

000  
001  
002054  
055  
056003 

# Common Fate Hough Transform for Concurrent Detection and Recognition

057  
058  
059004  
005  
006  
007  
008  
009  
010  
011060  
061  
062  
063  
064  
065012 

## Abstract

066  
067

This paper proposes a method for concurrent detection and recognition of multiple multi-class objects by extending the probabilistic Hough transform to incorporate grouping results of the voting elements. In the method, each object part casts votes about the center and the class label of the object that generates itself. These object parts that will vote are firstly grouped by their motion patterns. The votes of each object part are assigned different priors according to the votes of the other object parts in the same motion group. The traceability of the Hough image formed by these votes is greatly enhanced. In such a manner, the proposed method meets the challenges of separating near objects and separating similar objects, while using a codebook trained from images with a clustered background. Experimental results are provided to show the merit in terms of detection and recognition accuracy.

068  
069030  
031070  
071032 

## 1. Introduction

072  
073

The goal of the proposed method is to detect and recognize multiple multi-class objects on a given video frame. This is of importance for many applications, like surveillance and video analysis.

074  
075

One reason, the generalized Hough transform [1] is favorable for detection and/or recognition, is its robustness to partial deformation, besides eigen window methods [15], pictorial structures [4], constellation models [5], sliding window classifiers [10, 22], and branch-and-bound methods [22]. Object parts act as voting elements in the Hough transform based methods, in the form of keypoint descriptors [11], image patches [6, 16], or image regions [7]. Each voting element votes for hypotheses that generate itself. The votes from different voting elements are added up to form a Hough image. The peaks of the Hough image are considered as detection hypotheses with the height of each peak as the confidence of the corresponding hypothesis. The simplicity of the learning procedure is another reason that the Hough transform based methods are attractive. An appearance codebook of object parts is built from a set of images

076  
077

on each of which the object center is annotated. Each code encodes the appearance and the offset to the object center of an object part. This offset acts as a Hough vote at run time, when object parts from the current image are matched against the codebook.

078  
079

Besides the recent usage for various vision tasks [13, 21, 9, 17], the Hough transform based methods are extended for better detection performance. The implicit shape model [11, 12] is extended by notifying correspondences between the object parts and the hypotheses [2] for the detection of multiple near objects. The Hough transform is placed in a discriminative framework for object detection [13] in a way that the codes are assigned different weights by the co-occurrence frequency of their appearance and offset to the object center.

080  
081

The proposed method not only detects multiple objects, but also recognizes each detected object. For recognition, the class label of each training image is annotated besides the object center. Then each code contains three aspects of information: the appearance, the offset to the object center, and the class label. The Hough image is formed by adding up votes with class labels, and the peaks of the Hough image are detection hypotheses with class labels. For each object, its location and class label are found from the hypotheses. However, for realistic scenarios with multiple multi-class objects, the Hough image formed in this manner is difficult to search over for concurrent detection and recognition.

082  
083

One reason is that the object parts of the codebook are from multi-class objects, which even of the same class vary in poses. Thus the portion of the correct votes is low. For the traceability of the Hough image, the codebook needs to be very effective. So, the training images where objects appear as foreground need to have a very clean background. Otherwise, keypoints on the background act as noise and lead to false votes. The noisy votes make the portion of the correct votes even lower. Though computer graphics rendering provides a good source of clean-background images for the training of voting based methods [14], in order to use real-world images with a clustered background, manual efforts are needed to mark the foreground, or a very large number of training images are needed to build the codebook, which

084  
085086  
087088  
089090  
091092  
093094  
095096  
097098  
099100  
101102  
103104  
105106  
107

108 harms efficiency.

109 Even if a very effective codebook is successfully built,  
 110 the Hough image built upon a scene containing near objects  
 111 and similar objects is still difficult to search over. This is  
 112 due to two reasons: (1) votes casted by parts from two near  
 113 objects make the peaks corresponding to different objects  
 114 mixed up, and (2) different-class object parts are sometimes  
 115 very similar, and this leads to tough decisions on the class  
 116 label of the peaks.

117 To utilize the codebook built from training images with  
 118 a clustered background and to meet the challenges of de-  
 119 tecting near objects and recognizing similar objects, this pa-  
 120 per proposes a concurrent detection and recognition method  
 121 based on the common fate principle [20]. The principle is  
 122 one of the four visual perception principles as theorized by  
 123 gestalt psychologists. For humans, tokens moving coher-  
 124 ently are perceptually grouped, and this provides an intuition  
 125 to group the object parts (keypoint descriptors) that  
 126 will vote by their motion patterns. The motion of each key-  
 127 point is represented by a trajectory generated by tracking the  
 128 keypoint through frames. The keypoints are then grouped  
 129 using the pairwise similarities of their corresponding trajec-  
 130 tories. Among the votes of each object part, those which are  
 131 more “agreeable” by the votes of the other object parts in the  
 132 same motion group are assigned higher priors. This results  
 133 in correct votes being more likely to be assigned higher pri-  
 134 ors. These votes with different priors enhance the traceabil-  
 135 ity of the Hough image for concurrent detection and recog-  
 136 nition.

137 To the best of our knowledge, this is the first to incor-  
 138 porate motion information in the Hough transform based  
 139 methods for concurrent detection and recognition. Addi-  
 140 tionally, the proposed method has several appealing prop-  
 141 erties:

- 142 • The method realizes concurrent detection and recogni-  
 143 tion of multiple multi-class objects. The existence of  
 144 three types of objects makes the task challenging: near  
 145 objects, similar different-class objects, and multi-pose  
 146 same-class objects.
- 147 • Its ability to use a codebook trained by images with a  
 148 clustered background.
- 149 • The manner in which the grouping results by motion  
 150 is integrated into the Hough transform is very general,  
 151 and it can be used for incorporating other grouping re-  
 152 sults.

153 The rest of the paper is organized as follows: The prob-  
 154 abilistic Hough transform formulism for a concurrent de-  
 155 tection and recognition problem is given in section 2. In  
 156 Section 3 the formulism of the common fate Hough trans-  
 157 form is given. The inference for concurrent detection and

158 recognition is described in Section 4. Then experimental  
 159 results are given in section 5, and section 6 concludes.

## 160 2. Probabilistic Hough Transform

161 This section describes how a Hough image for concur-  
 162 rent detection and recognition is formed from object parts  
 163 observed on an image.

164 Let  $e$  denote an object part observed on the image. The  
 165 appearance of  $e$  is matched against the codebook, and  $e$  acti-  
 166 vates  $N$  best matched codes from the trained codebook.  
 167 Each code contains the appearance, its offset to the object  
 168 center, and the class label. According to the  $N$  matched  
 169 codes,  $e$  casts  $N$  votes. Each vote  $V_e$  is about the object  
 170 center that generates  $e$ . The position of the object center  
 171 casted by  $V$  is denoted by  $x_V$ , while the class label by  $l_V$ .  
 172 Based on the  $N$  votes of  $e$ , the probability that a position  $\tilde{x}$   
 173 is the center of an object with class label  $\tilde{l}$  is given by,

$$174 p(\tilde{x}, \tilde{l}|e) = \sum_{i=1}^N p(\tilde{x}, \tilde{l}|V_e^i)p(V_e^i|e). \quad (1)$$

175 Here  $p(\tilde{x}, \tilde{l}|V_e^i)$  is the likelihood of  $\tilde{x}$  being an object center  
 176 of class  $\tilde{l}$  based on  $V_e^i$ . And  $p(V_e^i|e)$  is the prior of  $V_e^i$ , given  
 177  $e$ .

178 The idea of the proposed method is that, the prior term,  
 179  $p(V_e^i|e)$ , is defined by the motion grouping results of all the  
 180 object parts.

181 The likelihood term is defined as,

$$182 p(\tilde{x}, \tilde{l}|V) = \begin{cases} 0 & \text{if } l_V \neq \tilde{l} \text{ or } |\tilde{x} - x_V| > d \\ G(\tilde{x}; x_V, \sigma) & \text{otherwise} \end{cases}. \quad (2)$$

183 Here  $G(\tilde{x}; x_V, \sigma)$  is a Gaussian function that fixes the spa-  
 184 cial gap between  $\tilde{x}$  and  $x_V$ .

185 Let  $M$  be the total number of object parts on the image,  
 186 and assume the object parts are mutually independent. Then  
 187 by marginalizing over all the object parts, the probability of  
 188  $\tilde{x}$  being the center of a  $\tilde{l}$ -class object is given by,

$$\begin{aligned} 189 p(\tilde{x}, \tilde{l}) &= \sum_{j=1}^M p(\tilde{x}, \tilde{l}|e_j)p(e_j) \\ 190 &= \sum_{j=1}^M \sum_{i=1}^N p(\tilde{x}, \tilde{l}|V_{e_j}^i)p(V_{e_j}^i|e_j)p(e_j). \end{aligned} \quad (3)$$

191 Usually, a uniform prior is assumed for each object part,  
 192 and  $p(e_j) = \frac{1}{M}$ . Then by considering  $p(\tilde{x}, \tilde{l})$  as the evalua-  
 193 tion score of the Hough space  $(\tilde{x}, \tilde{l})$ , the task of concurrent  
 194 detection and recognition converts to finding and then vali-  
 195 dating the local maxima of the Hough image.



216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
Figure 1. Effect of the proposed prior. The motion grouping results are given in (a), and different motion groups are marked with different colors. The object centers voted according to the 7 best matched codes are given in (b). Red circles are centers of pedestrians while blue circles are centers of bicycle riders. The voted centers given in (c) are the voted centers of the 7 votes with the highest defined priors from 35 votes. The object centers in (d) are voted by votes with priors higher than 0.1.

### 3. Common Fate Hough Transform

232 Since the codebook is trained by multi-class object parts  
233 and the training images have a clustered background, the  
234 best matched codes here cannot guarantee the traceability of  
235 the Hough image. This section describes how, by considering  
236 the motion grouping results of the object parts, different  
237 priors are assigned to the votes of each object part.

238 Let  $\gamma = \{g\}$  denote the grouping results, where  $g$  is a  
239 group of object parts, i.e.,  $e_m \in g$  and  $e_n \in g$ . Those votes  
240 of  $e_m$  which are more “agreeable” by the votes of the other  
241 objects in  $g$  are assigned higher priors.

242 Towards this end, the relationship between the votes of  
243  $e_m$  and the votes of  $e_n$  needs to be given in advance. This  
244 relationship is named support. The support from  $V_{e_n}$  to  $V_{e_m}$   
245 is defined by that based on  $V_{e_n}$ , the possibility  $V_{e_m}$ ’s voted  
246 center is correct, as,

$$247 S(V_{e_n} \rightarrow V_{e_m}) = p(\mathbf{x}_{V_{e_m}}, l_{V_{e_m}} | V_{e_n}), n \neq m.$$

248 Here  $p(\mathbf{x}_{V_{e_m}}, l_{V_{e_m}} | V_{e_n})$  is defined in (2). This measures  
249 the coherence of the two votes from different object parts.

250 Then, the support from  $e_n$  to  $V_{e_m}$  is defined by that  
251 based on  $e_n$ , the possibility  $V_{e_m}$ ’s voted center is correct,  
252 as,

$$253 S(e_n \rightarrow V_{e_m}) = p(\mathbf{x}_{V_{e_m}}, l_{V_{e_m}} | e_n) \\ 254 = \sum_{i=1}^N p(\mathbf{x}_{V_{e_m}}, l_{V_{e_m}} | V_{e_n}^i) p(V_{e_n}^i | e_n) \\ 255 = \sum_{i=1}^N S(V_{e_n}^i \rightarrow V_{e_m}) p(V_{e_n}^i | e_n), n \neq m.$$

256 And the support from  $g$  to  $V_{e_m}$  is defined by the possi-  
257 bility that  $V_{e_m}$ ’s voted center is correct based on the votes  
258 of all the other object parts but its belonging object part in

259  $g$ , as,

$$260 S(g \rightarrow V_{e_m}) = \sum_{e_i \in g - \{e_m\}} p(\mathbf{x}_{V_{e_i}}, l_{V_{e_m}} | e_i) p(e_i) \\ 261 = \frac{1}{M} \sum_{e_i \in g - \{e_m\}} S(e_i \rightarrow V_{e_m}).$$

262 Assuming all object parts in the same motion group are  
263 from the same object, which means motion grouping gives  
264 good results. The center position and the class label given  
265 by one vote shall be consistent with that given by the motion  
266 group. Thus for a particular vote of  $e_m$ , i.e.,  $V_{e_m}$ , a prior is  
267 assigned to it by considering its consistence with  $g$  and the  
268 consistence of  $e_m$ ’s other votes with  $g$ , as:

$$269 p(\tilde{V}_{e_m} | e_m) = \frac{S(g \rightarrow \tilde{V}_{e_m}) + \frac{\Delta}{N}}{\sum_{i=1}^N S(g \rightarrow V_{e_m}^i) + \Delta} \\ 270 = \frac{\sum_{e_j \in g - \{e_m\}} S(e_j \rightarrow \tilde{V}_{e_m}) + \frac{M\Delta}{N}}{\sum_{i=1}^N \sum_{e_j \in g - \{e_m\}} S(e_j \rightarrow V_{e_m}^i) + M\Delta} \\ 271 = \frac{\sum_{e_j \in g - \{e_m\}} \sum_{k=1}^N S(V_{e_j}^k \rightarrow \tilde{V}_{e_m}) p(V_{e_j}^k | e_j) + \frac{M\Delta}{N}}{\sum_{i=1}^N \sum_{e_j \in g - \{e_m\}} \sum_{k=1}^N S(V_{e_j}^k \rightarrow V_{e_m}^i) p(V_{e_j}^k | e_j) + M\Delta}. \quad (4)$$

272 Here,  $\Delta$  is a small constant for preventing zeros. Notice,  
273  $p(\tilde{V}_{e_m} | e_m)$  is defined using  $p(V_{e_j}^k | e_j)$ , the priors of  
274 the votes of the other object parts in  $g$ . In order to give  
275  $p(\tilde{V}_{e_m} | e_m)$ , uniform priors are firstly assigned to the votes  
276 of each object part in  $g$ , i.e.,  $p(V_{e_j}^k | e_j) = \frac{1}{N}$ . Then new  
277 priors are calculated based on the uniformly assigned priors.  
278 The priors of votes to form the Hough image are priors  
279 converged in iterations.

280 The grouping results  $\gamma = \{g\}$ , can be replaced by other  
281 grouping results of the voting elements. The proposed

324 method uses motion to group the voting elements. The  
 325 extended Hough transform with motion grouping results is  
 326 called the common fate Hough transform. The voted centers  
 327 voted by votes according to the best matched codes and  
 328 the centers voted by votes with the highest priors are shown  
 329 in Figure 1. The centers voted by votes with highest priors  
 330 defined here are more concentrated to the true object centers  
 331 than the centers voted according to the best matched codes.  
 332

### 333 3.1. Motion Grouping

334 In order to group the object parts by their motion patterns.  
 335 The object parts are tracked through frames before  
 336 and after the current frame to generate trajectories. The ob-  
 337 ject parts in this method are in the form of keypoint descrip-  
 338 tors. The Harris Corner [8] feature is chosen to represent the  
 339 object part, while for appearance, the region covariance [19]  
 340 of the image patch around the keypoint is used.  
 341

342 For each object part, a trajectory is generated by tracking  
 343 its corresponding Harris Corner by the KLT tracker [18].  
 344

345 To group the trajectories, pairwise similarities are firstly  
 346 defined. Let  $T_{\mathbf{e}_m}$  and  $T_{\mathbf{e}_n}$  denote two trajectories corre-  
 347 sponding to  $\mathbf{e}_m$  and  $\mathbf{e}_n$ . The first similarity between two  
 348 trajectories is defined as,

$$349 D_1(T_{\mathbf{e}_m}, T_{\mathbf{e}_n}) = \max_{i=1 \dots L} (Euclid(\mathbf{x}_{T_{\mathbf{e}_m}}^i, \mathbf{x}_{T_{\mathbf{e}_n}}^i)) .$$

350 Here, *Euclid* calculates the Euclidean distance,  $i$  is the  
 351 frame index, and  $L$  is the length of the overlapping part of  
 352 the two trajectories. To define the second similarity, the  $i$ th  
 353 directional vector of  $T$  is firstly defined as,  $\mathbf{d}_T^i = \mathbf{x}_T^{i+3} - \mathbf{x}_T^i$ .  
 354 Let  $\mathbf{a}_i = \mathbf{d}_{T_{\mathbf{e}_m}}^i$ ,  $\mathbf{b}_i = \mathbf{d}_{T_{\mathbf{e}_n}}^i$ ,  $a_i = \frac{\mathbf{a}_i \cdot \mathbf{b}_i}{\|\mathbf{a}_i\| \|\mathbf{b}_i\|}$ , and  $b_i = \frac{\mathbf{a}_i \cdot \mathbf{b}_i}{\|\mathbf{b}_i\|}$ .  
 355 Then the second similarity is defined as,

$$356 D_2(T_{\mathbf{e}_m}, T_{\mathbf{e}_n}) = \max_{i=1 \dots L-3} (\max(|\mathbf{a}_i - a_i \mathbf{a}_i|, |\mathbf{b}_i - b_i \mathbf{b}_i|)) .$$

357 To group the trajectories, the static points are firstly ex-  
 358 cluded. Inspired by [3], a minimal spanning tree of the tra-  
 359 jectories is built upon  $D_1$ , and split by cutting edges larger  
 360 than a threshold,  $D_{th}^1$ . For each element of the splitting re-  
 361 sults, a minimal spanning tree is built upon  $D_2$  and split by  
 362 cutting the edges larger than a threshold,  $D_{th}^2$ . This hierar-  
 363 chical procedure ensures that trajectories in the same group  
 364 have both small  $D_1$  and  $D_2$ .  
 365

366 Each trajectory corresponds to an object part, and the  
 367 grouping results of the trajectories correspond to grouping  
 368 results of the object parts.  
 369

### 370 3.2. Codebook

371 For training, Harris corners are extracted from the train-  
 372 ing images with the object center and the class label anno-  
 373 tated. In this method, region covariance is chosen to repre-  
 374 sent the appearance, which is defined as,  
 375

$$376 \mathbf{r} = \frac{1}{K-1} \sum_{i=1}^K (\mathbf{z}_i - \mu)(\mathbf{z}_i - \mu)^T .$$

377 Here,  $K$  is the number of pixels in the region, and  $\mathbf{z}_i$  is a  
 378 7-dimensional vector regarding the  $(x, y)$  coordinate of the  
 379 pixel, while  $\mu$  is the mean of  $\mathbf{z}_i$ . And  $\mathbf{z}(x, y)$  contains the  
 380 RGB color of the pixel and intensity gradient of the pixel,  
 381 as:  $r(x, y)$ ,  $g(x, y)$ ,  $b(x, y)$ ,  $|\frac{\partial I(x, y)}{\partial x}|$ ,  $|\frac{\partial I(x, y)}{\partial y}|$ ,  $|\frac{\partial^2 I(x, y)}{\partial x^2}|$ ,  
 382 and  $|\frac{\partial^2 I(x, y)}{\partial y^2}|$ .  
 383

384 The appearance similarity between  $\mathbf{r}_m$  and  $\mathbf{r}_n$  is given  
 385 by,  
 386

$$387 \rho(\mathbf{r}_m, \mathbf{r}_n) = \sqrt{\sum_{i=1}^7 \ln^2 \lambda_i} .$$

388 Here,  $\lambda_i$  is the generalized eigenvalue by solving the gen-  
 389 eralized eigenvalue problem,  $\lambda_i \mathbf{r}_m \mathbf{u}_i = \mathbf{r}_n \mathbf{u}_i$ ,  $\mathbf{u}_i \neq 0$ , with  
 390  $\mathbf{u}_i$  the eigenvector.  
 391

392 A square image patch around each keypoint is used to  
 393 represent the appearance of an object part. Six region co-  
 394 variances are generated for each image patch by using the  
 395 pixels of the top-left, the top-right, the bottom-left, the  
 396 bottom-right, the central, and all of the image patch. Then  
 397 besides the offset and the class label, a code contains six  
 398 region covariances. When an object part is matched against  
 399 the codebook, the similarity between the image patch of the  
 400 object part and a code is defined by the smallest similarity  
 401 of the corresponding region covariances.  
 402

## 4. Detection and Recognition

403 After forming the Hough image, the detection and recog-  
 404 nition hypotheses are validated. Let  $\mathbf{h} = \{H\}$  be the points  
 405 in the Hough space which are evaluated by  $p(\mathbf{x}_H, l_H)$  and  
 406 have  $p(\mathbf{x}_H, l_H) > 0$ . Inspired by [2], the hypotheses are  
 407 validated by an optimizing procedure. Let  $O$  be the number  
 408 of the points in  $\mathbf{h}$ . let  $u_i = 1$  or 0 indicate  $H_i$  being a true  
 409 object center or not. The problem is:  
 410

$$411 \arg \max_{u_i} \prod_{i=1}^O p^{u_i}(H_i) \iff \arg \max_{u_i} \sum_{i=1}^O u_i \ln(p(H_i)) .$$

412 Let  $v_{ij} = 1$  or 0 indicate  $e_j$  belongs to  $H_i$  or not, then  
 413

$$414 p(H_i) = \sum_{j=1}^M p(\mathbf{x}_{H_i}, l_{H_i} | \mathbf{e}_j) p(\mathbf{e}_j) \\ 415 = \frac{1}{M} \sum_{j=1}^M v_{ij} p(\mathbf{x}_{H_i}, l_{H_i} | \mathbf{e}_j) ,$$

432 and by assuming one object part belongs to and only belongs to one hypothesis, the problem is,  
 433  
 434

$$\begin{aligned} & \arg \max_{u_i, v_{ij}} \sum_{i=1}^O u_i \ln \left( \sum_{j=1}^M v_{ij} p(\mathbf{x}_{H_i}, l_{H_i} | \mathbf{e}_j) \right) \\ & \text{s.t. : } u_i = 0 \text{ or } u_i = 1, \forall i \\ & \quad v_{ij} = 0 \text{ or } v_{ij} = 1, \forall i, \forall j \\ & \quad \sum_{i=1}^O v_{ij} = 1, \forall j \\ & \quad \sum_{j=1}^M v_{ij} \leq u_i, \forall i . \end{aligned}$$

448 Following [2], the optimal result for the problem is inferred by greedy maximization. As described in Algorithm 1, the largest local maximum of all the local maxima  
 449 is chosen to be the center of a true object and then the object parts belonging to the chosen object center are excluded  
 450 from the object part set. A new Hough image where new objects  
 451 are found is formed using the remaining object parts. This procedure ends when the object part set is empty or the  
 452 confidence of the chosen object is lower than a threshold.  
 453  
 454  
 455  
 456  
 457

---

**Algorithm 1** Greedy Maximization
 

---

460 Let  $\varepsilon$  be the set of object parts,  $p_{th}$  be the low confidence  
 461 threshold to accept detection responses, and  $\hat{\mathbf{h}}$  be the local  
 462 maxima of  $\mathbf{h}$

```

 463 1: while  $\varepsilon \neq \emptyset$  do
 464 2:   Form  $\mathbf{h}$  with  $\varepsilon$ 
 465 3:   Generate  $\hat{\mathbf{h}}$  and select  $H_i \in \hat{\mathbf{h}}$ ,
 466      $u_i = 0$  and  $\forall H' \in \hat{\mathbf{h}}, p(\mathbf{x}_{H_i}, l_{H_i}) \geq p(\mathbf{x}_{H'}, l_{H'})$ 
 467 4:   if  $p(\mathbf{x}_{H_i}, l_{H_i}) \geq p_{th}$  then
 468 5:      $u_i \leftarrow 1$ 
 469 6:     for  $\mathbf{e}_j \in \varepsilon$  do
 470 7:       if  $\forall H' \in \hat{\mathbf{h}}, p(\mathbf{x}_{H_i}, l_{H_i} | \mathbf{e}_j) \geq p(\mathbf{x}_{H'}, l_{H'} | \mathbf{e}_j)$ 
 471         then
 472 8:            $v_{ij} \leftarrow 1$ 
 473 9:            $\varepsilon \leftarrow \varepsilon - \{\mathbf{e}_j\}$ 
 474 10:        end if
 475 11:      end for
 476 12:    else
 477 13:      for  $\mathbf{e}_j \in \varepsilon$  do
 478 14:         $v_{1j} \leftarrow 1$ 
 479 15:      end for
 480 16:       $\varepsilon \leftarrow \emptyset$ 
 481 17:    end if
 482 18:  end while
 483 19:  return all  $u_i, u_i = 1$ 
  
```

---

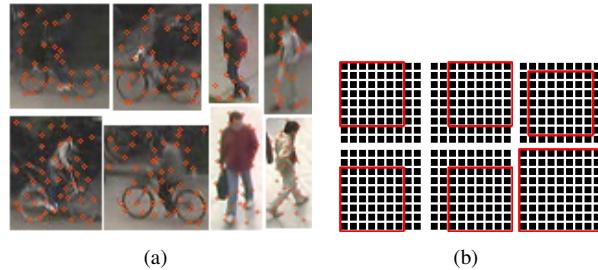


Figure 2. Examples of training images are given in (a). The red circles are the extracted Harris corners. Some keypoints fall on the background. The six rectangles in (b) mark the pixels of the  $9 \times 9$  image patch used for the six region covariances.

## 5. Experimental Results and Evaluation

In the experiments, the enhancement of the Hough image's traceability by the method is verified in terms of detection accuracy and recognition accuracy. The method is firstly tested on the *P-campus* dataset with [2] as benchmark, and then tested on a dataset of similar animals.

### 5.1. Campus Objects Detection and Recognition

**Dataset** The *P-campus* dataset contains two primary classes of foreground objects: pedestrians and bicycle riders. The frame size is  $720 \times 576$ . Among all the 401 continuous frames, 633 different-class ground truth bounding boxes are annotated on 79 frames. In this dataset, pedestrians and bicycle riders have in common the upper human body, and pedestrians appear in front, back, and side views.

**Implementation Settings** For training, 52 images of bicycle riders and 171 images of pedestrians are randomly selected. Harris corners are generated on the image, and training image examples are given in Figure 2(a). For appearance, six region covariances are generated for each keypoint using the  $9 \times 9$  image patch around it as shown in Figure 2(b). The appearance, the offset to the image (object) center, and the label of the training image are encoded, and the code is inserted into a codebook. The final codebook contains 5502 codes.

For motion grouping, each keypoint is tracked through 10 frames before and through 10 frames after the current frame. The similarity of two 21-point trajectories is defined using only the overlapping part. To set the two thresholds for motion grouping,  $D_1$  and  $D_2$  are measured for keypoint pairs of different objects.  $D_{th}^1$  is set that it is larger than only 10% of the measured  $D_1$ s, and so is  $D_{th}^2$ . By doing so, keypoints belonging to different objects are not likely to be grouped together. So, in one motion group, the keypoints are very likely to belong to the same object, as shown in Figure 3.

To form the Hough image, 35 best matched codes are chosen from the codebook for each object part. In (3),  $d$  and  $\sigma$  need to be given. The precision-recall curves are based on



Figure 3. Motion grouping results.

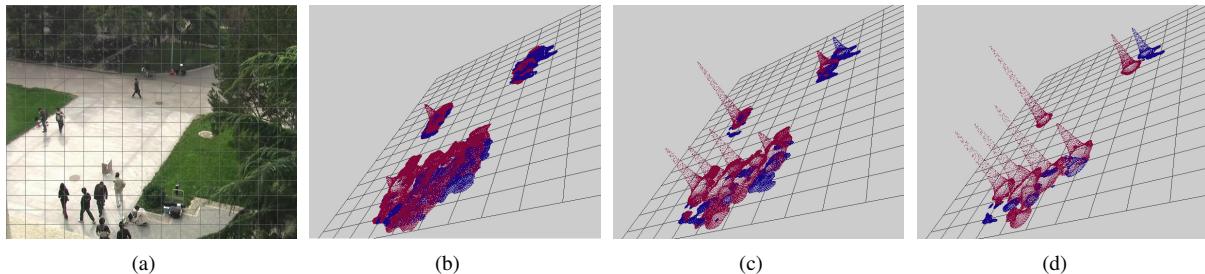


Figure 4. Example Hough images. Grids in (b), (c), and (d) correspond to the grids on the image coordinate. Red indicates pedestrians, while blue indicates bicycle riders. Hough image of (b) is formed by votes with uniform priors, Hough image of (c) is formed by votes with priors after 5 iterations, and Hough image in (d) is formed with converged priors.

$\sigma$ , while  $d$  is set to be 10. Here  $\sigma$  is the most important parameter, while changing  $d$  from 5 to 16 only results in a 1% difference in the detection rate. For Algorithm 1,  $p_{th}$  is set to be 0, which means all local maxima found in Algorithm 1 are accepted.

**Comparisons** For comparison, concurrent detection and recognition are done on the Hough images formed with and without motion grouping results. The same codebook and the same parameter settings are used for forming and searching over both Hough images. The votes of each object part are assigned uniform priors in the benchmark method, while priors defined in (4) are assigned in the proposed method. Example Hough images are shown in Figure 4. With the noisy codebook, peaks of different objects on the original Hough image not only mix up in position but also in class label.

The precision-recall curves are shown in Figure 5(a). An object is considered as correctly detected only if the distance from the ground truth to it is less than 10 pixels. In Figure 5(a), the correctly detected but wrongly recognized objects are considered as true positives, aiming at verifying the detection ability of the proposed method.

The confusion matrices are given in Figure 5(b). For clarity of the comparisons, the proposed method is compared with the benchmark method when they have nearly equal number of false alarms. To evaluate the concurrent detection and recognition ability, a class of “none” to represent missed detections and false alarms is manually added. For example, in Figure 5(b), 487 pedestrian instances are correctly detected and recognized by the proposed method,

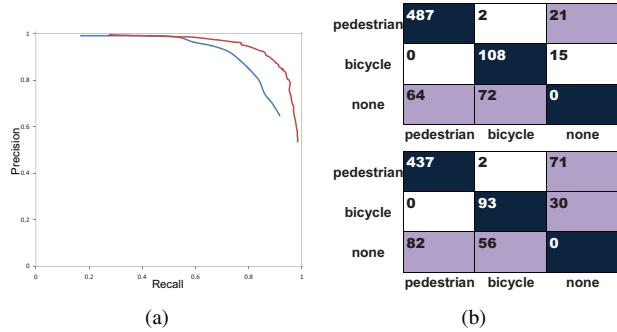


Figure 5. Precision-recall curves (red: the proposed method, blue: the benchmark method) and confusion matrices (upper: the proposed method, down: the benchmark method). In (b), pedestrian and bicycle miss detections are considered as wrongly recognized to be none. Pedestrian or bicycle false alarms are considered as belonging to the none class while wrongly recognized as pedestrians or bicycles.

2 are wrongly recognized to be bicycle riders, and 21 are miss-detected.

The proposed method is limited by its relying on reliable motion information. As shown in Figure 6, failures occur when objects enter the scene.

## 5.2. Big “Cats” Detection and Recognition

**Dataset** To further verify the proposed method, a mini dataset is built upon leopards and tigers of the family Felidae. The dataset contains 6 video clips of 9 leopards and 4 tigers. The frame size is  $640 \times 480$ . Both the animals are in the side view.

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663



Figure 6. Campus objects detection and recognition results. Red rectangles and blue rectangles mark correctly detected and recognized pedestrians and bicycle riders. Yellow rectangles mark missed detections. White rectangles mark correctly detected but not correctly recognized objects. Green rectangles mark false alarms. Black rectangles mark static objects, which are beyond the verification for the method.

667

668



Figure 7. Effect of the proposed prior. Red circles are voted center for leopards, while blue ones are voted centers for tigers. On the top are the motion grouping results. In the middle are the voted centers according to the best matched codes. On the bottom are the voted centers voted by votes with highest priors.

698      **Implementation settings** Most implementation settings  
699      are the same with the settings for campus object detection  
700      and recognition. For training, 5 leopards and 2 tigers are  
701      used. The size of the image patch around each keypoint is

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722

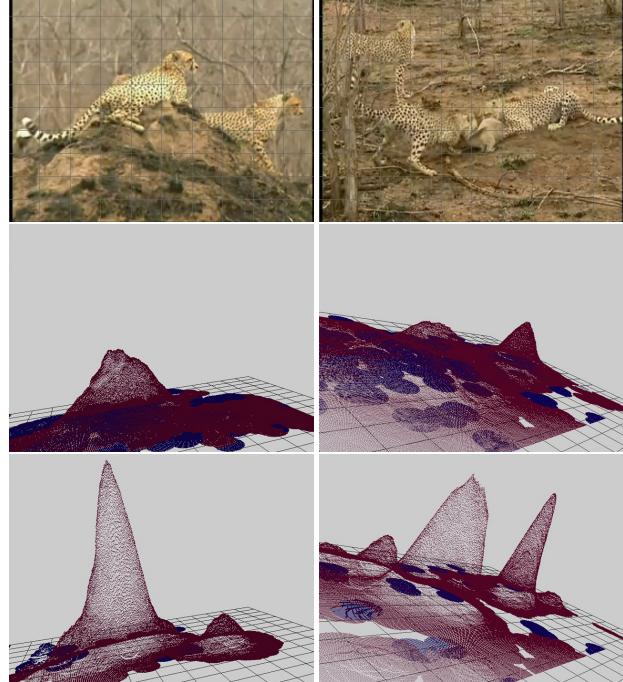


Figure 8. Example Hough images. On the top are the original images. In the middle are the Hough images formed by votes with uniform priors. On the bottom are the Hough images formed by votes with the proposed priors. Red indicates leopards, and blue indicates tigers. Note for the two leopards, there is no peak corresponding to the right one on the benchmark Hough image. For the three leopards, there is also no peak corresponding to the leopard in behind on the benchmark Hough image.

27×27.

**Comparisons** In Figure 7, the motion grouping results

723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744

745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755



Figure 9. Big “cats” detection and recognition results. Red crosses mark the centers for leopards and blue crosses mark the centers for tigers.

and how the voted centers are affected are given. Since parts from different positions of the leopard are very similar, the true center of a leopard is difficult to find from the voted centers of the object parts. In Figure 8, example Hough images are given to show the merit of the proposed prior by the ability to detect leopards. In Figure 9, the detection and recognition results are given. The proposed method successfully detects and recognizes all the leopards and tigers, while the benchmark method miss-detects three leopards.

## 6. Discussion and Conclusion

This paper proposes a method for concurrent detection and recognition by extending the probabilistic Hough transform with motion. The underlying assumption is that object parts moving coherently are likely to belong to the same object. Before voting, all the object parts are grouped by their motion patterns. The votes of each object part are given different priors according to the motion grouping results. In this manner, the proposed method enhances the traceability of the Hough image formed. Experiments show the merit of the method in terms of detection and recognition accuracy.

## References

- [1] D. Ballard. Generalizing the hough transform to detect arbitrary shapes. *Pattern Recognition*, 13(2):111–122, 1981. 1
- [2] O. Barinova, V. Lempitsky, and P. Kohli. On detection of multiple object instances using hough transforms. In *CVPR*, pages 2233–2240, 2010. 1, 4, 5
- [3] G. Brostow and R. Cipolla. Unsupervised bayesian detection of independent motion in crowds. In *CVPR*, pages I: 594–601, 2006. 4
- [4] P. Felzenszwalb and D. Huttenlocher. Pictorial structures for object recognition. *IJCV*, 61(1):55–79, January 2005. 1
- [5] R. Fergus, P. Perona, and A. Zisserman. Object class recognition by unsupervised scale-invariant learning. In *CVPR*, pages II: 264–271, 2003. 1
- [6] J. Gall and V. Lempitsky. Class-specific hough forests for object detection. In *CVPR*, pages 1022–1029, 2009. 1
- [7] C. Gu, J. Lim, P. Arbelaez, and J. Malik. Recognition using regions. In *CVPR*, pages 1030–1037, 2009. 1
- [8] C. Harris and M. Stephens. A combined corner and edge detector. In *Alvey Vision Conference*, pages 147–152, 1988. 4
- [9] J. Knopp, M. Prasad, G. Willems, R. Timofte, and L. Van Gool. Hough transform and 3d surf for robust three dimensional classification. In *ECCV*, pages VI: 589–602, 2010. 1
- [10] C. Lampert, M. Blaschko, and T. Hofmann. Beyond sliding windows: Object localization by efficient subwindow search. In *CVPR*, pages 1–8, 2008. 1
- [11] B. Leibe, A. Leonardis, and B. Schiele. Robust object detection with interleaved categorization and segmentation. *IJCV*, 77(1-3):259–289, May 2008. 1
- [12] B. Leibe and B. Schiele. Interleaved object categorization and segmentation. In *BMVC*, pages 759–768, 2003. 1
- [13] S. Maji and J. Malik. Object detection using a max-margin hough transform. In *CVPR*, pages 1038–1045, 2009. 1
- [14] S. Mohottala, S. Ono, M. Kagesawa, and K. Ikeuchi. Fusion of a camera and a laser range sensor for vehicle recognition. In *OTCBVS*, pages 16–23, 2009. 1
- [15] K. Ohba and K. Ikeuchi. Detectability, uniqueness, and reliability of eigen windows for stable verification of partially occluded objects. *PAMI*, 19(9):1043–1047, September 1997. 1
- [16] R. Okada. Discriminative generalized hough transform for object detection. In *ICCV*, pages 2000–2005, 2009. 1
- [17] M. Sun, G. Bradski, B. Xu, and S. Savarese. Depth-encoded hough voting for joint object detection and shape recovery. In *ECCV*, pages V: 658–671, 2010. 1
- [18] C. Tomasi and T. Kanade. Detection and tracking of point features. Technical report, Carnegie Mellon University, April 1991. 4
- [19] O. Tuzel, F. Porikli, and P. Meer. Region covariance: A fast descriptor for detection and classification. In *ECCV*, pages II: 589–600, 2006. 4
- [20] M. Wertheimer. Laws of organization in perceptual forms (partial translation). In W. B. Ellis, editor, *A Sourcebook of Gestalt Psychology*, pages 71–88. Harcourt, Brace, 1938. 2
- [21] A. Yao, J. Gall, and L. Van Gool. A hough transform-based voting framework for action recognition. In *CVPR*, pages 2061–2068, 2010. 1
- [22] T. Yeh, J. Lee, and T. Darrell. Fast concurrent object localization and recognition. In *CVPR*, pages 280–287, 2009. 1