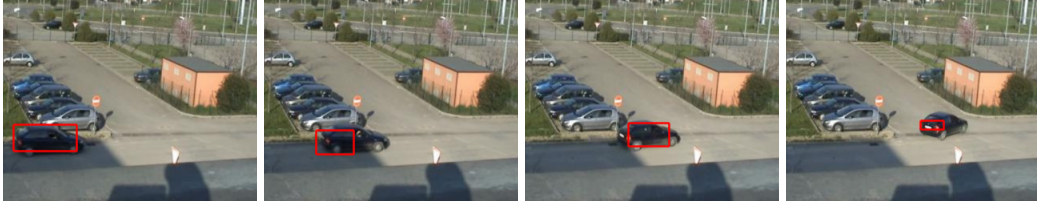


(c) Our method

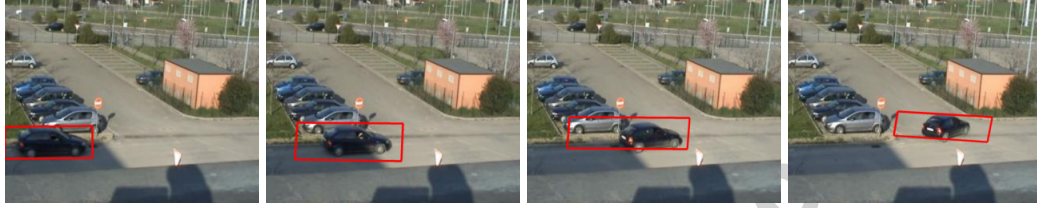
Figure 5: Tracking results on the ‘Winne’ sequence

340 results of several important frames are shown in Figure 5, together with cor-  
 341 responding results of the two compared methods. For this relatively simple  
 342 sequence, all three tracking results are satisfactory, however both the  $l_1$  algo-  
 343 rithm and the method of Nejhum et al. [30] tend to include more background  
 344 regions over time.

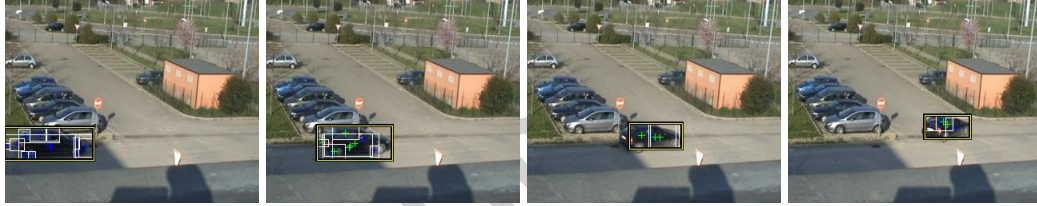
345 The second sequence exhibits a car driving away from a shadowed area.  
 346 It mainly includes variations in terms of scale (due to distance traveled) and  
 347 illumination (from shadow to shiny). Similarly, tracking results of four key  
 348 frames are provided in Figure 6, together with the corresponding tracking



(a) Nejhun et al. 2008 [30]



(b) Mei and Ling 2009 [8]



(c) Our method

Figure 6: Tracking results on the ‘Car’ sequence

349 results using the other two trackers. In this case, all three trackers follows  
 350 the general motion of the tracked object. However, the Nejhun et al. al-  
 351 gorithm and the  $l_1$  method tend to drift away, reflected by either shrunk  
 352 or enlarged templates. Clearly, our method performs best in terms of its  
 353 tracking accuracy.

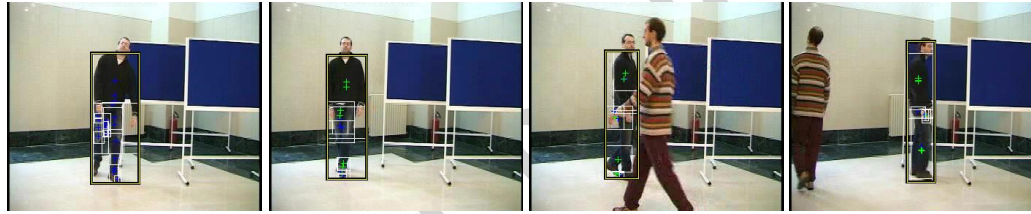
354 The third sequence includes two walking subjects, whose walking patterns  
 355 trigger occlusion. It is mainly used to examine the influence of occlusion and  
 356 pose changes on the tracking performance. The tracking results of several  
 357 important frames are shown in Figure 7, together with corresponding results



(a) Nejhum et al. 2008 [30]



(b) Mei and Ling 2009 [8]



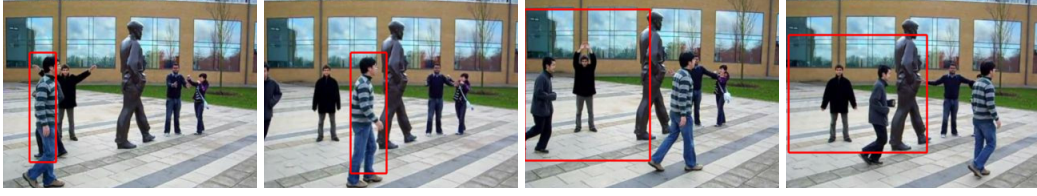
(c) Our method

Figure 7: Tracking results on the ‘Human’ sequence with occlusion

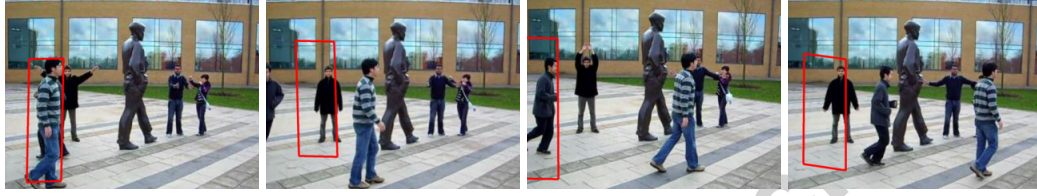
358 using the other two methods. Similarly, our method is still best in this case,  
 359 in terms of its accuracy and robustness.

360 The final sequence is the most challenging. It includes several people  
 361 performing different motions in a cluttered background, filmed by a moving  
 362 camera. It can be used to examine the effects of moving background, clutter  
 363 and occlusions on tracking. The results of representative frames are shown in  
 364 Figure 8, together with corresponding results using the two compared meth-  
 365 ods. For both the  $l_1$  and Nejhum et al. algorithms, either more backgrounds





(a) Nejhum et al. 2008 [30]



(b) Mei and Ling 2009 [8]



(c) Our method

Figure 8: Tracking results on the ‘Human’ sequence with moving camera

366 or the lost of tracking occur. However, our method does not suffer from this  
 367 at all.

368 In terms of computational complexity of our system, the total cost  $K$  is  
 369  $N_{frames} * (Cost_R + Cost_M)$ , where  $N_{frame}$  is the number of frame,  $Cost_R$  is  
 370 the cost for eliminating outliers, and  $Cost_M$  is the cost for the actual graph  
 371 matching. Due to our online graph updating mechanism,  $1/K$  of the cost is  
 372 spent on forward matching and the rest is spent on backward matching. Since  
 373 we use RANSAC to eliminate outlier matches, the associated complexity is  
 374  $O(T_{iter}(m^2 * N + 2 * N))$ . Here  $T_{iter}$  is the number of RANSAC iterations,  $m$

375 is the degrees of freedom in the estimated motion model and  $N$  is the number  
 376 of nodes in the graph. It is worth noting that graph matching is generally a  
 377 NP-hard problem and the decomposition process described in [28] is used to  
 378 for fast matching. Moreover, the Laplacian graph complexity analysis further  
 379 reduces the number of nodes in the graph by an order of magnitude, which  
 380 importantly enables our system to matching two frames within a second in  
 381 practice.

382 In summary, our method performs best on all these four sequences, in  
 383 terms of accuracy and robustness. In particular, our method can handle dif-  
 384 ferent challenging factors such as changes of scale and illumination, moving  
 385 camera, background clutter, and occlusion. On one hand, this benefits from  
 386 our hierarchical tree representation, which reflects structure and appearance  
 387 information of global object and its local parts (regions); on the other hand,  
 388 because the hierarchy description can be performed rapidly and reliably on-  
 389 line, the model of the object of interest can be adaptively updated over time  
 390 using our forward/background matching scheme, without suffering from the  
 391 problem of drifting that causes the loss of tracked objects.

## 392 5. Discussion and Conclusion

393 This paper has presented a novel and robust algorithm for visual tracking.  
 394 The key contribution of our method (and reason of its superior performance)  
 395 is the use of structural regional hierarchies to represent objects and a novel  
 396 graph matching and updating framework as the tracking mechanism.

397 Our tracker is robust owing to the rich underlying object descriptions  
 398 encoded as RAGs, which are tracked within a graph theoretic paradigm.



399 In terms of efficiency, building the RAG involves the use of several com-  
 400 monly practiced state-of-the-art techniques, where efficient implementations  
 401 are publicly available. Importantly, sizes of RAGs are reduced to a manage-  
 402 able level due to the hierarchy filtering technique used, which in turn largely  
 403 reduced the burden of the graph matching and updating mechanism. Track-  
 404 ing accuracy is also improved because of the hierarchical nature of our object  
 405 representations (Figure 4).

406 We have demonstrated the benefits of our method through a series of  
 407 qualitative results on real challenging sequences. The tracker is evaluated  
 408 against major challenging tracking problems such as illumination change,  
 409 occlusion, moving background and so on — a setup that was also used in the  
 410 two state-of-the-art trackers we compare with.

411 Finally, we also believe the idea of tracking RAG has important potential  
 412 benefits in higher-level tasks such as object matching and recognition across  
 413 videos.

#### 414 **Acknowledgment**

415 This work is in part supported by the Open Projects Program of Na-  
 416 tional Laboratory of Pattern Recognition, the Natural Science Basic Re-  
 417 search Plan in Shaanxi Province of China (2010JM8005) and the NSFC Grant  
 418 (61072105).

#### 419 **References**

- 420 [1] W. M. Hu, T. N. Tan, L. Wang, S. Maybank, A survey on visual surveil-  
 421 lance of object motion and behaviors, IEEE Transactions on Systems,  
 422 Man and Cybernetics - Part C 34 (3) (2004) 334–352.

- 423 [2] H. Yang, L. Shao, F. Zheng, L. Wang, Z. Song, Recent advances and  
424 trends in visual tracking: A review, *Neurocomput.* 74 (2011) 3823–3831.
- 425 [3] P. Arbeláez, M. Maire, C. Fowlkes, J. Malik, From contours to regions:  
426 An empirical evaluation, in: *International Conference on Computer Vi-*  
427 *sion and Pattern Recognition*, 2009, pp. 2294–2301.
- 428 [4] J. A. Bangham, R. W. Harvey, P. D. Ling, R. V. Aldridge, Morpholog-  
429 ical scale-space preserving transforms in many dimensions, *Journal of*  
430 *Electronic Imaging* 5 (1996) 283–299.
- 431 [5] N. Ahuja, S. Todorovic, Connected segmentation tree - a joint represen-  
432 tation of region layout and hierarchy, in: *International Conference on*  
433 *Computer Vision and Pattern Recognition*, 2008, pp. 1–8.
- 434 [6] S. Paris, F. Durand, A topological approach to hierarchical segmentation  
435 using mean shift, in: *International Conference on Computer Vision and*  
436 *Pattern Recognition*, 2007, pp. 1–8.
- 437 [7] Y. Z. Song, P. Arbelaez, P. M. Hall, C. Li, A. Balikai, Finding semantic  
438 structures in image hierarchies using Laplacian graph energy, in: *Euro-*  
439 *pean Conference on Computer Vision*, Vol. 4, 2010, pp. 694–707.
- 440 [8] X. Mei, H. B. Ling, Robust visual tracking using  $l_1$  minimization, in:  
441 *International Conference on Computer Vision*, 2009.
- 442 [9] C. Bibby, I. Reid, Robust real-time visual tracking using pixel-wise pos-  
443 teriors, in: *European Conference on Computer Vision*, 2008, pp. 831–  
444 844.

- 445 [10] H. Wang, D. Suter, K. Schindler, Effective appearance model and sim-  
446 ilarity measure for particle filtering and visual tracking, in: European  
447 Conference on Computer Vision, Vol. 3, 2006, pp. 606–618.
- 448 [11] Y. Shi, W. C. Karl, Real-time tracking using level sets, in: International  
449 Conference on Computer Vision and Pattern Recognition, Vol. 2, 2005,  
450 pp. 34–41.
- 451 [12] D. Lowe, Distinctive image features from scale-invariant key points, In-  
452 ternational Journal of Computer Vision 60 (2) (2004) 91–110.
- 453 [13] J. Matas, O. Chum, U. Martin, T. Pajdla, Robust wide baseline stereo  
454 from maximally stable extremal regions, in: British Machine Vision  
455 Conference, Vol. 1, 2002, pp. 384–393.
- 456 [14] K. Mikolajczyk, C. Schmid, Scale and affine invariant interest point  
457 detectors, International Journal of Computer Vision 60 (1) (2004) 63–  
458 86.
- 459 [15] K. Mikolajczyk, C. Schmid, Performance evaluation of local descriptors,  
460 IEEE Transactions on Pattern Analysis and Machine Intelligence 27 (10)  
461 (2005) 1615–1630.
- 462 [16] M. Grabner, H. Grabner, H. Bischof, Fast approximated SIFT, in: Asian  
463 Conference on Computer Vision, 2006, pp. 918–927.
- 464 [17] H. Bay, T. Tuytelaars, L. V. Gool, SURF: Speed up robust features, in:  
465 European Conference on Computer Vision, Vol. 1, 2006, pp. 404–417.

- 466 [18] M. Maire, P. Arbeláez, C. Fowlkes, J. Malik, Using contours to detect  
467 and localize junctions in natural images, in: International Conference  
468 on Computer Vision and Pattern Recognition, 2008, pp. 1–8.
- 469 [19] L. Raskin, M. Rudzsky, E. Rivlin, 3D human body-part tracking and  
470 action classification using a hierarchical body model, in: British Machine  
471 Vision Conference, 2009.
- 472 [20] B. Stenger, A. Thayananthan, P. Torr, R. Cipolla, Model-based hand  
473 tracking using a hierarchical Bayesian filter, *IEEE Transactions on Pat-*  
474 *tern Analysis and Machine Intelligence* 28 (9) (2006) 1372–1384.
- 475 [21] D.-N. Ta, W.-C. Chen, N. Gelfand, K. Pulli, SURFTrac: Efficient track-  
476 ing and continuous object recognition using local feature descriptors, in:  
477 International Conference on Computer Vision and Pattern Recognition,  
478 2009, pp. 2937–2944.
- 479 [22] H. Zhou, Y. Yuan, C. Shi, Object tracking using SIFT features and  
480 mean shift, *Computer Vision and Image Understanding* 113 (3) (2009)  
481 345–352.
- 482 [23] S. Tran, L. Davis, Robust object tracking with regional affine invariant  
483 features, in: International Conference on Computer Vision, 2007, pp.  
484 1–8.
- 485 [24] F. Tang, H. Tao, Object tracking with dynamic feature graph, in: Joint  
486 International Workshop on Visual Surveillance and Performance Evalu-  
487 ation of Tracking and Surveillance, 2005, pp. 25–32.

- 488 [25] W. He, T. Yakayoshi, H. Lu, S. Lao, SURF tracking, in: International  
489 Conference on Computer Vision, 2009, pp. 1586–1592.
- 490 [26] M. Donoser, H. Bischof, Efficient maximally stable extremal region  
491 (mser) tracking, in: International Conference on Computer Vision and  
492 Pattern Recognition, Vol. 1, 2004, pp. 553–560.
- 493 [27] I. Gutman, B. Zhou, Laplacian energy of a graph, *Linear Algebra and*  
494 *its applications* 414 (2006) 29–37.
- 495 [28] L. Torresani, V. Kolmogorov, C. Rother, Feature correspondence via  
496 graph matching: Models and global optimization, in: European Confer-  
497 ence on Computer Vision, 2008, pp. 596–609.
- 498 [29] H. G. Barrow, J. M. tenenbaum, R. C. Bolles, H. C. Wolf, Parametric  
499 correspondence and Chamfer matching: Two new techniques for image  
500 matching, in: International Joint Conference on Artificial Intelligence,  
501 Vol. 2, 1977, pp. 659–663.
- 502 [30] S. M. S. Nejhum, J. Ho, M.-H. Yang, Visual tracking with histograms  
503 and articulating blocks, in: International Conference on Computer Vi-  
504 sion and Pattern Recognition, 2008.

Chuan Li is a research officer in the Department of Computer Science, University of Bath. He holds a Bachelor Degree in Software Engineering (Zhejiang University, 2005) and a PhD Degree in Computer Science (University of Bath, 2010). He works at the Media Technology Research Center and the Center of Digital Entertainment, both at University of Bath.

Accepted manuscript

Liang Wang received the PhD degree from the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, in 2004. From July 2004 to January 2007, he was with the Department of Computing, Imperial College London, and the Institute for Vision Systems Engineering, Monash University, respectively. Currently, he is with the Department of Computer Science, University of Bath. He has widely published in major international journals such as IEEE Transactions in Pattern Analysis and Machine Intelligence, IEEE Transactions Image Processing, and Pattern Recognition and conferences such as ICDM, ICCV, and CVPR. His main research interests include pattern recognition, computer vision, image and video processing, machine learning, data mining, etc.

Peiyi Shen is a Professor in National School of Software at the XiDian Univ. He is also a member of IEEE. He was a research officer and a PhD student in the MTRC, Computer Science at the Univ. of Bath, and a research fellow in CVSSP at the Univ. of Surrey, studying and working under Prof. Philip Willis and Prof. Adrian Hilton.

He was also with Agilent Technologies in the USA, UK, Malaysia and Singapore from 2000 to 2003. He was also a postdoctoral research fellow in the School of Computing at the National Univ. of Singapore in 2000 after he took his first PhD in XIDIAN Univ. in 1999.

His research interests are in Computer Vision, Volume Visualization and ITS applications.



Peter Hall is Associate Professor in the Department of Computer Science at the University of Bath. He is also director of the Media Technology Research Centre, also at Bath. He founded to vision, video and graphics network of excellence in the UK, and has served on the executive committee of the British Machine Vision Conference since 2003. He has published extensively in Computer Vision, especially where it interfaces with Computer Graphics. More recently he is developing an interest in robotics.

Accepted manuscript

Yi-Zhe Song received both the B.Sc. (first class) and Ph.D. degrees in Computer Science from the Department of Computer Science, University of Bath, UK, in 2003 and 2008 respectively; prior to his doctoral studies, he obtained a Diploma (M.Sc.) degree in Computer Science from the Computer Laboratory, University of Cambridge, UK, in 2004. After his Ph.D., he continued in the same department to become a research and teaching fellow. Since 2011, he became a lecturer (assistant professor) at School of Electronic Engineering and Computer Science, Queen Mary, University of London.

Accepted manuscript











