

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011

## Abstract

Recognition of social styles of people are an interesting but not yet *poular* topic. In this paper, we explore the use of pre-trained convolutional neural network(CNN) features for social group recognition. A social group recognition framework is proposed. Pre-trained CNN is fine-tuned, and CNN features are used for both individual person images and global scene images. Our model shows promising results on urban tribes dataset, with 71.23% accuracy, which significantly outperforms previous results of 46%. We also find that there is a correlation between the probability of an image being in Imagenet classes and social group classes, and that better-recognized categories have more highly-correlated ImageNet categories. This gives us insight into the features extracted from pre-trained CNNs.

## 1. Introduction

In the past few years, there have been impressive progress in understanding semantic meaning of images, such as object recognition, scene recognition, and object detection. The power of Convolutional Neural Networks (CNNs) has is especially noticeable. However, the analysis the social features of images of groups of people has not attracted a great deal of research. Current search algorithms fail to capture information of personal styles or social characteristics of groups of people, but retrieve images with similar global appearance [1]. The analysis of groups of people is difficult in that the group categories are semantically ambiguous, and have high intra-class variance.

Recognition of groups of people from a social perspective provides many potential applications. With more accurate group searching results, more accurate recommendations can be made in social networks, and more relevant advertisement for particular groups of people benefits both consumers and sellers.

Kwak *et al.* studied this problem of group recognition ([2], [1]). They created an urban tribe dataset consisting of 11 categories, with about 100 labeled per class. They proposed a group recognition pipeline. Rather than classifying

## Urban Tribe

Anonymous WACV submission

Paper ID \*\*\*\*

isolated individuals in the group images, they focused on group features and models.

CNN architecture has been proved to achieve outstanding results in various computer vision tasks, and it's argued that deep architectures in CNN can capture visual features of different semantic level in the hidden units ([3]). Recently, features learnt for large scale recognition tasks with large amounts of training data have been used for new tasks, and the features outperform many conventional features in many tasks([3], [4]). This indicates that generic visual features may be obtained from pre-trained CNN model.

not sure about this paragraph The most frequently used pre-trained model was proposed in [5] trained with ImageNet dataset<sup>1</sup>. And the new tasks for which CNN features are used range from object recognition, scene recognition, to subcategory recognition. However, the recognition of style hasn't been researched yet. The social groups of individuals require the recognition of people's social characteristics, which is a problem of style, instead of category. Moreover, the ImageNet classes used for pre-training cover few human images, making the problem more interesting.

In this paper, we investigate the generalization ability of pre-trained CNN features to social group recognition. We propose a CNN feature based architecture for social group recognition. Our model takes in both individual features and global scene features fine-tuned from CNN pre-trained weights. Our result shows a boost of performance from the previous classification method provided by [1]. We show that both individual information and global scene information contribute to a social group's characteristics, and that different feature extraction schemes for individual and global information is necessary. We also show the role of adapting generic pre-trained CNN features to social group images with fine-tuning.

We further investigate why features extracted from pre-trained CNN are useful for the urban tribe recognition task. For an input image, there is a correlation between the probability of it being in ImageNet classes and being in urban tribes classes. Moreover, the degree of correlation is related to the recognition rate of different urban tribes classes

<sup>1</sup><http://image-net.org/challenges/LSVRC/2012/browse-synsets>

108 - better-recognized categories have more highly-correlated  
109 ImageNet categories. This may indicate that the “generic”  
110 features extracted by pre-trained CNN networks are not so  
111 generic. The network is trained to separate the ImageNet  
112 classes most, and if we use the features for a new classifica-  
113 tion task, the performance of the task is related to how well  
114 the new classes can be “mapped” to the ImageNet classes.  
115 However, the actual relationship between the two types of  
116 categories is still mysterious in most cases.  
117

## 118 2. Related Work

119 Convolutional Neural Networks(CNNs) with back-  
120 propogation were introduced around 1990’s by LeCun *et*  
121 *al.* [6]. Since then, CNNs have shown successful results on  
122 various computer vision tasks, such as hand-written digit  
123 classification ([7]), ImageNet Large Scale Visual Recog-  
124 nition Challenge ([5], [8]), object detection ([9],[8]), object  
125 localization ([9]).

126 Recently, many researchers have shown the utility of  
127 generalization of pre-trained CNN features on large dataset  
128 such as ImageNet. Krizhevsky *et al.* showed excellent gen-  
129 eralization of the pre-trained CNN features[3]. They kept  
130 all the layers of ImageNet-trained model using fixed pre-  
131 trained features except for the last softmax classifier, and  
132 achieved best results on Caltech-101 and Caltech-256. Don-  
133 ahue *et al.* used different layers of pre-trained CNN network  
134 as features and trained simple classifiers such as SVM and  
135 Logistic Regression, and outperformed the state-of-the-art  
136 on several vision challenges such as scene recognition and  
137 domain adaptation[4].

138 Kwak *et al.* created an urban tribe dataset consisting of  
139 11 classes[1]. The classes are defined from social group  
140 labels provided by Wikipedia. They selected the eight most  
141 popular categories from their list of subculture, and added  
142 three other classes corresponding to typical social venues  
143 in addition. For each class, images of groups of people  
144 were discovered with different search engines, and a broad  
145 range of images for each class were collected. Kwak *et al.*  
146 also provided a group description and several classification  
147 methods([2], [1]). Group description consists of person de-  
148 scriptors and global group descriptors: Six part of person is  
149 detected, and a set of predefined descriptors are computed  
150 for each part, including ratio of skin pixels, color informa-  
151 tion like RGB histograms, and HoG features; Global de-  
152 scriptors use a both low level and high level descriptors to  
153 describe the context and group properties of the image. Low  
154 lever features include color information, Gist, HoG and  
155 ratio of pixels of person, and high level descriptors includs  
156 proximity of persons, alignment or pose of the group, and  
157 scene layout of individuals. Two options of classification  
158 methods are provided: bag of parts-based classification and  
159 SVM-based classification.  
160

161 Categorizing the social groups of individuals belongs to

162 fine-grained classification task which recently draws more  
163 interest of the computer vision community. It aims at giving  
164 the fine-grained categories in a certain class. Fine-grained  
165 classification is more difficult than conventional classifica-  
166 tion tasks, because the categories are semantically as well as  
167 visually similar, and are even challenging for humans. The  
168 Fine-grained Challenge 2013 (FGComp) provided the data  
169 in several categories including aircraft, birds, cars, dogs,  
170 and shoes. [10] achieved the best result using classifier  
171 based on fisher vectors. However, CNN based methods us-  
172 ing [11] or [4] gave inferior results, especially when the  
173 bounding box of test data is unknown.

174 There is some research in analyzing social groups of peo-  
175 ple. [12] showed the visual structure of a group helps under-  
176 standing events. [13] showed social relationships modeling  
177 helps people recognition. [14] used both local and global  
178 factors for group level expression analysis.

## 179 3. Methods

180 This section describes the urban tribe dataset and elabo-  
181 rates on the model architecture.

### 182 3.1. Urban tribes dataset

183 Urban tribes are groups of people who have similar vi-  
184 sual appearances, personal style and ideals. The urban  
185 tribes dataset consists of 11 different categories: *biker*,  
186 *country*, *goth*, *heavy-metal*, *hip-hop*, *hipster*, *raver*, *surfer*,  
187 *club*, *formal*, *casual/pub*, with an average of 105 images  
188 from each category.

189 Unlike conventional visual classification problems, ur-  
190 ban tribe categories are more ambiguous and subjective.  
191 Also, each class contains a broad range of scenarios. The  
192 high intra-class variation of the urban tribe dataset makes  
193 the classification task challenging. The urban tribe dataset  
194 also has some interesting properties. The number of peo-  
195 ple in each urban tribe image varies. Members in one tribe  
196 often have similar visual styles, including their clothes,  
197 accessories, and even demeanor. For example, surfers pos-  
198 sibly carry surfboards, and the goth often have dark attire,  
199 makeup and hair. The environment they are in also con-  
200 tributes to each tribe characteristics: pictures of country  
201 tribes are more likely to be taken outdoors with grassland,  
202 while pictures of clubbers are often photographed in clubs  
203 with dim lightings.

### 204 3.2. Classification hierarchy

205 To utilize the properties of urban tribes fully, our fea-  
206 ture vector consists of both elements: individual features  
207 and environmental features. For each feature type, we use  
208 a similar extraction strategy. Individual features and envi-  
209 ronmental features are hierarchically combined to form the  
210 final decision function. The network hierarchy is shown in  
211 Figure 1.

216 For each group image, we represent the group  $G$  as the  
 217 combination of a set of people and the environment of the  
 218 scene. To give the prediction of class  $C$ , the individual fea-  
 219 ture vectors and scene feature vectors are extracted sepa-  
 220 rately.  
 221

222 For the individual feature vectors, first, individual can-  
 223 didate person images are detected with a poselet based  
 224 person detection algorithm. The candidate person images  
 225  $H = \{H_1, H_2, \dots, H_p\}$  are used as a whole instead of a set  
 226 of body part bounding boxes. Each candidate person is re-  
 227 sized to  $256 \times 256$ , and ten  $227 \times 227$  patches  $\{h_{ij}\}, i \in$   
 228  $\{1, 2, \dots, p\}, j = 1, 2, \dots, 10$  are extracted (patches from  
 229 four corner and the center, and their horizontal reflections).  
 230

231 Each Individual image patch  $h_{ij}$  then passes through  
 232 the Convolutional Neural Network for person images  
 233  $\text{CNN}_{\text{Person}}$ , generating activations from the 6th and 7th  
 234 hidden layer. The activations from 6th and 7th layer are  
 235 both 4096 dimensioned. They are concatenated to form an  
 236 8192-dimensional vector  $f_{ij}$ , where  $i \in \{1, 2, \dots, p\}, j \in$   
 237  $\{1, 2, \dots, 10\}$ .  
 238

239 The feature vectors are then fed into a multi-class  
 240  $\text{SVM}_{\text{Person}}$ . We use LIBLINEAR[15] to train the SVM  
 241 on individual patches, and to estimate probabilities for  
 242 each category given individual patch  $h_{ij}$ :  $\Pr_{ij}(C|h_{ij}), C \in$   
 243  $\{1, 2, \dots, c\}$ , where  $c$  is the number of classes in urban tribe  
 244 dataset. The individual patches  $h_{ij}$  in one group image are  
 245 usually highly correlated. Therefore, in order to obtain a re-  
 246 liable probability estimate from the noisy yet correlated set  
 247 of probabilities  $\Pr_{ij}$ , a simple but effective average pooling  
 248 is performed to  $\Pr_{ij}$ :

$$\Pr_{\text{People}}(C|H_1, \dots, H_p) = \frac{1}{10p} \sum_{i,j} \Pr_{i,j}(C|h_{ij}) \quad (1)$$

249  $\Pr_{\text{People}}(C|H_1, \dots, H_p)$  is the probability estimate of class  
 250  $C$  given the set of people candidate images  $H$ .  
 251

252 On the other hand, the entire environment in the scene  
 253 image, denoted by  $S$ , is directly utilized for probability es-  
 254 timation. The procedure to generate probability estimate of  
 255 class  $C$  given the environment as a whole  $\Pr_{\text{Scene}}(C|S)$  is  
 256 similar with that for  $\Pr_{\text{People}}(C|H)$ . The difference is that  
 257 the input of Convolutional Network is  $227 \times 227$  patches  
 258 extracted from the entire scene image, and the fine-tuned  
 259 Convolutional network:  $\text{CNN}_{\text{Scene}}$  and SVM:  $\text{SVM}_{\text{Scene}}$   
 260 are trained with the training set of entire scene images. Sev-  
 261 eral different strategies to extract patches from scene images  
 262 and corresponding Convolutional Neural Network architec-  
 263 tures are explained in Section 3.3.  
 264

265 Therefore, the probability estimate of a class  $C$  given  
 266 observation of scene  $S$  is:  
 267

$$\Pr_{\text{Scene}}(C|S) = \frac{1}{K} \sum_{k=1}^K \Pr_k(C|s_k) \quad (2)$$

268 where  $K$  is the number of scene patches extracted from one  
 269 group image, and this number varies with different patch  
 270 extraction strategies.  $s_k$  is the  $k$ th scene patches.  $\Pr_k(C|s_k)$   
 271 is the probability for class  $C$  given  $k$ th scene patch. Average  
 272 pooling is still used here, because the assumption of high  
 273 correlation in patches holds.  
 274

275 Now we have the estimates of two kinds of conditional  
 276 probability  $\Pr_{\text{People}}(C|H)$  and  $\Pr_{\text{Scene}}(C|S)$ . We make a  
 277 strong assumption that the two types of features are inde-  
 278 pendent, and that the prior probability distribution of the  
 279 urban tribes  $\Pr(C)$  is a uniform distribution. The classifica-  
 280 tion problem can be expressed as maximizing the objective  
 281 function:  
 282

$$L = \arg \max_{i=1, \dots, c} \Pr(C = i|G) \quad (3)$$

284 where  
 285

$$\begin{aligned} \Pr(C = i|G) &= \Pr(C = i|H, S) \\ &= \frac{\Pr_{\text{People}}(C = i|H_1, \dots, H_p) \cdot \Pr_{\text{Scene}}(C = i|S)}{\Pr(C = i)} \\ &\propto \Pr_{\text{People}}(C = i|H_1, \dots, H_p) \cdot \Pr_{\text{Scene}}(C = i|S) \end{aligned} \quad (4)$$

286 and  $L$  is the predicted label for the group image.  
 287

### 3.3. Convolutional network feature extraction

288 It is shown in many experiments that a set of weights of  
 289 convolutional network trained from ImageNet can generate  
 290 a set of generic visual features.  
 291

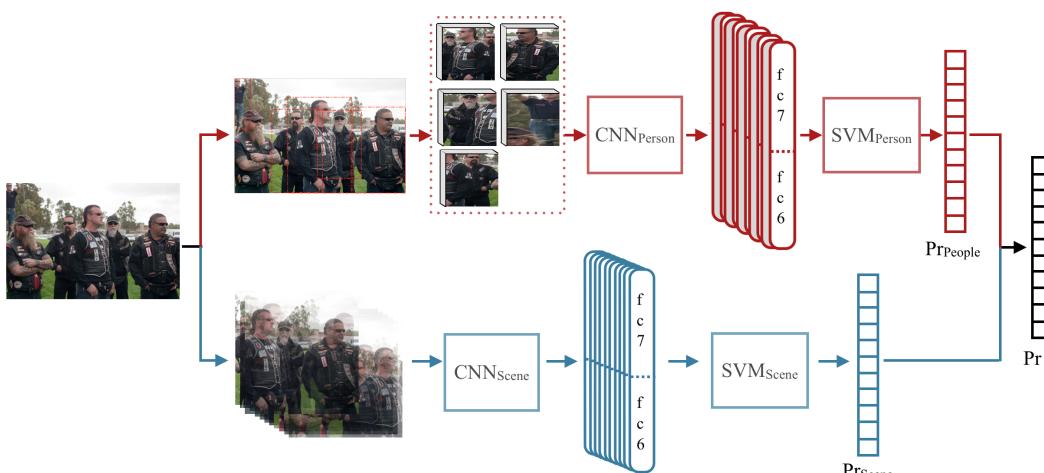
292 Following [11]'s work, we use the network framework  
 293 called Caffe. The network architecture is described in [5],  
 294 which won the ImageNet Large Scale Visual Recognition  
 295 Challenge 2012. We take the activations from the 6th and  
 296 7th hidden layer of the convolutional neural network, which  
 297 are two fully connected layers before the class prediction  
 298 layer. We also take the activations from 6th or 7th layer  
 299 alone as comparison. We choose these two layers, because  
 300 as the layers ascend, the features extracted show increasing  
 301 invariance and semantic meaning.  
 302

303 We use the pre-trained set of weights of the network re-  
 304 leased by Caffe as the initial parameters of our network. The  
 305 pre-trained model was trained on ImageNet ILSVRC-2012,  
 306 and all images are first resized to  $256 \times 256$  before they can  
 307 be used as inputs to the network.  
 308

#### 3.3.1 Pre-processing of the dataset

309 The urban tribe dataset is a relatively small dataset, and both  
 310 people candidate crops and scene images are of various res-  
 311 olution. Our convolutional neural networks requires con-  
 312 stant input image size of  $227 \times 227$ , so pre-processing of  
 313 the dataset is necessary.  
 314

315 There are several strategies to make one image compati-  
 316 ble with the CNN:  
 317

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341342 Figure 1: Architecture of classification algorithm using  $\text{Nets}_{SDense}$  which is introduced in Section 3.3.3. The upper half  
343 estimates the probability given people candidate images, and the lower half estimates the probability given the entire scene.  
344 Dense crop  $\text{CNN}_{Scene}$  and distorted crop  $\text{CNN}_{Person}$  are used, as shown in Section 3.3.1.  
345  
346

1. *Distort cropping:* As in [5], resize the image to a fixed resolution of  $256 \times 256$ , and crop five  $227 \times 227$  patches (from four corner and the center) and their horizontal reflections to generate ten patches from one single image. This way, the aspect ratio of the original images are lost, but for each crop, the portion it takes from the original image is fixed, so that the amount of information all the crops have is relatively stable.
2. *Sparse cropping:* Keep the aspect ratio of the original image, resize the shorter side to 256, and then crop five  $227 \times 227$  patches and their horizontal reflections as mentioned. This method avoids distortion of the image and objects in it, but the crops will possibly lose much information when the aspect ratio of the original image is far away from 1.
3. *Dense cropping:* Keep the aspect ratio of the original image, resize the shorter side to 256, and then densely crop multiple  $227 \times 227$  patches and their horizontal reflections. This way, the information of original image is kept by dense cropping process, and the distortion is avoided. The number of crops attained with this method is larger than the previous two methods.

371 

### 3.3.2 Network Fine-tune

373 Although the pre-trained network from [11] can already  
374 generalize well to many datasets, the urban tribe dataset has  
375 its unique property. It emphasizes certain visual features  
376 such as certain clothing styles, while pays less attention to  
377 other visual features. Also, it emphasizes style of the object

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

rather than distinct category. To rearrange the importance of different features and adjust the features to adapt to urban tribe dataset, the network can be fine-tuned.

The dataset used for fine-tune is the same set used for SVM training. The input patches of the Convolutional network is of size  $227 \times 227$ . The initial Convolutional network has 1000 outputs in final layer, corresponding to 1000 class-wise probability predictions. In our fine-tune process, the last layer is replaced by 11 probability prediction outputs, and the initial weights of the last layer connection are initialized to have zero mean gaussian distribution. Back propagation is used, and the learning rate is set to be small so that the fine-tune process adapts the extracted features to urban tribe dataset while preserving the initial property in general: the initial learning rate used for pre-training is 0.01; We set the initial fine-tune learning rate of the parameters except for the last layer as 0.001, and keep the initial learning rate for last layer 0.01, because the last layer is not pre-trained.

421 

### 3.3.3 Choices of network combination

Scene images and individual images have different properties, and need different strategies of pre-processing and separate fine-tuning. For scene network  $\text{CNN}_{Scene}$  and scene images, due to the small size of the dataset, we use the *dense cropping* technique, the third technique in Section 3.3.1, to increase the dataset. For person network  $\text{CNN}_{Person}$  and corresponding input, we use the *distort cropping* technique, because the subimages have normally long height and short width, and the second and third strategies using squared

432 crops of a person image will lose much information, no  
 433 matter which location we choose to crop them; whereas the  
 434 first method ensure each crop keeps the essential features  
 435 for classification.  
 436

437 The combination of dense crop  $\text{CNN}_{\text{Scene}}$  and distorted  
 438 crop  $\text{CNN}_{\text{Person}}$  are denoted as  $\text{Nets}_{\text{SDense}}$ .  
 439

440 We also construct other combination of networks for  
 441 comparison:  
 442

- 443 1.  $\text{Nets}_{\text{NoTune}}$ : Directly use the pre-trained network by  
 [11] for both scenes and persons, and use the *distort*  
*cropping* technique (distorted crops) as input patches  
 for both networks. This choice of cropping strategy is  
 in consistent with the way the network is pre-trained.  
 444
- 445 2.  $\text{Nets}_{\text{SSparse}}$ : Use the *sparse cropping* strategy for  
 scene features, and the *distort cropping* strategy for  
 person features. .  
 446
- 447 3.  $\text{Nets}_{\text{SDistort}}$ : Use the *distort cropping* strategy for  
 both scene features and person features.  
 448

## 4. Experiments and Results

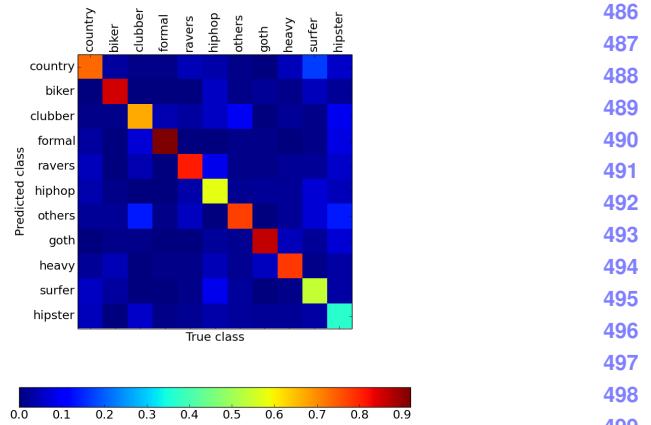
449 In this section, the performance of the proposed classi-  
 450 fication scheme is evaluated and analyzed. In the experi-  
 451 ments, six rounds of 5-fold cross validation are performed,  
 452 therefore we have 30 training experiments in total. Dataset  
 453 is partitioned into 5 equal sized subsets, containing one-fifth  
 454 of the data points from each category. One of the subsets is  
 455 used as test set, and the remained 4 subsets are used as train-  
 456 ing data. In the fine-tune procedure, there are 6000 training  
 457 iterations, and learning rate is decreased by ten times after  
 458 each 1000 iterations.  
 459

### 4.1. Urban tribe classification performance

460 Table 1 shows the comparison of performance using dif-  
 461 ferent approaches. The 30 segmentations of datasets are  
 462 used for all the approaches tested in this section, and 30 test  
 463 results are averaged for each approach. The standard error  
 464 is shown with the accuracy in Table 1. We also compare our  
 465 result with the result achieved by [1] using their best model.  
 466 The advantage of CNN pre-trained features is obvious.  
 467

468 The confusion matrix is shown in Figure 2 for  
 469  $\text{Nets}_{\text{SDense}}$ , and all the 30 training experiments are aver-  
 470 aged. We can observe there is a obvious difference of dif-  
 471 ficulty of different categories. Class *formal* has accuracy as  
 472 high as about 90%, while class *hipster* is the most difficult  
 473 class, having less than 60% accuracy.  
 474

475 Comparing the result of using different features in the  
 476 same approach shows the necessity of every step of our  
 477 architecture. In results using  $\text{Nets}_{\text{SDense}}$  with concatenated  
 478 features, average accuracy for each candidate person is low  
 479 as 47.10%. Average pooling of candidate person probability  
 480 estimates produces a large accuracy increase of about 20%.  
 481



498 Figure 2: Confusion matrix for classification results with  
 499  $\text{Nets}_{\text{SDense}}$ , using people and scene features.  
 500

501 Accuracy using the entire scene only results in 67.26% ac-  
 502 curacy. Combining probabilities  $\Pr_{\text{People}}(C|H_1, \dots, H_p)$   
 503 and  $\Pr_{\text{Scene}}(C|S)$  achieves accuracy as high as 71.22%,  
 504 which verifies the complementary role of people candidate  
 505 feature and environment feature in a group image.  
 506

507 We also compare vertically the accuracy of different ap-  
 508 proaches, to show the role of network feature concatena-  
 509 tion. Using only 7th layer or 6th layer activation from  
 510 the networks  $\text{Nets}_{\text{SDense}}$  produces decent results, showing  
 511 both layers’ activations can generate high semantic fea-  
 512 tures. Concatenating both layers’ activations increases the  
 513 accuracy by 0.5%, indicating the slight information loss of the  
 514 7th fully connected layer.  
 515

516 To see the role of fine-tuning, we can compare the result  
 517 of  $\text{Nets}_{\text{NoTune}}$  and  $\text{Nets}_{\text{SDistort}}$ . These two approaches  
 518 both use resizing that causes distortion, and they only vary  
 519 in fine-tune procedure. There is a large performance im-  
 520 provement with fine-tune, both for person features and  
 521 scene features. This shows the benefit of adapting existing  
 522 generic Convolutional Network to specific dataset.  
 523

524 Nets<sub>SDistort</sub>, Nets<sub>SSparse</sub>, and Nets<sub>SDense</sub> use different  
 525 patch extraction strategies. Note that we use the same dis-  
 526 torted patch extraction method for person images, as men-  
 527 tioned in Section 3.3.3, while we use three different meth-  
 528 ods for scene images. The results for scene images show  
 529 the advantage of keeping the aspect ratio of scene images,  
 530 and the slight advantage of using dense crops. However, the  
 531 final results with People+Scene for the three methods don’t  
 532 have significant differences, this is due to the combination  
 533 with people information.  
 534

Table 1: Performance of different approaches using different information.

Accuracy (%)	Individual candidate	People	Entire scene	People+Scene
Nets <sub>SDistort</sub> with concatenated features	$39.99 \pm 0.30$	$64.09 \pm 0.63$	$62.77 \pm 0.52$	$69.28 \pm 0.50$
Nets <sub>SDense</sub> with fc7 features	$47.07 \pm 0.30$	$67.09 \pm 0.52$	$65.05 \pm 0.42$	$70.74 \pm 0.47$
Nets <sub>SDense</sub> with fc6 features	$45.84 \pm 0.34$	$66.20 \pm 0.46$	$67.08 \pm 0.51$	$70.43 \pm 0.46$
Nets <sub>SDense</sub> with concatenated features	<b><math>47.10 \pm 0.34</math></b>	<b><math>67.29 \pm 0.50</math></b>	<b><math>67.26 \pm 0.50</math></b>	$71.22 \pm 0.46$
Nets <sub>SSparse</sub> with concatenated features	<b><math>47.10 \pm 0.34</math></b>	<b><math>67.29 \pm 0.50</math></b>	$66.81 \pm 0.42$	<b><math>71.23 \pm 0.49</math></b>
Nets <sub>SDistort</sub> with concatenated features	<b><math>47.10 \pm 0.34</math></b>	<b><math>67.29 \pm 0.50</math></b>	$65.35 \pm 0.37$	$71.15 \pm 0.50$
SVM <sub>8</sub> [1]	-	-	-	46(std: 2)

## 4.2. Convolutional Network feature analysis

**Yufei: Is this section necessary?** We use t-distributed stochastic neighbor embedding technique [16] to visualize the power of different layers' features in convolutional network. Two dimensional embedding of the high dimensional CNN feature space is extracted, and we plot the data as 2-d points with different colors indicating different classes they belong to. Powerful features will pull data of different semantic classes more apart.

We randomly choose one training-test partition and corresponding fine-tuned network parameters, and to avoid overfitting effect, we examine the test set of Nets<sub>SDense</sub> approach. In Figure 3, each data is plotted as a dot in each figure. The three columns correspond to data separation in first layer, fourth layer and seventh layer respectively. The first row visualizes features of all test data in scene CNN CNN<sub>Scene</sub>, and second row picks three classes from the first row. The third row is features of CNN<sub>Person</sub>, and the three classes are picked for the last row. The three classes chosen in second and last row is (goth:red, heavy:green, surfer:blue).

In both CNN<sub>Scene</sub> and CNN<sub>Person</sub>, there is a clear trend of class separation. As the layer ascends, , the data from same class are more concentrate, and inter-class distance are larger. The selected three classes show the trend more clearly: they essentially form three clusters in the seventh layer. As shown in the second and last rows, person data is more challenging: The semantic classes are not as separate as with the scene features.

## 4.3. Urban tribe classes vs. ImageNet classes

It's recently being acknowledged that CNN pre-trained features are generic and can be used for new tasks. In this section we check the relationship between the new tasks and original ImageNet task, using the urban tribe dataset, which gives us some insight about the features extracted from pre-trained network.

The urban tribe dataset contains groups of people, and the important features for categorization are mainly human related features, such as attire, make up, posture and expressions. However, the ILSVRC dataset used for pre-training

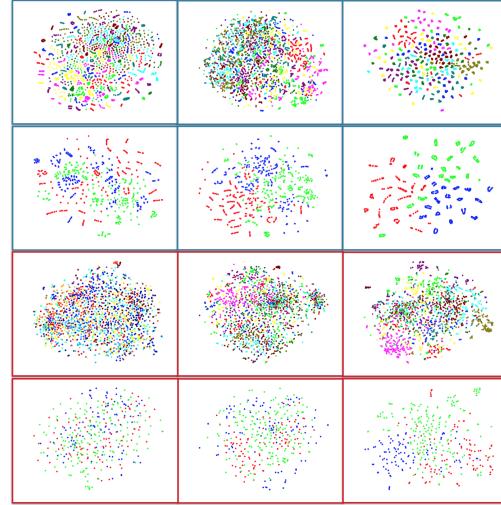


Figure 3: Feature visualization of Nets<sub>SDense</sub>: CNN<sub>Scene</sub> (first two rows) and CNN<sub>Person</sub>(second two rows). Rows: first: all scene images; second: all data from goth/heavy/surfer urban tribes; third: all test person images; fourth: all test data from goth/heavy/surfer urban tribes. Columns: first: first convolutional layer; second: fourth convolutional layer; third: seventh fully convolutional layer

contains few images of humans. Instead of examining the output of the layers directly, we try to find the relationship between the 1000 classes in ILSVRC dataset and the classes in two image sets: urban tribe dataset, and candidate person images extracted from urban tribe dataset.

We use the parameters of pre-trained CNN network as our feature extraction model, and train a softmax layer on top of its 7th layer to predict the probabilities of one input image(either scene image or candidate person image) being in certain urban-tribe class  $\text{Pr}(l_{\text{urban}})$ , where  $l_{\text{urban}}$  pre-trained is the 11 urban tribe categories. The output layer is trained for 3000 training iterations. We can also use the output of the pre-trained CNN network to predict probabilities of one input image being in certain ImageNet category, denoted as  $\text{Pr}(l_{\text{ImageNet}})$ , where  $l_{\text{ImageNet}}$  is the 1000 Im-

648  
649  
650  
651  
ageNet categories. We use one round of 5-fold cross validation, and use all the images in urban tribe dataset for analysis.  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661

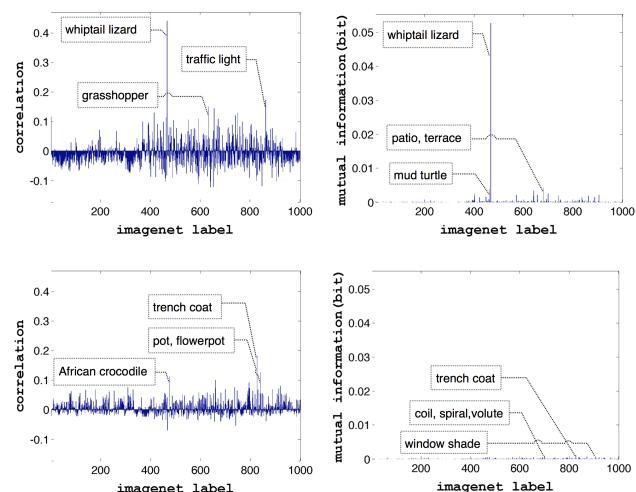
We first check the relationship of  $\Pr(l_{urban})$  and  $\Pr(l_{ImageNet})$  of candidate person images. We calculate the correlation coefficient of  $\Pr(l_{urban})$  and  $\Pr(l_{ImageNet})$  for all 1000 ImageNet classes, denoted as  $R(\Pr(l_{urban}), \Pr(l_{ImageNet}))$ . We also calculate the 1000 mutual information of the predicted score of urban tribe and ImageNet (where predicted score is 1 if the predicted label is the category being tested, 0 otherwise), denoted as  $I(l_{urban}, l_{ImageNet})$ .

In Figure 4, we choose two urban classes: *biker* and *hipster*, and plot the correlation coefficient  $R$  and mutual information  $I$ . The first row shows the result of *biker*, and second row *hipster*. For *biker*, there are several impulses in correlation plot, and one significant impulse in mutual information plot. *whiptail lizard* has high correlation with *biker*. Meanwhile, for *hipster*, which is the most difficult class, the correlation coefficient and mutual information are both low for all ImageNet classes.

To confirm the correlation, we use  $\Pr(l_{ImageNet})$  directly as features, substituting the concatenated fc7fc6 features, and use the Nets<sub>SDistort</sub> approach for classification. The accuracy is 52.68%. This decent result indicates the relationship between  $l_{urban}$  and  $l_{ImageNet}$ . Then, we check  $l_{ImageNet}$  with highest correlation coefficients with  $l_{urban}$ . In Figure 5, we choose four  $l_{urban}$ : *formal*, *ravers*, *goth*, *hipster*, and choose some examples of candidate person images that have both high  $\Pr(l_{ImageNet})$  and high  $\Pr(l_{urban})$ . We also show examples of the images in  $l_{ImageNet}$ . We can see some of the shared features between corresponding person images and ImageNet images, for example, similar shape for women tops and stingrays (Figure 5b).

There is a correlation between class-wise accuracy of predicted  $l_{urban}$  and the degree of relationship between predicted  $\Pr(l_{urban})$  and  $\Pr(l_{ImageNet})$ , as shown in Figure 6. Class-wise accuracy is calculated for candidate person images (Figure 6a, 6b) or scene images (Figure 6c, 6d). For each  $l_{urban}$ , the maximum correlation/mutual information over 1000 ImageNet classes are used to indicate the degree of its relationship with  $l_{ImageNet}$ .

The correlation between ImageNet class and urban tribe class and its relationship with class-wise recognition rate may indicate that the “generic” features extracted by pre-trained CNN networks are not so generic. The network is trained to separate the ImageNet classes most, and if we use the features for a new classification task, the performance of the task is related to how well the new classes can be “mapped” to the ImageNet classes.



702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
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 4: Correlation and mutual information of  $l_{urban}$  and 1000  $l_{ImageNet}$ . The first row is for  $l_{urban}=\text{Biker}$ . The second row is for  $l_{urban}=\text{Hipster}$ . Top three ImageNet classes are marked.

## 5. Conclusion

In this work, we proposed a framework for social group recognition. The framework takes in both individual and global features. Features are extracted from fine-tuned CNN networks which has been pre-trained on ImageNet dataset, and then combined. Our results showed the success of our framework by achieving much better result than the previous work.

We also investigated into the pre-trained CNN features. Both visualization and numeric results showed the generalization ability of pre-trained CNN features to features of people’s social styles, which has little shared features with the ImageNet object categories. Meanwhile, we found that there is a correlation between the probability of an image being in Imagenet classes and social group classes, and that better-recognized categories are more correlated with ImageNet categories.

In the future work, we intend to improve the classification performance by adapting convolutional networks more to social groups datasets. The relationship between ImageNet categories and urban tribes classes also brings forward an interesting future topic.

## References

- [1] I. S. Kwak, A. C. Murillo, P. Belhumeur, S. Belongie, and D. Kriegman, “From bikers to surfers: Visual recognition of urban tribes,” in *British Machine Vision Conference (BMVC)*, (Bristol), September 2013.
- [2] A. C. Murillo, I. S. Kwak, L. Bourdev, D. Kriegman, and S. Belongie, “Urban tribes: Analyzing group photos from a

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Figure 5: selected urban tribe classes and the corresponding highest correlation  $l_{ImageNet}$ . Upper nine images: candidate person images with high  $Pr(l_{ImageNet})$  and high  $Pr(l_{urban})$ . Lower eight images: example of images in  $l_{ImageNet}$ .

(a) Individual accuracy - maximum of  $R$       (b) Individual accuracy - maximum of  $I$       (c) Scene accuracy - maximum of  $R$       (d) Scene accuracy - maximum of  $I$

Figure 6: The relationship between class-wise recognition rate and maximum of correlation  $R(\Pr(l_{urban}), \Pr(l_{ImageNet}))$ , class-wise recognition rate and maximum of mutual information  $I(l_{urban}, l_{ImageNet})$

social perspective,” in *CVPR Workshop on Socially Intelligent Surveillance and Monitoring (SISM)*, (Providence, RI), June 2012.

- [3] M. D. Zeiler and R. Fergus, “Visualizing and understanding convolutional networks.,” *CoRR*, vol. abs/1311.2901, 2013.
- [4] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, “Decaf: A deep convolutional activation feature for generic visual recognition.,” *CoRR*, vol. abs/1310.1531, 2013.
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks.,” in *NIPS* (P. L. Bartlett, F. C. N. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, eds.), pp. 1106–1114, 2012.
- [6] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, “Backpropagation applied to handwritten zip code recognition,” *Neural Computation*, vol. 1, pp. 541–551, 1989.
- [7] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” in *Proceedings of the IEEE*, vol. 86(11), pp. 2278–2324, 1998.

- [8] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. R. Ioffe, J. Shlens, and Z. Wojna, “Going Deeper with Convolutions,” 2014, arXiv:1409.4842.
- [9] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, “Overfeat: Integrated recognition, localization and detection using convolutional networks,” in *International Conference on Learning Representations (ICLR2014)*, 2014.
- [10] P.-H. Gosselin, N. Murray, H. Jégou, and F. Perronnin, “Inria+xerox@fgcomp: Boosting the fisher vector for fine-grained classification,” Tech. Rep. 0, INRIA, December 2013.
- [11] Y. Jia, “Caffe: An open source convolutional architecture for fast feature embedding.” <http://caffe.berkeleyvision.org/>, 2013.
- [12] A. C. Gallagher and T. Chen, “Understanding images of groups of people.,” in *CVPR*, pp. 256–263, IEEE, 2009.
- [13] G. Wang, A. C. Gallagher, J. Luo, and D. A. Forsyth, “Seeing people in social context: Recognizing people and social relations,” 2013.

8

864	tionships.,” in <i>ECCV</i> (K. Daniilidis, P. Maragos, and N. Paragios, eds.), vol. 6215, pp. 169–182, Springer, 2010.	918
865		919
866		920
867	[14] A. Dhall, J. Joshi, I. Radwan, and R. Goecke, “Finding hap-	921
868	piest moments in a social context.,” in <i>ACCV 2</i> (K. M. Lee,	922
869	Y. Matsushita, J. M. Rehg, and Z. Hu, eds.), vol. 7725,	923
870	pp. 613–626, Springer, 2012.	924
871	[15] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-	925
872	J. Lin, “Liblinear: A library for large linear classification,”	926
873	<i>Journal of Machine Learning Research</i> , vol. 9, pp. 1871–	927
874	1874, 2008.	928
875	[16] G. H. Laurens van der Maaten, “Visualizing data using	929
876	t-sne,” <i>Journal of Machine Learning Research</i> , vol. 9,	930
877	pp. 2579–2605, November 2008.	931
878		932
879		933
880		934
881		935
882		936
883		937
884		938
885		939
886		940
887		941
888		942
889		943
890		944
891		945
892		946
893		947
894		948
895		949
896		950
897		951
898		952
899		953
900		954
901		955
902		956
903		957
904		958
905		959
906		960
907		961
908		962
909		963
910		964
911		965
912		966
913		967
914		968
915		969
916		970
917		971