



**HAL**  
open science

## Discriminative part model for visual recognition

Ronan Sicre, Frédéric Jurie

► **To cite this version:**

Ronan Sicre, Frédéric Jurie. Discriminative part model for visual recognition. [Research Report] GREYC CNRS UMR 6072, Université de Caen. 2015. hal-01132389v1

**HAL Id: hal-01132389**

**<https://inria.hal.science/hal-01132389v1>**

Submitted on 17 Mar 2015 (v1), last revised 25 Aug 2015 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Discriminative part model for visual recognition

Ronan Sicre and Frédéric Jurie

CNRS UMR 6072 – University of Caen Basse-Normandie – ENSICAEN – France

Email: {ronan.sicre, frederic.jurie}@unicaen.fr

---

## Abstract

The recent literature on visual recognition and image classification has been mainly focused on Deep Convolutional Neural Networks [1] and their variants, which has resulted in a significant progression of the performance of these algorithms. Nevertheless, these recent advances should not conceal the fact that part-based models are expected to outperform approaches that code images as a whole, because of the flexibility such models offer. Based on this hypothesis, this article introduces a new algorithm for image recognition allowing to model image categories as a collection of distinctive parts, discovered automatically. These parts are matched across images while learning their visual model and are finally pooled to provide images signatures. The so-obtained parts are free of any appearance constraints and are optimized to allow the distinction between the categories to be recognized, in an optimal way. A key ingredient of the approach is a *softassign*-like matching algorithm that simultaneously learns the model of each part and automatically assigns image regions to the model's parts. Once the model of the category is trained, it can be used to classify new images by finding image's regions similar to the learned parts and encoding them in a single compact signature. The approach is experimentally validated by showing that using neural code as low-level image features allows to go beyond the performance given by Deep Convolutional Neural Networks, hence providing state-of-the-art results on several publicly available datasets.

*Keywords:*

## 9 **1. Introduction**

10 The arrival of effective approaches based on Deep Convolutional Neural  
11 Networks (Deep CNN), such as the remarkable work of Krizhevsky *et al.*  
12 [1] has been perceived as a new trend in image classification, relegating the  
13 not so distant approaches such as the bag-of-words [2, 3, 4] or the even more  
14 recent Fisher vectors [5] to what some consider now to be a legacy of previous  
15 time.

16 Since then, the literature on *image classification* – the task consists in  
17 predicting whether an image contains an object or, more generally, a visual  
18 concept based on the content of the image – has benefited from a revival of  
19 interest because of the new perspective Deep CNN provides (*e.g.* [6, 7, 8, 7, 9],  
20 to cite only a few recent of them).

21 However, even if Deep CNN obtain very good performances, most of the  
22 recent approaches can be considered as *holistic* in the sense that the CNN  
23 architecture is fed with the whole image, resulting in a single image vector.  
24 One can believe this is a limitation, as scenes (and therefore images) can be  
25 seen as spatial arrangements of objects or parts, and as the decomposition of  
26 images into distinctive parts can results in more expressive and discriminative  
27 models [10, 11, 12].

28 One motivation of this paper is hence to bring together the advantages of  
29 Deep CNN and part-based model. The results achieved by Oquab[13] *et al.*  
30 constitute one interesting step toward that end. They indeed shown that it  
31 is possible to transfer image representations learned with CNNs trained on  
32 large datasets to different tasks, even in presence of limited training data.  
33 Their method uses ImageNet pre-trained layers of CNN to compute mid-level  
34 image signature and can be utilized as an efficient feature encoding system.  
35 We use this framework as an alternative to Bag-of-words (BOW) or Fisher  
36 vector to encode image regions.

37 Another key issue raised by the representation of images in the context  
38 of image classification, is how to efficiently use geometric information and,  
39 as aforementioned, how to decompose images into stable and distinctive re-  
40 gions. While the early works were building on pure bag-of-words *e.g.* [2],  
41 which consists of pooling the visual features without using their spatial co-  
42 ordinates in any way, it has been shown later (*e.g.* by [4]) that performance  
43 can be significantly improved by encoding separately a set of multiple (pos-  
44 sibly overlapping) regions, which constitutes a first step toward the use of  
45 geometry. Using fixed regions (usually image quad-trees) is obviously limited  
46 as the corresponding implicit segmentations of the image is not adapted to  
47 the image’s content. Several more recent works such as [3, 14, 15] have intro-  
48 duced more flexibility by adapting the shape/position of the regions, but a  
49 strong limitation of these works is that the layout of images is still supposed  
50 to be fixed, for a given category.

51 The proposed work starts with the observation that images within a given  
52 category can have very different layouts or spatial organization, even if they  
53 can be interpreted globally as sharing the same meaning. In line with this  
54 observation, several recent works have shown that categories can be efficiently  
55 represented by a set of distinctive regions either called *parts* or *fragments*  
56 [16, 10, 11, 12]. For example, if ‘car’ images can be recognized because of the  
57 joint presence of ‘wheel’, ‘road’ or ‘window’-like parts, the position of these  
58 regions can be any as long as they are in the image. This idea of introducing  
59 some invariance (or alignment) with respect to the position of the parts have  
60 been successfully utilized in the Deformable Part Model of [17]. However, in  
61 the case of image classification the relative position of the parts is much less  
62 constrained than in the case of object detection.

63 In reaction to these observations and concerns, another motivation of our  
64 work is precisely to propose a new way to describe images by a set of parts  
65 that are aligned across images by construction, without having to use strong  
66 geometric constraints between them. This is achieved by proposing a new  
67 model for categories, which is based on the fact that (i) a category is defined

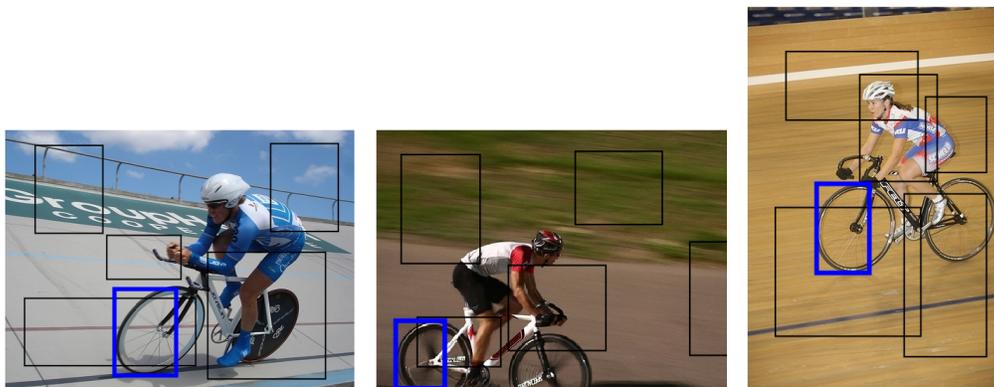


Figure 1: Our system aims at discovering distinctive parts (blue boxes) from a set of regions (black boxes) randomly extracted from images of a category.

68 by a set of  $K$  parts (ii) these parts are distinctive in the sense that they  
 69 occur more frequently in the image of the category than in those from other  
 70 categories (iii) the presence of regions visually similar to the model's parts  
 71 is expected in the images of the category. These definitions are implemented  
 72 into an objective function which is optimized during a learning stage. The  
 73 objective function relies on a *match* function which automatically discovers  
 74 and relates model's parts to image regions. Training can be achieved from  
 75 a set of images describing the category to be recognized, without having to  
 76 provide any extra annotations. In particular, bounding boxes revealing ob-  
 77 jects locations are not necessary. During training, a part classifier is learned  
 78 in conjunction with the alignment of parts to image regions. In a second  
 79 time, these classifiers can be used to build a global visual descriptor of im-  
 80 ages, which combines the signatures of the regions discovered in the image.  
 81 More precisely, the paper proposes 3 representations: one is obtained by ag-  
 82 gregating the Deep CNN signatures of the different image regions, another  
 83 consists in aggregating the scores of individual part classifiers while the third  
 84 encodes the distinctive regions of an image with a Fisher vector.

85 The proposed approach is experimentally validated on three classification  
 86 datasets. First, Willow [18] aims at classifying 7 human actions in still im-

87 ages, while the goal of Boats Datasets is to classify 5 different categories of  
88 boats. Finally the MIT 67 dataset [19] contains images of 67 types of scenes  
89 which are to be recognized. These experiments show that not only the pro-  
90 posed method outperform Deep CNN but also that it offers state-of-the-art  
91 results on the very competitive MIT 67 dataset.

92 The rest of the paper is organized as follows. Related work is presented in  
93 Section 2, while Section 3 provides details on the proposed system that learns,  
94 aligns, and encodes distinctive parts. Finally, the experimental validation is  
95 given in Section 4, before concluding the paper.

## 96 2. Related Work

97 *Image classification.* has received a large attention from the computer vision  
98 community, *e.g.* see the abundant literature related to the Pascal VOC [20]  
99 and ImageNet [21] challenges. A large part of the modern approaches follow  
100 the bag-of-word model [2], composed of a 4 step pipeline: 1) extraction of  
101 local image features, 2) encoding of local image descriptors, 3) pooling of  
102 encoded descriptors into a global image representation, 4) training and clas-  
103 sification of pooled image descriptors for the purpose of object recognition.  
104 Several studies evaluated the influence of the first step: the low level features  
105 *e.g.* gradient, shape, color, and texture descriptors, such as [22], while other  
106 proposed combining different levels (low - mid - high) of information [23].  
107 Regarding the second step: image encoding; Fisher vectors [5] were consid-  
108 ered as achieving state-of-the-art performance, in many cases. The third,  
109 pooling, step is also shown to provide improvements, and spatial and feature  
110 space pooling techniques have been widely investigated [24, 4]. Moreover,  
111 [3, 14] have recently proposed two different strategies for embedding spa-  
112 tial information into the bag-of-words framework. Finally, regarding the last  
113 step of the pipeline, discriminative classifiers such as Support Vector Ma-  
114 chines (SVM) are widely accepted as the reference in terms of classification  
115 performance.

116 During the last months, the deep CNN approaches have been successfully

117 applied to large-scale image classification datasets, such as ImageNet [21] [1],  
118 obtaining state-of-the-art results significantly above Fisher vectors or bag-  
119 of-words approaches. These networks have a much deeper structure than  
120 standard representations, including several convolutional layers followed by  
121 fully connected layers, resulting in a very large number of parameters that  
122 have to be learned from training data. By learning these networks parameters  
123 on large image datasets, a structured representation can be extracted at an  
124 intermediate to a high-level, depending on the extracted layers [25, 26]. Deep  
125 CNN representation have been recently combined with VLAD descriptors [27]  
126 or Fisher vectors [9]

127 *Mid-level features.* Several authors have shown the importance of adding  
128 intermediate representations [28], also referred as the mid-level features, for  
129 leveraging the performance. We observe three mains trends of mid-level  
130 description in the recent literature: hand-crafted, learned, and unsupervised  
131 features. *Hand-crafted mid-level features* aim at encapsulating information  
132 on groups of pixels such as superpixels [29, 30], patches [31] or segments  
133 [32]. These descriptors are computed similarly for any given image and do  
134 not require any learning. On the other hand, a large variety of *learned*  
135 *mid-level features* have been proposed. One of the original method was the  
136 Deformable Part Model, proposed by [17]. We can also mention the semantic  
137 attributes [33, 34] which have received a lot of interest. Within the learned  
138 mid-level features techniques, we observe a large variety in terms of learning  
139 data utilized. While some feature are based on extra training data such as  
140 labeled fragments [35], sketch tokens [36] or pre-trained object detectors [37],  
141 most methods use a standard split of training and testing data to learn the  
142 distinctive features, as the *structural element patch model* [38] or the *blocks*  
143 *that shout* [11]. Finally, regarding *unsupervised mid-level features*, the work  
144 of [39] aims at detecting distinctive patches in an image dataset without any  
145 label information.

146 *Learned mid-level features.* Our work aims at learning distinctive parts with-  
147 out extra annotations. Therefore, closely related work includes the De-

148 formable Part Model (DPM) [17]. The DPM models categories by using  
149 a mixture of parts and classify image regions as object vs non object regions.  
150 Classifiers are applied to a representation in which the parts are aligned, by  
151 shifting the parts with respect to the root filter. However, for image classifi-  
152 cation, the variability of parts positions as well as the variation of appearance  
153 within a category makes the problem different. Our work also bears simi-  
154 larities with [16], which tries to discover the *fragments* that maximize the  
155 mutual information between the category and the presence of the fragment  
156 in the image. However, [16] suffers from that (i) contrarily to [17], part are  
157 just image patches and not discriminative classifiers (ii) the decision is made  
158 by verifying the presence of the fragments in the image, instead of training  
159 a classifier taking fragment descriptors as input. Our approach takes the  
160 advantages of both approaches without having their drawbacks.

161 More recently, [10] proposes a learning framework for the automatic dis-  
162 covery of image’s parts, assuming that partial correspondence between in-  
163 stances of a category are available. These partial correspondences allow  
164 the training of part detectors, used in a first time to extract candidates re-  
165 gions. While we share the same motivations, our approach does not require  
166 any supervision. In addition, it is worth mentioning [11] and [12] which  
167 both propose algorithms for learning parts that are good representatives of  
168 a given category. Our work follows the same objectives, without the local-  
169 ization constraints imposed by [12] and the large computation requirement  
170 and unoptimized encoding of [11]. This work finally shows the importance of  
171 mid-level information and justifies its use to improve recognition capabilities.

172 This article is an extension of [40] providing a richer description of the  
173 related works, more details and improvement of the method as well as a much  
174 more experimental validation.

### 175 **3. Proposed method**

176 From a general point of view, the overall approach consists in three steps:  
177 (i) a learning step during which some category’s distinctive parts are discov-

178 ered, (ii) a representation step in which a global signature of the image is  
179 computed, on the basis of parts presence in the image, (iii) a classification  
180 step relying on a linear SVM classifier. The originality of the work is in  
181 the discovery of category’s distinctive parts and their use in the encoding of  
182 images (two first steps), which are the subject of this section, and not in the  
183 classification step which is the most classic.

184 The model we propose for representing image categories consists in a  
185 collection of  $K$  distinctive parts defined by their visual appearance, without  
186 any geometric relationships between them. It is expected that positive images  
187 (with respect to a given category) contain regions visually similar to these  
188 parts (considered as instances of the parts) while there are fewer of them in  
189 negative images. The distance from an image to the class is then defined as  
190 a function of the set of distances between image regions and parts (*e.g.* using  
191 max pooling).

192 As aforementioned, the main contribution of this paper lies in the method  
193 allowing to automatically discover distinctive parts in the images of a given  
194 category and thus to learn the model of this category. These parts are further  
195 aligned with images regions, which are utilized to produce images signatures.  
196 Signatures are subsequently used in a standard classification framework.

197 This section first presents our part-based model and its associated cost  
198 function, which is to be optimized during learning. In a second time, we  
199 explain how the parameters of the model can be learned using an iterative  
200 framework inspired from the *softassign* algorithm. Then, more details are  
201 given on the algorithm initialization step. Finally, we explain how images  
202 signatures can be computed using the learned model.

### 203 *3.1. Part model and objective function*

204 First, let us introduce some notations. We assume having a set of images  
205 belonging to the category to be modeled, considered as positive training  
206 images and denoted as  $\mathcal{I}^+$ .  $|\mathcal{I}^+|$  represents the number of positive images. In  
207 the same way,  $\mathcal{I}^-$  is the set of (negative) images belonging to other categories.  
208 The whole training set is denoted as  $\mathcal{I} = \mathcal{I}^+ \cup \mathcal{I}^-$  and contains  $|\mathcal{I}|$  images.

209 From each image  $I \in \mathcal{I}$ , we extract a dense random set of image regions  
 210 denoted as  $\mathcal{R}_I$ . Each region  $r$  is represented by its signatures  $x_r$ , which is,  
 211 in practice, the bag-of-words or CNN, representation of the region. More  
 212 details on the description choices are discussed in section 3.4. The model  
 213 of the category includes a set of parts denoted as  $\mathcal{P}$ . The number of parts,  
 214  $K = |\mathcal{P}|$ , is fixed. In the following,  $p \in \mathcal{P}$  denotes one of these parts.

215 As explained before, our model relies on three assumptions: first, the  
 216 model is supposed to be composed of a set of  $K$  different parts. Second, it  
 217 is expected that each part of the model is present in each positive image.  
 218 Third, parts should be representatives of the category, which means that  
 219 they should occur more frequently in positive images than in negative ones.

220 We implement the second constraint by introducing the match function  
 221  $m(r, p)$  associating model parts and image regions, and by imposing that  
 222  $\forall I \in \mathcal{I}^+ r \in \mathcal{I}$  and  $\forall p \in \mathcal{P}$ ,  $\sum_{r \in I} m(r, p) = 1$ . The match function is defined  
 223 as:

$$m(r, p) = \begin{cases} 1 & \text{if region } r \text{ is assigned to part model } p \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

224 In practice, the match function can be seen as a binary matrix with one  
 225 row per part and one column per image region. We add the first constraint  
 226 ensuring that an image region can be assigned to at most one part, which is  
 227 written as:  $\forall I \in \mathcal{I}^+, \forall r \in I, \sum_{p \in \mathcal{P}} m(r, p) \leq 1$ .

228 Regarding the third assumption, which states that regions should be dis-  
 229 criminative, one way to achieve this would be to measure to which extent  
 230 each part can be matched with regions from the negative set, and promote  
 231 those occurring more on positive images. However, such process would be  
 232 very costly. Therefore, as suggested by [11], we use the LDA technique of  
 233 [41], which consists in learning once and for all a universal model of negative  
 234 patches. In practice, the parameter vector  $w$  of a part classifier, correspond-

235 ing to the part  $p$ , is defined simply as:

$$w(p, m) = \Sigma^{-1} \left( \frac{\sum_{r \in I, \forall I \in \mathcal{I}^+} m(r, p) \times x_r}{\sum_{r \in I, \forall I \in \mathcal{I}^+} m(r, p)} - \frac{\sum_{r \in I, \forall I \in \mathcal{I}} x_r}{|r \in I, \forall I \in \mathcal{I}|} \right), \quad (2)$$

236 where  $\Sigma$  is the covariance matrix obtained by taking the whole set of re-  
 237 gions from both positive and negative images. Consequently, the part models  
 238  $w(p, m)$  are fully defined once the match function is defined. In addition, the  
 239 similarity between a region  $r$  and a part  $p$  of the model can be computed as  
 240  $w^T(p, m) \times x_r$ .

241 The model is thus fully defined by giving the match function  $m(r, p)$ . Fol-  
 242 lowing the afore mentioned constraints, we define the optimal match function,  
 243 denoted as  $\hat{m}$ , as the one maximizing:

$$\left\{ \begin{array}{l} \hat{m} = \arg \max_m \sum_{p \in \mathcal{P}} \sum_{I \in \mathcal{I}^+} \sum_{r \in I} m(r, p) \times w^T(p, m) \times x_r \\ s.t. \forall I \in \mathcal{I}^+, \forall p \in \mathcal{P}, \sum_{r \in I} \hat{m}(r, p) = 1 \\ s.t. \forall I \in \mathcal{I}^+, \forall r \in I, \sum_{p \in \mathcal{P}} \hat{m}(r, p) \leq 1. \end{array} \right. \quad (3)$$

244 Learning this model hence consists in the (combinatoric) optimization  
 245 of Eq. (3). Finding the global optimum is not computationally feasible,  
 246 nevertheless we propose to adapt the point matching algorithm of [42] to  
 247 obtain an approximate solution, as explained in the following section. This  
 248 algorithm was first introduced to solve simultaneously the correspondence  
 249 problem as well as the pose estimation of 3D and 2D data. In [42], two sets  
 250 of points  $X_j$  and  $Y_k$  are related by a geometric transformation. Both sets  
 251 can contain outliers. The *match matrix*  $m_{jk}$  is defined as the correspondence  
 252 matrix such that  $m_{jk} = 1$  if point  $X_j$  corresponds to point  $Y_k$  and 0 otherwise.  
 253 The problem is further presented as finding the *pose* (*i.e.* the geometric  
 254 transformation) and the corresponding match matrix  $m_{jk}$  that best relates  
 255 the two sets of points. These two problems are finally solved simultaneously  
 256 using an iterative process aiming at minimizing an energy function.

257 *3.2. Learning with softassign*

258 Now, our main goal is to efficiently find a good (sub-optimal) solution  
 259 of the objective function given by Eq. (3). If we ignore, for the moment,  
 260 the inequality constraint (last constraint of Eq. 3), then the match matrix  
 261  $m$  can be seen as a permutation matrix. We use the *deterministic annealing*  
 262 method of [43] to turn our combinatoric problem into a continuous one, mak-  
 263 ing the optimization simpler and more efficient. The key idea is to minimize  
 264 a sequence of objective functions controlled by a parameter  $\beta$  representing  
 265 the inverse temperature of the system. By increasing the parameter, the  
 266 objective functions leans towards the discrete function.

267 The constraints are then relaxed from a permutation matrix constraints  
 268 to *doubly stochastic matrix* constraints, meaning that every rows and columns  
 269 of the matrix should sum up to 1 (see [42] for more explanations). Therefore,  
 270 the computation of the match function can be achieved iteratively using the  
 271 *softmax* formulation:

$$\forall I \in \mathcal{I}^+, \forall r \in I, m(r, p) = \frac{\exp(\beta \times w^T(p, m^*) \times x_r)}{\sum_{r \in I} \exp(\beta \times w^T(p, m^*) \times x_r)}. \quad (4)$$

272 Where  $w(p, m^*)^T \times x_r$  is the score function relating the similarity between  
 273 the part  $p$  and the region  $r$  of the image  $I$ , using the match function  $m^*$   
 274 computed at the previous iteration. Such a formulation does produce values  
 275 in the interval  $[0, 1]$ , which is expected. Furthermore, when  $\beta \rightarrow \infty$ , there  
 276 will be one region per image for which  $m(r, p) = 1$ , while for the other ones  
 277  $m(r, p) = 0$ , therefore satisfying the first constraint.

278 In practice, we experimentally observed that the previous formulation  
 279 tends to favor parts converging to the same ‘mean part’ for small values of  
 280  $\beta$ . Therefore, we utilized the following formulation, which leads to better  
 281 performance as it encourages sparser representations.  $\forall I \in \mathcal{I}^+$  and  $\forall r \in I$ ,

$$m^\dagger(r, p) = \exp(\beta((w^T(p, m^*) \times x_r) - \max_{\forall r \in I} (w^T(p, m^*) \times x_r))). \quad (5)$$

282 In addition, the match matrix  $m$  has also to satisfy the doubly stochastic

283 constraints. This can be achieved by using Sinkhorn (see more details in  
 284 [42]), by iteratively normalizing rows and columns, see Algorithm 1.

285 Up to this point, we ignored the inequality constraint stating that  $\forall I \in \mathcal{I}^+$   
 286 and  $\forall r \in I, \sum_{p \in \mathcal{P}} m(r, p) \leq 1$ . Gold *et al.* [42] turned the inequality con-  
 287 straint into an equality constraint by adding a slack variable [44]. However,  
 288 unlike Gold *et al.* [42], our problem is not symmetrical. In order to handle the  
 289 inequality constraint, we add non-linearities to the process by setting to zero  
 290 the very low values (*inferior to  $10^{-7}$  in practice*), of  $m^\dagger$  right after its calcula-  
 291 tion, see Equation 5. Then, the following process is the normalization, which  
 292 distributes the weights, except for the null values that remain unchanged.  
 293 Therefore, the normalized match-matrix satisfies the previous constraints:  
 294  $\forall I \in \mathcal{I}^+, \forall p \in \mathcal{P}, \sum_{r \in I} m^\dagger(r, p) = 1$  and  $\forall I \in \mathcal{I}^+, \forall r \in I, \sum_{p \in \mathcal{P}} m^\dagger(r, p) \leq 1$ .  
 295 For example, if a region obtains a very low score for all parts, a column  
 296 of the match-matrix  $m^\dagger(r, p)$  is set to 0. In other words, image regions not  
 297 matching any parts, with very low scores, will not contribute to any parts  
 298 and while the parameter  $\beta$  is increased the selection will be more strict and  
 299 more regions will be discarded. This process further allows to speed up the  
 300 normalizations in the algorithm.

### 301 3.3. Providing initial correspondences between parts and image regions

302 The learning process allowing to learn distinctive parts by iteratively  
 303 refining the match function  $m$ , as presented in the previous section, is a  
 304 process requiring to know  $m$  from the previous iteration, and hence raises  
 305 the question of the initialization of  $m$ .

306 It seems reasonable to think that because the optimization process is not  
 307 convex, the algorithm will perform better if the initial part to regions cor-  
 308 respondences already involve discriminative regions. To select these initial  
 309 discriminative regions, we first extract the signatures  $x_r$  of the regions sam-  
 310 pled from positive training images. These signatures are then clustered, using  
 311 K-means. Then, we use again the LDA acceleration of [41] to learn initial  
 312 classifiers. For each cluster, the classifier  $w$  is defined as  $w = \Sigma^{-1}(\bar{x} - \mu_0)$

313 where  $\bar{x}$  is the average of the signatures within the cluster and  $\mu_0$  and  $\Sigma$  the  
 314 overall mean and covariance matrix.

315 These classifiers are further applied on the regions of the training images.  
 316 Maximum responses to the classifiers are then selected per image and aver-  
 317 aged over positive and negative subsets, giving us the two scores  $s^+$  and  $s^-$ ,  
 318 for a given cluster  $j$ , defined as:

$$\begin{aligned} s_j^+ &= \frac{1}{|\mathcal{I}^+|} \sum_{r^* \in \mathcal{I}^+} w_j^T x_{r^*} \\ s_j^- &= \frac{1}{|\mathcal{I}^-|} \sum_{r^* \in \mathcal{I}^-} w_j^T x_{r^*}. \end{aligned} \quad (6)$$

319 Where  $\forall I \in \mathcal{I}^+$ ,  $r^* = \arg \max_{r \in I} (w_j^T x_r)$ . Then, we denote as  $C_p$  the  $K$   
 320 clusters having the largest  $s_j^+ / s_j^-$  ratios, which are selected as initial discrim-  
 321 inative regions. These initial regions are further used to compute the initial  
 322 part classifier  $w(p, m_0)$  as :  $w(p) \leftarrow \Sigma^{-1}(C_p - \mu_0)$ , used to compute the initial  
 323 match matrix  $m_0(r, p)$ .

### 324 3.4. Computing region and image signatures

#### 325 3.4.1. Image region signatures

326 First, we would like to comment on the patch signatures  $x_r$ , used in the  
 327 learning process. We note that these descriptors must be compact, *i.e.* no  
 328 more than a few thousand dimensions, to allow the learning to be effective.  
 329 In fact, we remind that each image is represented by a few thousand of these  
 330 patches, or regions. Therefore, we first used the simple BOW description  
 331 using  $k = 1000$  clusters, as in [12]. Later, following [25], we extracted the  
 332 seven-th layer of the CNN representation. This intermediate representation  
 333 offers much higher results, as we can see in section 4, and allows a better  
 334 comparison to the current best performing methods. Therefore, we build  
 335 two systems, or pipelines, the standard one or BOW based on SIFT and the  
 336 CNN-based pipeline.

#### 337 3.4.2. Image signatures

338 Once the model is learned, images signatures can be computed using the  
 339 distinctive parts of the model. Let us denote as  $I$  an image to be encoded. We

340 first extract a set of random regions  $r \in I$  and compute their corresponding  
 341 descriptors  $x''_r$ . We can measure to which extent each region is similar  
 342 to one of the model parts by using the scoring function defined previously  
 343 by Eq. (2) as  $w^T(p, m) \times x_r$ , where  $m$  is the match function learned during  
 344 training. Then, we pool the per part similarities to produce a global signature  
 345 of the image.

346 We propose three different pooling/encoding strategies: the Bag-of-parts  
 347 inspired from [11] and two novel approaches so-called the *Fisher-on-parts* and  
 348 the *CNN-on-parts*.

349 *Encoding images with Bag-of-parts.* To compute the bag-of-parts (BOP), the  
 350 per parts scores are computed for each extracted region on an image. The  
 351 signature of the image is then given by aggregating, for each part of the  
 352 model, the average and the maximum of the region scores over the image.  
 353 Namely, if  $p_j$  is one of the  $K$  parts of our model, the signature of the image  
 354  $I$  will be represented by the two following components:

$$\frac{\sum_{r \in I} w^T(p_j, m) \times x_r}{|r \in I|} \quad \text{and} \quad \max_{r \in I} w^T(p_j, m) \times x_r. \quad (7)$$

355 When the problem is a multi-class problem, we do the same for each class  
 356 and aggregate the results. Therefore, we obtain a  $2 \times K \times C$ -dimensional  
 357 descriptor, where  $C$  is the number of classes.

358 *Encoding images with Fisher-on-parts.* Fisher-on-parts (FOP) aims at en-  
 359 coding together the maximum response of each part in an image  $I_t$ . As  
 360 in BOP, scores are computed for each region. Then, instead of aggregating  
 361 average and maximum scores as for the BOP, the maximum scoring region  
 362  $r^*$  for the part  $p$  is selected, as follows:

$$r^* = \arg \max_{r \in I} w^T(p, m) \times x_r. \quad (8)$$

363 Finally, a Fisher vector is computed on the area of the image covered by the  
 364  $K$  selected regions  $r^*$ . Therefore, the final FOP descriptors is  $2 \times G \times D \times C$ -  
 365 dimensional vector, where  $G$  is the number of Gaussian in the mixture model

**Initialization:**  $w(p) \leftarrow \Sigma^{-1}(C_p - \mu_0)$

**while**  $\beta \leq \beta_f$  **do**

**while**  $m^\dagger$  not converged or # of iteration  $\leq I_0$  **do**

**update match matrix by softassign**

Compute  $m^\dagger(r, p)$ , based on Eq. 5

**while**  $\hat{m}^\dagger$  not converged or # of iteration  $\leq I_1$  **do**

$\forall I \in \mathcal{I}^+$

Update  $\hat{m}$  by normalizing rows

$\hat{m}_1^\dagger(r, p) \leftarrow \frac{\hat{m}_0^\dagger(r, p)}{\sum_{r \in I} \hat{m}_0^\dagger(r, p)}$

Update  $\hat{m}$  by normalizing columns

$\hat{m}_0^\dagger(r, p) \leftarrow \frac{\hat{m}_1^\dagger(r, p)}{\sum_{p \in \mathcal{P}} \hat{m}_1^\dagger(r, p)}$

**end**

**update parts using LDA**

Compute  $w(p, m_0^\dagger)$ , based on Eq. 2.

**end**

$\beta \leftarrow \beta_r \beta$

**end**

**Algorithm 1:** Algorithm for learning the mach function.

366 of the Fisher vector,  $D$  is the dimensionality of SIFT descriptors and  $C$  the  
 367 number of categories.

368 *Encoding images with CNN-on-parts.* In this case, regions are encoded with  
 369 CNN features and scores are obtained for each region of an image, as for  
 370 the Bag-of-part signature. For each part or the model, the region giving  
 371 the highest score (see previous paragraphs) is selected and it's descriptor  
 372 kept. All the descriptors so selected are further concatenated resulting in a  
 373  $D^* \times K \times C$ . Where  $D^* = 4096$  is the dimension of the CNN descriptor.

## 374 4. Experiments

375 This section presents an experimental validation of the proposed ap-  
 376 proach. We start by describing the datasets used in our experiments; then we

377 introduce some baseline algorithms used for comparison purposes and give  
378 the details of our implementation; finally the performance obtained with the  
379 proposed approach are exposed and compared to baselines and state-of-the  
380 art algorithms.

#### 381 *4.1. Datasets*

382 Three classification datasets are utilized to experimentally validate the  
383 proposed approach: The Willow actions dataset [18], the Boats Dataset, and  
384 the MIT 67 scenes dataset [19].

385 *The Willow actions dataset [18]* is a dataset for action classification on  
386 unconstrained consumer images from the Internet. The dataset contains 911  
387 images split into 7 classes of common human actions, e.g. ‘running’, ‘riding  
388 cycle’, etc. There are at least 108 images per actions, with 70 images used as  
389 training and the rest as testing images. We note that the dataset also offers  
390 bounding boxes fitted on humans performing the actions. In our case, we  
391 perform the test *without* using these bounding boxes, as we want to detect  
392 the relevant parts of images automatically without any prior knowledge on  
393 the scenes.

394 *The MIT 67 scenes dataset [19]* is composed of 67 categories of indoor  
395 scenes. These categories include stores (e.g. bakery, toy store), home (e.g.  
396 kitchen, bedroom), public spaces (e.g. library, subway), leisure (e.g. restau-  
397 rant, concert hall), and work (e.g. hospital, TV studio). Some scenes can be  
398 best characterized by their global layout (corridor), or by the objects they  
399 contain (bookshop). Each category has around 80 images for training and  
400 20 for testing.

401 *The RECONSURVE Boats Classification Dataset*<sup>1</sup> is composed of 2,877  
402 images divided in 5 categories of boats (e.g. boating, fishing, merchant ship,  
403 tanker, passenger).

404 In the following, the performance on the three datasets is measured using  
405 the mean Average Precision (mAP).

---

<sup>1</sup>can be downloaded from <https://jurie.users.greyc.fr>

406 *4.2. Comparisons to baseline approaches*

407 Our approach is compared to different state-of-the-art approach of the  
408 literature.

409 On one hand we report results obtained with Bag-of-words and Fisher  
410 vectors computed on the dense root SIFT, on the whole image (see [45] for  
411 details). In addition, we used Fisher vectors computed on spatial pyramids,  
412 using the two first layers *i.e.*  $1 \times 1$  and  $2 \times 2$  segments. A SVM classifier is  
413 then trained on the train set and applied to the test images, following the  
414 standard procedures for such image classification tasks.

415 We introduce a second baseline inspired by the work of [13]. A CNN  
416 is first trained with CAFFE on ImageNet (for experiments on the Willow  
417 action dataset) or ImageNet+Places datasets [26] (for the MIT 67 Scenes),  
418 and the penultimate layer of the network is used as an image descriptor.  
419 The images of the target dataset (*i.e.* MIT 67 or Willow) are then encoded  
420 using this CNN-based descriptor and processed with a standard linear SVM  
421 classifier framework.

422 *4.3. Details of our classification pipeline*

423 As explained in the previous section, the proposed algorithm relies on  
424 two steps: a first step in where the parts are learned and a second one  
425 in which a global signature of the image is computed, using the selected  
426 parts. In the first stage, two different encodings of the regions are evaluated  
427 (SIFT based bag-of-words and CNN features). In the second steps 3 different  
428 encoding/pooling strategies are considered: (i) Fisher-on-parts consisting in  
429 computing SIFT-based Fisher vectors on the selected regions, (ii) Bag-of-  
430 parts in which the scores of each part classifier are aggregated to form the  
431 image descriptor and finally (iii) CNN-on-parts in which the CNN descriptors  
432 of each region are concatenated to form the image descriptor (see section 3.4  
433 for a description of these image descriptors).

434 Once images descriptors are computed they are processed by a linear  
435 SVM such as usually done for these classification tasks. We remind that

436 the originality of the work is in the encoding of the image and not in the  
437 classification step which is standard.

438 In the following paragraphs, we give the details of the implementation  
439 used in this experimental validation.

440 *Extraction of image regions.* For each image, a set of regions is generated by  
441 randomly sampling 2,000 regions per image, over the entire image. We note  
442 that for the CNN-based pipeline, only 1,000 regions are extracted to save  
443 time. The scale and aspect ratio of these regions are randomly chosen, but  
444 regions are constrained to have a size of at least 5% of the image size and  
445 aspect ratio should belong to  $[0.5; 2]$ .

446 *Regions descriptors.* As said above, two types of regions descriptors con-  
447 sidered. For the BOW-based region descriptors, dense SIFT features are  
448 extracted within the regions to be encoded, using VLFEAT [46]. We use the  
449 default 4 scales, and sample points every 3 pixels. The SIFT features are  
450 further square-rooted to get rootSIFT features and the feature dimension is  
451 reduced to 80 using PCA, as suggested by [45]. Then each region is charac-  
452 terized using a 1,000-dimensional bag-of-words. These choices are standard  
453 for this type of problem [45].

454 Regarding the CNN descriptors, we use the 7-th layer of the CNN pro-  
455 posed by [47], resulting in a 4,096-dimensional vector. The CNN architecture  
456 is the standard CAFFE architecture [47], and the network is learned on Im-  
457 ageNet for the Willow action dataset or learned on the hybrid ImageNet  
458 and Places datasets (see [26]) for the MIT 67 Scenes dataset. Note that the  
459 same description method is used to compute region descriptors within the  
460 CNN-on-parts descriptor (such as defined section 3.4).

461 *Parameters of the learning algorithm.* Regarding the learning algorithm, we  
462 empirically set the parameters as suggested by [42]:  $\beta = 0.41$ ,  $\beta_r = 1.245$ ,  
463  $\beta_f = 1.2$ ,  $I_0 = 4$ ,  $I_1 = 30$  (see Algorithm 1 for the definition of these  
464 parameters). The algorithm iterates over the estimation of  $m$  until the sum

Table 1: Results for the (SIFT) Bag-of-parts showing the influence of the initialization and the optimization processes, on the Willow dataset.

<b>Method</b>	BOP (random init)	BOP (salient regions init)	BOP on salient regions
<b>mAP</b>	0.460	0.510	0.467

Table 2: Performance of the CNN-based Bag-of-parts descriptors. Left-hand side: using salient regions, Right-hand side: using the proposed learning approach.

<b>Method</b>	BOP on salient regions	BOP (proposed approach)
<b>Willow</b>	0.656	0.766
<b>MIT 67</b>	0.555	0.788

465 over  $m$  of the absolute difference between two iterations is smaller than  $\epsilon =$   
 466 0.005.

#### 467 4.4. Results

468 In this section, we first comment on the quantitative results then show  
 469 some qualitative results, i.e. visualization of learned parts, in Figures 4 and  
 470 5. As said above, the performance is measured using the mean Average  
 471 Precision (mAP).

472 *Initial matchings between model parts and image regions.* First, we evalu-  
 473 ate the impact of the initialization step in the part-learning process, on the  
 474 Willow dataset. Results are given Table 1. The objective is to measure the  
 475 contribution of the initial set of correspondences between parts and regions,  
 476 such as described in section 3.3, and to compare it with a simple random ini-  
 477 tialization of the parts/correspondences. If we randomly initialize the match  
 478 function we observe a mAP of 46.0% (with the SIFT based Bag-of-parts  
 479 encoding). Adding the proposed initialization (based on salient regions) im-  
 480 proves the mAP by 5% (51.0%).

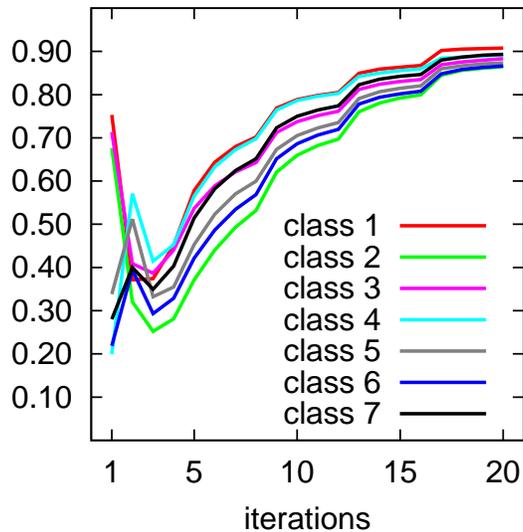


Figure 2: Convergence of the match matrix  $m$  for each class of the Willow actions dataset.

481 In addition, to prove the usefulness of the proposed model versus a simple  
 482 selection of discriminative regions, we evaluated the performance obtained  
 483 by initializing the match function with salient regions (using the method  
 484 proposed in section 3.3), without performing any subsequent optimization,  
 485 *i.e.* without learning  $m$  but keeping the correspondences between parts and  
 486 salient regions as they initially are. If we just use the salient regions, the  
 487 performance drops to 46.7%; we did the same observation for the CNN-based  
 488 pipeline, see Table 2 and Figure 3.

489 The experiments demonstrate that the proposed algorithm improves sig-  
 490 nificantly over a simple selection of discriminative parts, and that a good  
 491 initialization of our algorithm is better than a random initialization of the  
 492 part to region correspondences.

493 *Convergence.* Theses experiments aim at understanding how the match ma-  
 494 trix converges towards a sparse matrix hard-assigning regions to parts. Figure  
 495 2 represents this convergence process, by showing the ratio  $\frac{\sum m_{i,j}^2}{K \times |\mathcal{I}^+|}$ , where  $K$   
 496 is the number of parts and  $|\mathcal{I}^+|$  the number of positive images. The ratio

497 should be of 1 in case of hard assignments. We note that there is a consistent  
498 drop in the first few iterations, as the initial parts are not (yet) constraint  
499 to be generative, *i.e.* a parts should be observed in every positive image of  
500 the specific class. Finally, we observe a small step each time the temperature  
501 parameter is updated.

502 Overall, the convergence is behaving as expected.

503 *Bag-of-word based representation.* The SIFT-Bag-of-parts and Fisher-on-parts  
504 pipelines are then evaluated on the three datasets, see Table 3 and Table 5.  
505 For Willow actions, the performance of the two baseline algorithms (Bag-of-  
506 words and Fisher vectors) are respectively of 50.0% and 58.1%. One can note  
507 that the (SIFT) Bag-of-parts slightly outperforms the standard Bag-of-word.  
508 More interestingly, the proposed Fisher-on-parts representation outperforms  
509 Fisher vectors by more than 3%. Please note that the proposed approach  
510 does not use any extra annotations, contrarily to most of the proposed ap-  
511 proaches (*e.g.* [12] which uses the bounding boxes). This explains why we do  
512 not provide any comparisons with these methods, as they would be mean-  
513 ingless.

514 The Boats dataset also shows improvements on both the (SIFT) Bag-of-  
515 parts and the Fisher-on-parts. Specifically, we observe more than 6% and  
516 2% mAP increase over BOW and Fisher vectors respectively.

517 Concerning MIT 67, we first observe that our (SIFT) Bag-of-parts en-  
518 coding offers better performances than the Bag-of-word model as well as  
519 the Bag-of-parts proposed in [11]. We also notice that our Fisher-on-parts  
520 improves on the two previous methods. However, we do not obtain better  
521 performance than the Fisher vectors extracted on the full image. We believe  
522 that this result is due to the fact that the MIT 67 requires a lot of con-  
523 text information to recognize scenes, while our Fisher-on-parts encoding acts  
524 as a pooling system that encapsulates most information on the foreground.  
525 Combining Fisher-on-parts with Fisher vectors on the whole image (with  
526 SPM) gives a mAP of 60.0%, which is significantly better than any other  
527 approaches.

Table 3: Results on Willow and Boats dataset. See text for details.

Method	Willow (mAP)	Boats (mAP)
Bag-of-words [45]	0.500	0.673
Fisher vectors [45]	0.581	0.827
Bag-of-parts	0.510	0.741
Fisher-on-parts	0.614	0.852

Table 4: Results on Willow dataset, using the CNN-based pipeline.

Method	Willow (mAP)
CNN on full image	0.763
Bag-of-parts (CNN)	0.766
CNN-on-parts (CNN)	0.816
Bag-of-parts & CNN-on-parts (CNN)	0.819

528 *CNN-based pipeline.* In these experiments, the regions are described by CNN  
 529 features. First, we evaluate Figure 3 the impact of the number of parts used  
 530 to describe images, on Willow. It is interesting to note that for (CNN)  
 531 Bag-of-parts and CNN-on-parts, performances are stable for any number of  
 532 parts between 25 to 400 parts. Furthermore, utilizing only 10 parts offers  
 533 reasonable performance. However, if we compute the Bag-of-parts on the  
 534 initialization parts (salient regions), *i.e.* without learning  $m$ , we note that  
 535 having more than 100 parts slightly reduces the performances.

536 We also observed a consistent improvement of the (CNN) bag-of-parts  
 537 and CNN-on-parts over the CNN on the full image, as shown Table 4 and  
 538 Table 5. For the Willow action dataset the CNN-on-part offers the largest  
 539 improvement, while the (CNN) Bag-of-part is the best performing method  
 540 on the MIT 67 scenes dataset. This result supports our observation with  
 541 the Fisher-on-parts that a foreground pooling effect is very advantageous on  
 542 Willow actions, while contextual information is better for MIT scenes.

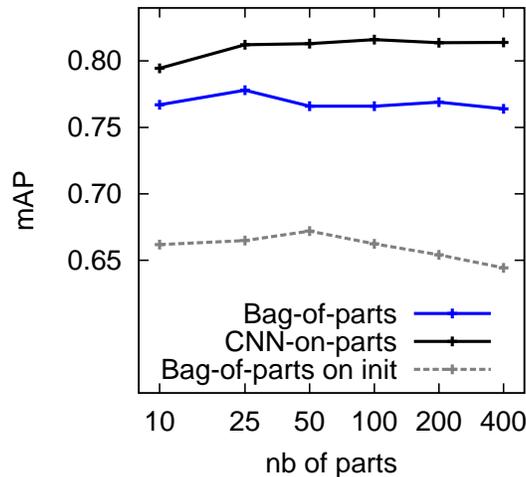


Figure 3: Scores for Willow actions as a function of the number of parts.

543 Interestingly, we also observed that combining the two proposed encoding  
 544 methods – by doing a simple concatenation of the (CNN) Bag-of-part and  
 545 CNN-on-part representations – makes the performance even better, produc-  
 546 ing performance higher than any reported method to our knowledge (80% of  
 547 mAP on MIT67).

548 These experiments show that our descriptors, based on distinctive parts  
 549 learning, are capable of incorporating mid-level information and produce  
 550 richer representations.

## 551 5. Conclusions

552 In this paper, we propose a new algorithm to recognize images by model-  
 553 ing categories as set of distinctive parts that are discovered automatically and  
 554 aligned across images, while learning their visual model. The parts that are  
 555 discovered are free of any appearance constraint and allow the distinction  
 556 between the categories to be recognized. We show how to use the softas-  
 557 sign matching algorithm, to simultaneously learn the part models and assign  
 558 image regions to model’s parts, starting from an initial set of randomly ex-

Table 5: Results on MIT 67 scenes dataset. See text for details.

<b>Method</b>	<b>mAP</b>
bag-of-words [45]	0.345
Fisher vectors [45]	0.550
bag-of-parts of [11]	0.373
our bag-of-parts	0.401
Fisher-on-parts	0.549
Fisher-on-parts based combination	0.600
CNN on full image [26]	0.726
Bag-of-parts (CNN)	0.788
CNN-on-parts (CNN)	0.778
Bag-of-parts & CNN-on-parts (CNN)	0.801

559 tracted image regions. Based on the part model, signatures are computed  
 560 to describe images. Finally, the proposed algorithm is validated on three  
 561 different datasets on which state-of-the-art performances are obtained.

## 562 Acknowledgment

563 This work is partly funded by the RECONSURVE project.

## 564 References

- 565 [1] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with  
 566 deep convolutional neural networks, in: Advances in neural information  
 567 processing systems, 2012, pp. 1097–1105.
- 568 [2] G. Csurka, C. R. Dance, L. Fan, J. Willamowski, C. Bray, Visual cate-  
 569 gorization with bags of keypoints, in: Intl. Workshop on Stat. Learning  
 570 in Comp. Vision, 2004.

- 571 [3] J. Krapac, J. Verbeek, F. Jurie, Modeling spatial layout with Fisher vec-  
572 tors for image categorization, in: International Conference on Computer  
573 Vision, 2011.
- 574 [4] S. Lazebnik, C. Schmid, J. Ponce, Beyond bags of features: Spatial  
575 pyramid matching for recognizing natural scene categories, in: Com-  
576 puter Vision and Pattern Recognition, 2006 IEEE Computer Society  
577 Conference on, Vol. 2, IEEE, 2006, pp. 2169–2178.
- 578 [5] F. Perronnin, J. Sanchez, T. Mensink, Improving the fisher kernel for  
579 large-scale image classification, in: ECCV, 2010.
- 580 [6] K. Chatfield, K. Simonyan, A. Vedaldi, A. Zisserman, Return of the  
581 devil in the details: Delving deep into convolutional nets, arXiv preprint  
582 arXiv:1405.3531.
- 583 [7] M. D. Zeiler, R. Fergus, Visualizing and understanding convolutional  
584 networks, in: Computer Vision–ECCV 2014, Springer, 2014, pp. 818–  
585 833.
- 586 [8] K. He, X. Zhang, S. Ren, J. Sun, Spatial pyramid pooling in deep con-  
587 volutional networks for visual recognition, in: Computer Vision–ECCV  
588 2014, Springer, 2014, pp. 346–361.
- 589 [9] L. Liu, C. Shen, L. Wang, A. van den Hengel, C. Wang, Encoding  
590 high dimensional local features by sparse coding based fisher vectors, in:  
591 Advances in Neural Information Processing Systems, 2014, pp. 1143–  
592 1151.
- 593 [10] S. Maji, G. Shakhnarovich, Part discovery from partial correspondence,  
594 in: Computer Vision and Pattern Recognition, IEEE, 2013, pp. 931–938.
- 595 [11] M. Juneja, A. Vedaldi, C. V. Jawahar, A. Zisserman, Blocks that shout:  
596 Distinctive parts for scene classification, in: Computer Vision and Pat-  
597 tern Recognition, 2013.

- 598 [12] G. Sharma, F. Jurie, C. Schmid, et al., Expanded parts model for human  
599 attribute and action recognition in still images, in: *Computer Vision and*  
600 *Pattern Recognition*, 2013.
- 601 [13] M. Oquab, L. Bottou, I. Laptev, J. Sivic, Learning and transferring  
602 mid-level image representations using convolutional neural networks, in:  
603 *Computer Vision and Pattern Recognition*, 2014.
- 604 [14] J. Krapac, J. Verbeek, F. Jurie, Learning tree-structured descriptor  
605 quantizers for image categorization, in: *BMVC 2011-British Machine*  
606 *Vision Conference*, BMVA Press, 2011, pp. 47–1.
- 607 [15] H. E. Tasli, R. Sircé, T. Gevers, A. A. Alatan, Geometry-constrained  
608 spatial pyramid adaptation for image classification, in: *ICIP*, 2014.
- 609 [16] M. Vidal-Naquet, S. Ullman, Object recognition with informative fea-  
610 tures and linear classification, in: *International Conference on Computer*  
611 *Vision*, 2003, pp. 281–288.
- 612 [17] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, D. Ramanan, Object  
613 detection with discriminatively trained part-based models, *Trans. PAMI*  
614 32 (9) (2010) 1627–1645.
- 615 [18] V. Delaitre, I. Laptev, J. Sivic, Recognizing human actions in still im-  
616 ages: a study of bag-of-features and part-based representations., in:  
617 *BMVC*, Vol. 2, 2010, p. 7.
- 618 [19] A. Quattoni, A. Torralba., Recognizing indoor scenes, in: *Computer*  
619 *Vision and Pattern Recognition*, 2009.
- 620 [20] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, A. Zisserman,  
621 *The PASCAL Visual Object Classes Challenge 2012 Results*.
- 622 [21] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet:  
623 A large-scale hierarchical image database, in: *Computer Vision and*  
624 *Pattern Recognition*, IEEE, 2009, pp. 248–255.

- 625 [22] N. Pinto, Y. Barhomi, D. Cox, J. DiCarlo, Comparing state-of-the-art  
626 visual features on invariant object recognition tasks, in: Applications of  
627 Computer Vision (WACV), 2011 IEEE Workshop on, 2011.
- 628 [23] S. Zheng, Z. Tu, A. Yuille, Detecting object boundaries using low-, mid-,  
629 and high-level information, in: Computer Vision and Pattern Recognition,  
630 2007.
- 631 [24] Y. Boureau, N. Le Roux, F. Bach, J. Ponce, Y. LeCun, Ask the locals:  
632 multi-way local pooling for image recognition, in: ICCV, 2011.
- 633 [25] M. Oquab, L. Bottou, I. Laptev, J. Sivic, et al., Learning and trans-  
634 ferring mid-level image representations using convolutional neural net-  
635 works, Computer Vision and Pattern Recognition.
- 636 [26] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, A. Oliva, Learning Deep  
637 Features for Scene Recognition using Places Database., in: Advances in  
638 Neural Information Processing Systems, 2014.
- 639 [27] Y. Gong, L. Wang, R. Guo, S. Lazebnik, Multi-scale orderless pooling of  
640 deep convolutional activation features, arXiv preprint arXiv:1403.1840.
- 641 [28] D. Parikh, Recognizing jumbled images: the role of local and global in-  
642 formation in image classification, in: International Conference on Com-  
643 puter Vision, IEEE, 2011, pp. 519–526.
- 644 [29] J. Tighe, S. Lazebnik, Superparsing: scalable nonparametric image pars-  
645 ing with superpixels, in: ECCV 2010, Springer, 2010.
- 646 [30] R. Sicre, T. E. Tasli, T. Gevers, Superpixel based angular differences as  
647 a mid-level image descriptor, in: ICPR, 2013.
- 648 [31] B. Fernando, E. Fromont, T. Tuytelaars, Effective use of frequent item-  
649 set mining for image classification, in: Computer Vision–ECCV 2012,  
650 Springer, 2012, pp. 214–227.

- 651 [32] J. Carreira, R. Caseiro, J. Batista, C. Sminchisescu, Semantic segmen-  
652 tation with second-order pooling, in: *Computer Vision–ECCV 2012*,  
653 Springer, 2012, pp. 430–443.
- 654 [33] Z. Niu, G. Hua, X. Gao, Q. Tian, Context aware topic model for scene  
655 recognition, in: *Computer Vision and Pattern Recognition*, IEEE, 2012,  
656 pp. 2743–2750.
- 657 [34] Y. Su, F. Jurie, Improving image classification using semantic attributes,  
658 *International journal of computer vision* 100 (1) (2012) 59–77.
- 659 [35] Z. Liao, A. Farhadi, Y. Wang, I. Endres, D. Forsyth, Building a dictio-  
660 nary of image fragments, in: *Computer Vision and Pattern Recognition*,  
661 IEEE, 2012, pp. 3442–3449.
- 662 [36] J. J. Lim, C. L. Zitnick, P. Dollár, Sketch tokens: A learned mid-level  
663 representation for contour and object detection, *Computer Vision and*  
664 *Pattern Recognition*, 2013.
- 665 [37] I. Endres, K. J. Shih, J. Jiaa, D. Hoiem, Learning collections of part  
666 models for object recognition, *Computer Vision and Pattern Recogni-*  
667 *tion*, 2013.
- 668 [38] J. Chua, I. Givoni, R. Adams, B. Frey, Learning structural element patch  
669 models with hierarchical palettes, in: *Computer Vision and Pattern*  
670 *Recognition*, IEEE, 2012, pp. 2416–2423.
- 671 [39] S. Singh, A. Gupta, A. A. Efros, Unsupervised discovery of mid-level  
672 discriminative patches, in: *ECCV*, Springer, 2012, pp. 73–86.
- 673 [40] R. Sicre, F. Jurie, Discovering and aligning discriminative mid-level fea-  
674 tures for image classification, in: *Pattern Recognition (ICPR)*, 2014  
675 22nd International Conference on, IEEE, 2014, pp. 1975–1980.
- 676 [41] J. M. B. Hariharan, D. Ramanan, Discriminative decorrelation for clus-  
677 tering and classification, in: *ECCV*, 2012.

- 678 [42] S. Gold, A. Rangarajan, C.-P. Lu, S. Pappu, E. Mjolsness, New algo-  
679 rithms for 2d and 3d point matching:: pose estimation and correspon-  
680 dence, *Pattern Recognition* 31 (8) (1998) 1019–1031.
- 681 [43] D. Geiger, F. Girosi, Parallel and deterministic algorithms from mrfs:  
682 Surface reconstruction, *Trans. PAMI* 13 (5).
- 683 [44] V. Chvatal, *Linear programming*. 1983, WH Freeman and Company,  
684 New York.
- 685 [45] K. Chatfield, V. Lempitsky, A. Vedaldi, A. Zisserman, The devil is in the  
686 details: an evaluation of recent feature encoding methods, in: *British*  
687 *Machine Vision Conference*, 2011.
- 688 [46] A. Vevaldi, B. Fulkerson, Vlfeat an open and portable library of com-  
689 puter vision algorithms, in: *ACM Multimedia*, 2010.
- 690 [47] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick,  
691 S. Guadarrama, T. Darrell, Caffe: Convolutional architecture for fast  
692 feature embedding, *arXiv preprint arXiv:1408.5093*.



Figure 4: Visualization of parts locations for "riding cycle", "playing instruments", and "riding horse".



Figure 5: This figure shows the highest scoring regions for a set of parts learned for the *riding horse* action. Each row corresponds to a part.