

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

# Deeply Optimized Hough Transform: Application to Action Segmentation

Adrien CHAN-HON-TONG<sup>1</sup>, Catherine ACHARD<sup>2</sup>, Laurent LUCAT<sup>1</sup>, and  
Patrick SAYD<sup>1</sup>

<sup>1</sup> CEA, LIST, DIASI, Laboratoire Vision et Ingénierie des Contenus, FRANCE  
{adrien.chan-hon-tong, laurent.lucat, patrick.sayd}@cea.fr

<sup>2</sup> Institut des Systèmes Intelligents et Robotique, UPMC, FRANCE  
catherine.achard@upmc.fr

**Abstract.** Hough-like methods (Implicite Shape Model, Hough forest, ...) have been successfully applied in multiple computer vision fields like object detection, tracking, skeleton extraction or human action detection. However, these methods are known to generate false positives. To handle this issue, several works like Max-Margin Hough Transform (MMHT) or Implicit Shape Kernel (ISK) have reported significant performance improvements by adding discriminative parameters to the generative ones introduced by the Implicit Shape Model (ISM). In this paper, we propose to use only discriminative parameters that are globally optimized according to all the variables of the Hough transform. To this end, we abstract the common vote process of all Hough methods into linear equations, leading to a training formulation that can be solved using linear programming solvers.

Our new Hough Transform significantly outperforms the previous ones on HoneyBee and TUM datasets, two public databases of action and behaviour segmentation.

**Keywords:** Hough Transform, Learning, Action Segmentation

## 1 Introduction

The Hough Transform has first been introduced to detect lines in picture. The main idea of this method is to perform the detection not directly in the picture space but in the line parameter space (Hough space) where each line in the image is mapped into a single point. This method has subsequently been extended to parametric objects [1], and non-parametric objects [9] (eg. car, pedestrian, sport activities, ...). For non parametric objects, the Hough Transform first learns a probabilistic-like parametrization of the objects on a training database, and, then performs the detections as a local problem in the corresponding Hough space.

Due to this property of local detection, Hough Transform is a very fast process both in theory (time complexity theory) and practice. For this reason, it has been applied in context of real-time system like [6] for skeleton extraction

## II CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

39 and more generally in multiple computer vision fields like tracking [5], object  
40 detection [4], human action detection [18], image segmentation [10] or human  
41 action segmentation [19].

42 In the context of temporal signals segmentation and recognition, the Hough  
43 Transform is composed of three steps:

- 44 1: Feature extraction and quantization to form codewords
- 45 2: Each codeword votes for each time (in a large neighborhood) and each label  
46 according to a specific learned weight
- 47 3: All the votes are agglomerated to form the Hough score from which segmen-  
48 tation decisions are taken

49 More formally, the Hough Transform (step 2-3) is based on a function  $\theta()$   
50 that links codewords, time displacements (quantified into a finite set) and labels  
51 to vote weights. Thus, a codeword  $w$  extracted at time  $t$  votes with a weight  
52  $\theta(w, l, \Delta_t)$  for the hypothesis that a label  $l$  is present at time  $t + \Delta_t$  (this weight  
53 does not depend on the time  $t$  but only on the relative time displacement  $\Delta_t$ ).  
54 Hence, given a set of localized codewords  $W = \{w, t\}$ , the Hough score  $\mathcal{H}$  for  
55 the label  $l$  at the time  $\bar{t}$  is:

$$\mathcal{H}(\bar{t}, l) = \sum_{(w, t) \in W} \theta(w, l, \bar{t} - t) \quad (1)$$

56 and, the decision about the label (in  $\mathcal{L}$ ) at time  $\bar{t}$  is given by:

$$\hat{l}(\bar{t}) = \max_{l \in \mathcal{L}} (\mathcal{H}(\bar{t}, l)) \quad (2)$$

57 Hence, all the purpose of the training is to select values for  $\theta(w, l, \Delta_t)$  that  
58 will provide correct decisions when following the equations (1) (2) at testing  
59 time. Several works, recalled in section 2, propose to improve the generative  
60 votes used by the Implicit Shape Model (ISM method) by introducing a partial  
61 discriminative optimization process during the vote estimation step. In section  
62 3, we propose to extend these methods by optimizing globally all the votes in a  
63 discriminative way. With this new learning process, our Hough method signifi-  
64 cantly outperforms previous ones on two public datasets of signal segmentation  
65 (the Honeybee dataset [12] and the TUM dataset [15]) as reported in section 4,  
66 before the conclusion in section 5.

## 67 2 State of the Art

68 In this section, we present the different published methods to select the vote  
69 weight during the training step of Hough Transform.

### 70 2.1 Implicit Shape Model

71 In the ISM [9], the Hough Transform (the set of  $\theta()$ -values) is based on gener-  
72 ative weights. Let  $\mathcal{P}(l, \Delta_t | w)$  be the probability that the label at time  $t + \Delta_t$

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE. III

73 is  $l$ , knowing that a codeword  $w$  has been extracted at time  $t$ . This probabil-  
74 ity is estimated with statistics on the training dataset and is supposed to be  
75 independent of  $t$  (it just depends on  $l, \Delta_t$  and  $w$ ). Then, the weights are given  
76 by:

$$\theta_{ISM}(w, l, \Delta_t) = \mathcal{P}(l, \Delta_t | w) \quad (3)$$

77 In practice, the probability  $\mathcal{P}(l, \Delta_t | w)$  is estimated by:

$$\mathcal{P}(l, \Delta_t | w) \approx \frac{N(l, \Delta_t, w)}{N(w)} \quad (4)$$

78 where  $N(l, \Delta_t, w)$  is the number of times a label  $l$  has been seen with a dis-  
79 placement  $\Delta_t$  from a codeword  $w$  and  $N(w)$  is the number of occurrences of the  
80 codeword  $w$ .

81 These ISM-based weights have several advantages (eg. parameter-free train-  
82 ing, robustness to over-training), but they suffer from several drawbacks. In  
83 particular, all codewords and training examples have the same importance and  
84 are considered independently from each other. Two methods, MMHT [11] and  
85 ISK [20] have been introduced to solve these drawbacks.

## 86 2.2 Max-Margin Hough Transform

87 In MMHT [11], a coefficient is introduced for each codeword to weight the ISM  
88 values, resulting in:

$$\theta_{MMHT}(w, l, \Delta_t) = \lambda_w \times \theta_{ISM}(w, l, \Delta_t) = \lambda_w \times \mathcal{P}(l, \Delta_t | w) \quad (5)$$

89 The weights  $\lambda_w$  give more or less importance to the different codewords  $w$  accord-  
90 ing to their discriminative power. They are learnt simultaneously in a discrim-  
91 inative way through an optimisation process similar to support vector machine  
92 (SVM) training [3].

## 93 2.3 Implicite Shape Kernel

94 In ISK [20], the votes are also based on the ISM generative ones, but some  
95 coefficients are introduced to weight the different training examples. Hence, ISK  
96 training leads to:

$$\theta_{ISK}(w, l, \Delta_t) = \sum_i \lambda_i \times \mathcal{P}_i(l, \Delta_t | w) \quad (6)$$

97 where  $\mathcal{P}_i(l, \Delta_t | w)$  is an estimation of the probability  $\mathcal{P}(l, \Delta_t | w)$  based only on  
98 the training example  $i$ . The weights  $\lambda_i$  are learnt simultaneously in a discrimi-  
99 native way using a specific kernel-SVM training [20].

100 MMHT and ISK report experimental improvements over ISM by adding dis-  
101 criminative parameters. This trend is also supported by [17] (we call this method  
102 Scaled Implicit Shape Model SISM).

IV CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

103 **2.4 Scaled Implicit Shape Model**104 This method [17] is also based on ISM but introduces a weighting coefficient for  
105 each displacement, resulting in:

$$\theta_{SISM}(w, l, \Delta_t) = \lambda_{\Delta_t} \times \mathcal{P}(l, \Delta_t | w) = \lambda_{\Delta_t} \times \theta_{ISM}(w, l, \Delta_t) \quad (7)$$

106 As in [11, 20], the weights  $\lambda_{\Delta_t}$  are learnt simultaneously in discriminative way  
107 through a SVM training.108 **2.5 Hough forest**109 To our knowledge ISM [9] and the presented extensions [11, 17, 20] are the only  
110 published methods to estimate the weights of Hough Transform. More precisely,  
111 these methods define links between codewords and votes. There are, of course,  
112 various ways to select the features and the codewords, like, the Hough forest  
113 methods which are major methods of the state of the art. Hough forests use ISM  
114 votes, but the mapping between features (usually data patches) and codewords (a  
115 leaf in a weak binary classifier tree) is constructed such that all training features  
116 associated with a same codeword are expected to come from training examples  
117 with a same label. Several works, like [4], report that this automatic feature  
118 mapping process associated with ISM votes leads to significant experimental  
119 improvements against codewords obtained without learning, by K-means for  
120 example.121 However, in this paper, we focus on the optimisation of the weights used  
122 during the vote process and so to the link between codewords and votes which is  
123 generic whatever the features and codewords used. Thus, the proposed method  
124 can be employed in the Hough forest context by substituting the weights esti-  
125 mated by ISM by the weights optimized by our proposed method.126 The common point between MMHT, ISK and SISM is that they add dis-  
127 criminative parameters to the generative ones introduced by the ISM. In this  
128 paper, we propose to use only discriminative votes strongly optimized. We call  
129 this method Deeply Optimized Hough Transform (DOHT).130 **3 Deeply Optimized Hough Transform**131 The goal of the training process is to establish a correspondence between code-  
132 words and weights. While the ISM methods only use generative weights, MMHT,  
133 ISK and SISM introduce discriminative parameters optimized according to code-  
134 words, training examples or displacements. Using these methods that optimize  
135 only one parameter of  $\theta(w, l, \Delta_t)$ , a small number of coefficients  $\lambda$  have to be  
136 determined. So the optimization process can be solved using SVM. We propose  
137 in this paper to optimize all the weights in a global way, according to all the  
138 parameters of  $\theta(w, l, \Delta_t)$  in multi-class context. In this way, we do not use ISM  
139 values and the method becomes deeply discriminative. The problem is that the  
140 number of unknown parameters is more important and their optimisation using

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE. V

141 SVM becomes intractable. So, we propose to reformulate the problem such that  
142 it becomes linear according to the unknown coefficients.

143 The goal is to define a function  $\theta()$  such that for all training examples (whose  
144 set is denoted  $\mathcal{T}$ ) and all times  $\bar{t}$ , the predicted label  $\hat{l}$  is the real one  $l^*$  (known  
145 on the training data).

146 Considering the definition of the predicted label  $\hat{l}$  (eq. (2)), our problem  
147 formulation is equivalent to :

$$\forall \bar{t}, l \neq l^*(\bar{t}), \mathcal{H}(\bar{t}, l) < \mathcal{H}(\bar{t}, l^*(\bar{t})) \quad (8)$$

148 by dividing  $\theta()$  by the minimal gap, this is equivalent to

$$\forall \bar{t}, l \neq l^*(\bar{t}), \mathcal{H}(\bar{t}, l) + 1 \leq \mathcal{H}(\bar{t}, l^*(\bar{t})) \quad (9)$$

149 and, using equation 1,

$$\forall \bar{t}, l \neq l^*(\bar{t}), \left( \sum_{(w,t) \in W} \theta(w, l, \bar{t} - t) \right) + 1 \leq \left( \sum_{(w,t) \in W} \theta(w, l^*(\bar{t}), \bar{t} - t) \right) \quad (10)$$

150 Hence, the constraints on the function  $\theta()$  are naturally linear. As in [3], to  
151 manage noisy training data, a soft margin framework is applied. For this purpose,  
152 some variables  $\xi$  are introduced leading to:

$$\forall \bar{t}, l \neq l^*(\bar{t}), \sum_{(w,t) \in W} \theta(w, l, \bar{t} - t) + 1 - \xi(\bar{t}) \leq \sum_{(w,t) \in W} \theta(w, l^*(\bar{t}), \bar{t} - t) \quad (11)$$

153 with the objective function:  $\min_{\theta \geq 0, \xi \geq 0} \left( \sum_{\bar{t}} \xi(\bar{t}) \right)$ .

154 To prevent over-fitting, a regularity term is added to the objective function  
155 as in [13]. It penalizes the gap between  $\theta()$  and the uniform votes (0 here). A  
156 coefficient  $\Upsilon$  regulates the trade off between the attachment to data and the  
157 regularity as in [3, 13]. In addition, as  $\theta(w, l, \Delta_t)$  and  $\theta(w, l, \Delta_t + \delta)$  should be  
158 close for a small  $\delta$  and for all  $w$  and  $l$ , we regularly quantify all possible  $\Delta_t$   
159 values.

160 Finally, the problem to solve is formulated as:

$$\begin{aligned} & \min_{\theta \geq 0, \xi \geq 0} \left( \sum_{(w,l,\Delta_t)} \theta(w, l, \Delta_t) + \Upsilon \sum_{\bar{t}} \xi(\bar{t}) \right) \\ & \text{under constraints: } \forall W \in \mathcal{T}, \bar{t}, l \in \mathcal{L} \setminus \{l^*(\bar{t})\}, \\ & \sum_{(w,t) \in W} (\theta(w, l^*(\bar{t}), \bar{t} - t) - \theta(w, l, \bar{t} - t)) + \xi(\bar{t}) \geq 1 \end{aligned} \quad (12)$$

161 As previously stated, a significant difference between MMHT,ISK or SISIM  
162 and DOHT is that we optimize simultaneously all values  $\theta(w, l, \Delta_t)$  of the theta  
163 function, and not only some variables in order to improve the  $\theta_{ISM}$  function.

VI CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

164 Hence, our set of variables is indexed by codewords  $w$  (as in HHMT), displace-  
 165 ments  $\Delta_t$  (as in SISM) and also by labels  $l$  as we consider a multi-class context  
 166 and not only a binary context (MMHT, ISK, SISM). These differences are sum-  
 167 marized in table 1. An other difference between our method and MMHT, SISM  
 168 or SVM is that the penalization of the gap between  $\theta()$  and 0 is measured in  
 169  $L_1$ -norm and not in  $L_2$ -norm. The  $L_1$ -norm allows to obtain linear equations  
 170 and so, to solve the problem efficiency (for example using the solver CPLEX<sup>3</sup>  
 171 available freely for academic purpose) as it is a linear program which is a well  
 172 studied problem in literature (eg. [8]).

methods	$\theta$	variables
ISM [9]	$\theta_{ISM}(w, l, \Delta_t) = \mathcal{P}(l, \Delta_t w)$	-
MMHT [11]	$\theta_{MMHT}(w, l, \Delta_t) = \lambda_w \times \mathcal{P}(l, \Delta_t w)$	$\lambda_w$
ISK [20]	$\theta_{ISK}(w, l, \Delta_t) = \sum_i (\lambda_i \times \mathcal{P}_i(l, \Delta_t w))$	$\lambda_i$
SISM [17]	$\theta_{SISM}(w, l, \Delta_t) = \lambda_{\Delta_t} \times \mathcal{P}(l, \Delta_t w)$	$\lambda_{\Delta_t}$
DOHT (our)	$\theta_{DOHT}(w, l, \Delta_t) = \lambda_{w,l,\Delta_t}$	$\lambda_{w,l,\Delta_t}$

$\mathcal{P}(l, \Delta_t|w)$  is the probability that the label at time  $t + \Delta_t$  is  $l$  knowing that a  
 codeword  $w$  has been extracted at time  $t$ .  $\mathcal{P}_i(l, \Delta_t|w)$  is the same probability  
 estimated using only the training example  $i$ .

**Table 1.** The different learning methods of the Hough Transform

173 In the next section, we evaluate the different methods (ISM, HHMT, SISM,  
 174 DOHT) in action segmentation or behavior segmentation contexts. As ISK is  
 175 only adapted to detection and can not be straightforwardly extended to segmen-  
 176 tation, we can not compare it to the others methods.

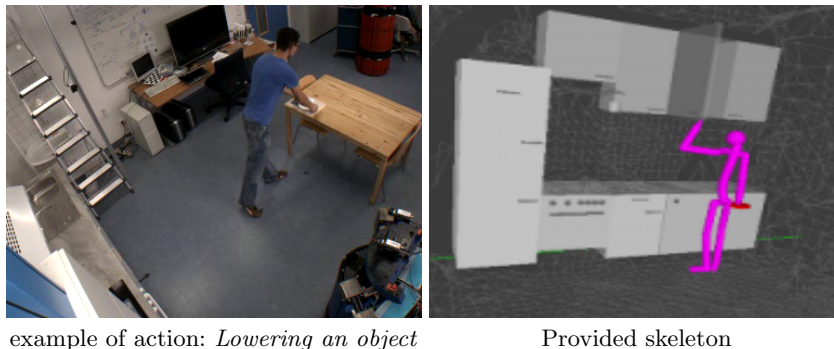
177 **4 Experimental Results**

178 Experiments have been conducted on the TUM [15] and Honeybee [12] datasets.  
 179 These datasets are well designed for segmentation as each frame (here, frames  
 180 and times are equivalent) is associated with a label.

181 **4.1 Application to Human Action Segmentation**

182 TUM is a multi-sensor dataset and in particular it contains skeleton streams  
 183 (fig. 1). It is composed of 19 sequences around 2 minutes each containing 9  
 184 kinds of actions (each action is a label) like *Lowering an object*, *Opening a*  
 185 *drawer* performed by 5 peoples. To provide results comparable to [19], the same  
 186 experimental protocol is applied for splitting data between training and testing  
 187 set, and, results are given in terms of accuracy (number of correctly labelled  
 188 frames divided by the total number of frames).

<sup>3</sup> [www-01.ibm.com/software/websphere/products/optimization/academic-initiative/](http://www-01.ibm.com/software/websphere/products/optimization/academic-initiative/)

example of action: *Lowering an object*

Provided skeleton

**Fig. 1.** TUM dataset [15]

189 As [19] reports better performances using skeleton features (than visual or  
 190 visual plus skeleton ones), we decide to consider only skeleton based features.  
 191 Hence, the input signal of our algorithm is the 3D positions of each articulation  
 192 at each time.

193 We use the same preprocessing (features and codewords) than the bag-of-  
 194 gestures from [2] which achieves the best published performance on this dataset  
 195 (with a manual segmentation). First, the positions are normalized (positions are  
 196 expressed in a system of coordinate linked to the subject to be invariant to camera  
 197 viewpoint, global body position, rotation and size). Then, we consider short  
 198 temporal series of 3D positions of each articulation as features: let the vector  
 199  $(p_1, \dots, p_T)$  be the normalized trajectory of one articulation, then, we consider  
 200 the vector  $(p_{t-\tau}, \dots, p_{t+\tau})$  as a feature extracted at time  $t$ . Similar features are  
 201 also considered in [14, 19, 16] which report the efficiency of interest points tra-  
 202 jectories for human action recognition. Finally, all these features are clustered  
 203 by K-means. The cluster centers defines the codebook and features are mapped  
 204 to their nearest codeword.

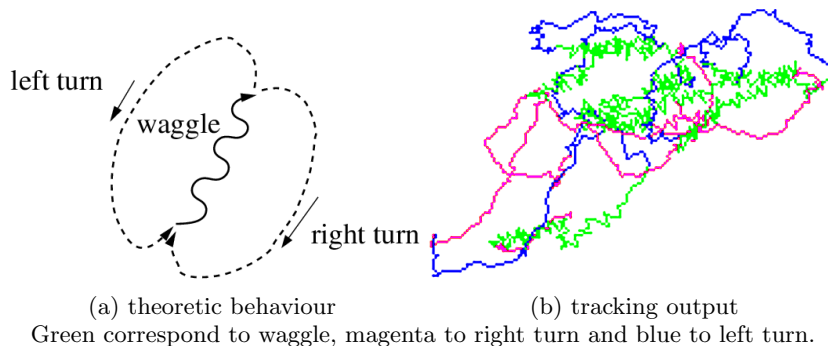
205 More precisely, we consider the 8 main articulations: feet, hands, knees, el-  
 206 bows with  $\tau = 6$ . The quantization with K-means is performed independently for  
 207 each articulation with  $K = 10$ , resulting in 80 codewords. The few parameters  
 208 of this experiments ( $\tau$ ,  $K$ ,  $\mathcal{Y}$  and the quantification granularity of  $\Delta_t$  for the  
 209 optimization process (see section 3)) empirically provide the best performances.  
 210 Results of this experiment are presented in table 2.

211 In this experiment, DOHT significantly outperforms ISM, MMHT and SISM  
 212 and achieves equivalent performance than a SVM based on the same features and  
 213 codeword applied on the optimal segmentation (obtained from the ground truth)  
 214 from [2]. Hence, for this dataset, we achieve equivalent performance than the best  
 215 published (82.6% against 84.3%) without using the optimal segmentation.

VIII CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

216 **4.2 Application to Behaviour Segmentation**

217 Experiments have also been conducted on the Honeybee dataset [12]. The Honey-  
 218 bee dataset provides tracking output of honey bees having 3 kinds of behaviours  
 219 (each behaviour corresponds to a label) correlated with their trajectories (figure  
 220 2). It composed of 6 large sequences. To provide results comparable to [12], the  
 221 same leave-one-out cross validation is applied. A global measure is obtained by  
 222 averaging accuracy from all runs.

**Fig. 2.** Honeybee dataset [12]

223 The input signals in this dataset are the sequences of bee 2D positions and  
 224 orientations  $(x_t, y_t, \alpha_t)$ . As in the previous experiment, normalized short tempo-  
 225 ral series of (2D here) positions are considered as features. Let us call  $R(\beta)$  the  
 226 matrix of the 2D rotation of angle  $-\beta$  and  $p(t) = (x_t, y_t)$ , then we consider the  
 227 vector  $(R(\alpha_t)(p_{t-\tau} - p_t), \dots, R(\alpha_t)(p_{t+\tau} - p_t))$  as the feature extracted at time  
 228  $t$ . All these features are clustered using K-means. The cluster centers defines the  
 229 codebook and features are mapped to their nearest codeword.

230 More precisely, short series of size  $\tau = 0, 3, 6$  are considered in this experi-  
 231 ment. K-means is performed independently for each  $\tau$  with  $K = 16, 32, 64$  re-  
 232 spectively, resulting in 112 codewords. The few parameters of this experiments  
 233 empirically provide the best performances. Results of this experiment are pre-  
 234 sented in table 2.

235 In this experiment, DOHT significantly outperforms ISM, MMHT and SISM.

236 In addition, DOHT achieves equivalent performances than the best published  
 237 results [7]. In [7], a multi-class SVM is applied on each temporal windows (with  
 238 similar kind of features and codewords). Then, segmentation is computed using  
 239 dynamic programming. As scores are computed on each temporal windows, this  
 240 method is **quadratic** in the maximal length of an activity while our is **linear**.  
 241 This quadratic property is a common drawback caused by performing scoring  
 242 as a global problem. Hence, for this dataset, we achieve equivalent performances



243 (86.5% against 89.3%) than the best published results while being significantly  
244 faster.

Method	Accuracy on TUM	Accuracy mean on Honeybee
ISM [9]	58.4	71.9
MMHT [11]	69.6	78.8
SISM [17]	68.5	77.5
DOHT (our)	<b>82.6</b>	<b>86.5</b>

**Table 2.** Global results on TUM [15] and Honeybee [12]

## 245 5 Conclusion

246 In this paper, we propose to use Hough transform to segment and recognize  
247 temporal series. In a non parametric context, the training of Hough transform  
248 consists to properly select the weights used in the voting process. The simple way  
249 (Implicit Shape Model) consists in computing some probabilities on the training  
250 database, leading to a generative model. Some methods (Max-Margin Hough  
251 Transform, Implicit Shape Kernel) propose to add some parameters optimized  
252 on a training database in a discriminative way. In this article, we propose to  
253 skip the first step based on a generative model and to globally learn all the  
254 parameters of the Hough transform on the training database, resulting a deeply  
255 discriminative model. This required to reformulate the voting process to express  
256 it in a linear form in order to use linear programming solvers.

257 We performed several experiments on public datasets where the Hough Trans-  
258 form trained with our method significantly outperforms other Hough Transform  
259 methods and provides equivalent results than best published results for these  
260 datasets while being significantly faster than the corresponding algorithms.

261 In future works, we will evaluate our method on other contexts eg. object  
262 segmentation in image, video spatio-temporal segmentation, automatic speech  
263 segmentation, sign language segmentation.

## 264 References

- 265 1. D.H. Ballard. Generalizing the hough transform to detect arbitrary shapes. *Pattern*  
266 *recognition*, 1981.
- 267 2. Adrien Chan-Hon-Tong, Nicolas Ballas, Catherine Achard, Bertrand Delezoide,  
268 Laurent Lucat, Patrick Sayd, and Françoise Prêteux. Skeleton point trajectories  
269 for human daily activity recognition. In *Proceedings of International Conference*  
270 *on Computer Vision Theory and Application*, 2013.
- 271 3. C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 1995.

X CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

- 272 4. J. Gall and V. Lempitsky. Class-specific hough forests for object detection. In  
273 *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference*  
274 *on*, pages 1022–1029. IEEE, 2009.
- 275 5. J. Gall, A. Yao, N. Razavi, L. Van Gool, and V. Lempitsky. Hough forests for  
276 object detection, tracking, and action recognition. *IEEE Transactions on Pattern*  
277 *Analysis and Machine Intelligence*, 2011.
- 278 6. R. Girshick, J. Shotton, P. Kohli, A. Criminisi, and A. Fitzgibbon. Efficient regres-  
279 sion of general-activity human poses from depth images. In *IEEE International*  
280 *Conference on Computer Vision*. IEEE, 2011.
- 281 7. M. Hoai, Z.Z. Lan, and F. De la Torre. Joint segmentation and classification of  
282 human actions in video. In *IEEE Conference on Computer Vision and Pattern*  
283 *Recognition*. IEEE, 2011.
- 284 8. N. Karmarkar. A new polynomial-time algorithm for linear programming. In *Pro-*  
285 *ceedings of the sixteenth annual ACM symposium on Theory of computing*. ACM,  
286 1984.
- 287 9. B. Leibe, A. Leonardis, and B. Schiele. Combined object categorization and seg-  
288 mentation with an implicit shape model. In *Workshop on Statistical Learning in*  
289 *Computer Vision*, 2004.
- 290 10. Sreedevi M and Jenopaul P. An efficient image segmentation using hough trans-  
291 formation. In *Asian Journal of Information Technology*, 2011.
- 292 11. S. Maji and J. Malik. Object detection using a max-margin hough transform. In  
293 *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009.
- 294 12. S.M. Oh, J.M. Rehg, T. Balch, and F. Dellaert. Learning and inferring motion pat-  
295 terns using parametric segmental switching linear dynamic systems. *International*  
296 *Journal of Computer Vision*, 2008.
- 297 13. J. Shawe-Taylor, P.L. Bartlett, R.C. Williamson, and M. Anthony. Structural risk  
298 minimization over data-dependent hierarchies. *IEEE Transactions on Information*  
299 *Theory*, 1998.
- 300 14. J. Sun, X. Wu, S. Yan, L.F. Cheong, T.S. Chua, and J. Li. Hierarchical spatio-  
301 temporal context modeling for action recognition. In *IEEE Conference on Com-*  
302 *puter Vision and Pattern Recognition*. IEEE, 2009.
- 303 15. M. Tenorth, J. Bando, and M. Beetz. The tum kitchen data set of everyday  
304 manipulation activities for motion tracking and action recognition. In *IEEE 12th*  
305 *International Conference on Computer Vision Workshops*. IEEE, 2009.
- 306 16. Heng Wang, Alexander Kläser, Cordelia Schmid, and Liu Cheng-Lin. Action  
307 Recognition by Dense Trajectories. In *IEEE Conference on Computer Vision &*  
308 *Pattern Recognition*, Colorado Springs, United States, 2011.
- 309 17. Paul Wohlhart, Samuel Schulter, Martin Kostinger, Peter Roth, and Horst Bischof.  
310 Discriminative hough forests for object detection. In *BMVC*, 2012.
- 311 18. A. Yao, J. Gall, and L. Van Gool. A hough transform-based voting framework  
312 for action recognition. In *IEEE Conference on Computer Vision and Pattern*  
313 *Recognition*, 2010.
- 314 19. Angela Yao, Juergen Gall, Gabriele Fanelli, and Luc Van Gool. Does human action  
315 recognition benefit from pose estimation? In *Proceedings of the British Machine*  
316 *Vision Conference*. BMVA Press, 2011.
- 317 20. Yimeng Zhang and Tsuhan Chen. Implicit shape kernel for discriminative learning  
318 of the hough transform detector. In *BMVC*, 2010.