

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 [4]. 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 ([5], [4]). 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 ([6]). 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([6], [7]). 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 [8] 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 [4]. 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 al.*  
121 [9]. Since then, CNNs have shown successful results on var-  
122 ious computer vision tasks, such as hand-written digit clas-  
123 sification ([10],[3]), ImageNet Large Scale Visual Recog-  
124 nition Challenge ([8], [11]), object detection ([12],[11]), ob-  
125 ject localization ([12]), and many modified architecture of  
126 CNNs have been proposed([2][3][1]).

127 Recently, many researchers have shown the utility of  
128 generalization of pre-trained CNN features on large dataset  
129 such as ImageNet. Krizhevsky *et al.* showed excellent gen-  
130 eralization of the pre-trained CNN features[6]. They kept  
131 all the layers of ImageNet-trained model using fixed pre-  
132 trained features except for the last softmax classifier, and  
133 achieved best results on Caltech-101 and Caltech-256. Don-  
134 ahue *et al.* used different layers of pre-trained CNN network  
135 as features and trained simple classifiers such as SVM and  
136 Logistic Regression, and outperformed the state-of the-art  
137 on several vision challenges such as scene recognition and  
138 domain adaptation[7].

139 Kwak *et al.* created an urban tribe dataset consisting of  
140 11 classes[4]. The classes are defined from social group  
141 labels provided by Wikipedia. They selected the eight most  
142 popular categories from their list of subculture, and added  
143 three other classes corresponding to typical social venues  
144 in addition. For each class, images of groups of people  
145 were discovered with different search engines, and a broad  
146 range of images for each class were collected. Kwak *et al.*  
147 also provided a group description and several classification  
148 methods([5], [4]). Group description consists of person de-  
149 scriptors and global group descriptors: Six part of person is  
150 detected, and a set of predefined descriptors are computed  
151 for each part, including ratio of skin pixels, color informa-  
152 tion like RGB histograms, and HoG features; Global de-  
153 scriptors use a both low level and high level descriptors to  
154 describe the context and group properties of the image. Low  
155 lever features include color information, Gist, HoG and  
156 ratio of pixels of person, and high level descriptors includes  
157 proximity of persons, alignment or pose of the group, and  
158 scene layout of individuals. Two options of classification  
159 methods are provided: bag of parts-based classification and  
160 SVM-based classification.  
161

162 Categorizing the social groups of individuals belongs to  
163 fine-grained classification task which recently draws more  
164 interest of the computer vision community. It aims at giving  
165 the fine-grained categories in a certain class. Fine-grained  
166 classification is more difficult than conventional classifica-  
167 tion tasks, because the categories are semantically as well as  
168 visually similar, and are even challenging for humans. The  
169 Fine-grained Challenge 2013 (FGComp) provided the data  
170 in several categories including aircraft, birds, cars, dogs,  
171 and shoes. [13] achieved the best result using classifier  
172 based on fisher vectors. However, CNN based methods us-  
173 ing [14] or [7] gave inferior results, especially when the  
174 bounding box of test data is unknown.

175 There is some research in analyzing social groups of peo-  
176 ple. [15] showed the visual structure of a group helps under-  
177 standing events. [16] showed social relationships modeling  
178 helps people recognition. [17] used both local and global  
179 factors for group level expression analysis.

## 180 3. Methods

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

### 183 3.1. Urban tribes dataset

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

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

### 205 3.2. Classification hierarchy

206 To utilize the properties of urban tribes fully, our fea-  
207 ture vector consists of both elements: individual features  
208 and environmental features. For each feature type, we use  
209 a similar extraction strategy. Individual features and envi-  
210 ronmental features are hierarchically combined to form the

216 final decision function. The network hierarchy is shown in  
 217 Figure 1.  
 218

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

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

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

239 The feature vectors are then fed into a multi-class  
 240  $\text{SVM}_{\text{Person}}$ . We use LIBLINEAR[18] 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  
 246 reliable probability estimate from the noisy yet correlated set  
 247 of probabilities  $\Pr_{ij}$ , a simple but effective average pooling  
 248 is performed to  $\Pr_{ij}$ :

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

254  $\Pr_{\text{People}}(C|H_1, \dots, H_p)$  is the probability estimate of class  
 255  $C$  given the set of people candidate images  $H$ .  
 256

257 On the other hand, the entire environment in the scene  
 258 image, denoted by  $S$ , is directly utilized for probability es-  
 259 timation. The procedure to generate probability estimate of  
 260 class  $C$  given the environment as a whole  $\Pr_{\text{Scene}}(C|S)$  is  
 261 similar with that for  $\Pr_{\text{People}}(C|H)$ . The difference is that  
 262 the input of Convolutional Network is  $227 \times 227$  patches  
 263 extracted from the entire scene image, and the fine-tuned  
 264 Convolutional network:  $\text{CNN}_{\text{Scene}}$  and SVM:  $\text{SVM}_{\text{Scene}}$  are  
 265 trained with the training set of entire scene images. Sev-  
 266 eral different strategies to extract patches from scene images  
 267 and corresponding Convolutional Neural Network architec-  
 268 tures are explained in Section 3.3.  
 269

Therefore, the probability estimate of a class  $C$  given

observation of scene  $S$  is:

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

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

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

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

where

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

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

### 3.3 Convolutional network feature extraction

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

Following [14]'s work, we use the network framework  
 called Caffe. The network architecture is described in [8],  
 which won the ImageNet Large Scale Visual Recognition  
 Challenge 2012. We take the activations from the 6th and  
 7th hidden layer of the convolutional neural network, which  
 are two fully connected layers before the class prediction  
 layer. We also take the activations from 6th or 7th layer  
 alone as comparison. We choose these two layers, because  
 as the layers ascend, the features extracted show increasing  
 invariance and semantic meaning.

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

#### 3.3.1 Pre-processing of the dataset

The urban tribe dataset is a relatively small dataset, and both  
 people candidate crops and scene images are of various res-

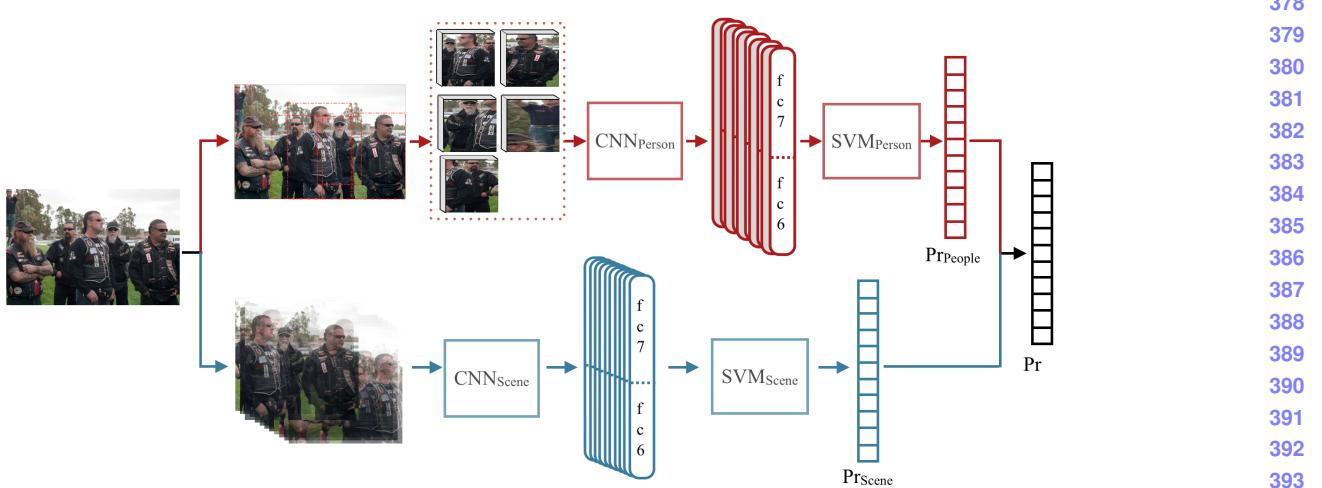


Figure 1: Architecture of classification algorithm using  $\text{Nets}_{SDense}$  which is introduced in Section 3.3.3. The upper half estimates the probability given people candidate images, and the lower half estimates the probability given the entire scene. Dense crop  $\text{CNN}_{Scene}$  and distorted crop  $\text{CNN}_{Person}$  are used, as shown in Section 3.3.1.

olution. Our convolutional neural networks require constant input image size of  $227 \times 227$ , so pre-processing of the dataset is necessary.

There are several strategies to make one image compatible with the CNN:

1. *Distort cropping*: As in [8], 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.

### 3.3.2 Network Fine-tune

Although the pre-trained network from [14] can already generalize well to many datasets, the urban tribe dataset has its unique property. It emphasizes certain visual features such as certain clothing styles, while pays less attention to other visual features. Also, it emphasizes style of the object 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.

### 3.3.3 Choices of network combination

Scene images and individual images have different properties, and need different strategies of pre-processing and sep-

432 arate fine-tuning. For scene network  $\text{CNN}_{\text{Scene}}$  and scene  
 433 images, due to the small size of the dataset, we use the *dense*  
 434 *cropping* technique, the third technique in Section 3.3.1, to  
 435 increase the dataset. For person network  $\text{CNN}_{\text{Person}}$  and  
 436 corresponding input, we use the *distort cropping* technique,  
 437 because the subimages have normally long height and short  
 438 width, and the second and third strategies using squared  
 439 crops of a person image will lose much information, no  
 440 matter which location we choose to crop them; whereas the  
 441 first method ensure each crop keeps the essential features  
 442 for classification.

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

445 We also construct other combination of networks for  
 446 comparison:

- 447 1.  $\text{Nets}_{NoTune}$ : Directly use the pre-trained network by  
 448 [14] for both scenes and persons, and use the *distort*  
 449 *cropping* technique (distorted crops) as input patches  
 450 for both networks. This choice of cropping strategy is  
 451 in consistent with the way the network is pre-trained.
- 452 2.  $\text{Nets}_{SSparse}$ : Use the *sparse cropping* strategy for  
 453 scene features, and the *distort cropping* strategy for  
 454 person features. .
- 455 3.  $\text{Nets}_{SDistort}$ : Use the *distort cropping* strategy for  
 456 both scene features and person features.

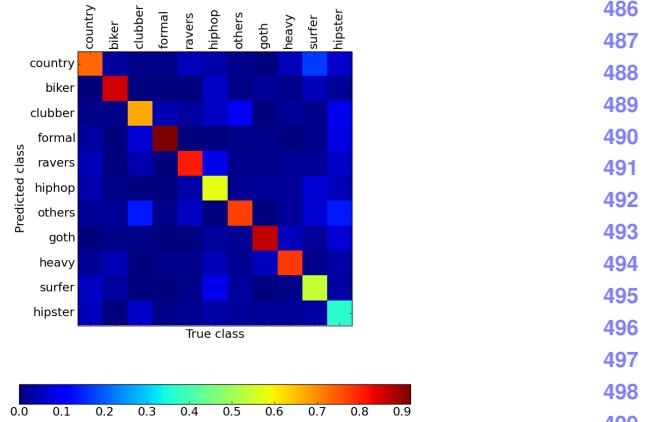
## 4. Experiments and Results

463 In this section, the performance of the proposed classi-  
 464 fication scheme is evaluated and analyzed. In the experi-  
 465 ments, six rounds of 5-fold cross validation are performed,  
 466 therefore we have 30 training experiments in total. Dataset  
 467 is partitioned into 5 equal sized subsets, containing one-fifth  
 468 of the data points from each category. One of the subsets is  
 469 used as test set, and the remained 4 subsets are used as train-  
 470 ing data. In the fine-tune procedure, there are 6000 training  
 471 iterations, and learning rate is decreased by ten times after  
 472 each 1000 iterations.

### 4.1. Urban tribe classification performance

473 Table 1 shows the comparison of performance using  
 474 different approaches. The 30 segmentations of datasets are  
 475 used for all the approaches tested in this section, and 30 test  
 476 results are averaged for each approach. The standard error  
 477 is shown with the accuracy in Table 1. We also compare our  
 478 result with the result achieved by [4] using their best model.  
 479 The advantage of CNN pre-trained features is obvious.

480 The confusion matrix is shown in Figure 2 for  
 481  $\text{Nets}_{SDense}$ , and all the 30 training experiments are aver-  
 482 aged. We can observe there is a obvious difference of dif-  
 483 ficulty of different categories. Class *formal* has accuracy as



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

504 high as about 90%, while class *hipster* is the most difficult  
 505 class, having less than 60% accuracy.

506 Comparing the result of using different features in the  
 507 same approach shows the necessity of every step of our  
 508 architecture. In results using  $\text{Nets}_{SDense}$  with concatenated  
 509 features, average accuracy for each candidate person is low  
 510 as 47.10%. Average pooling of candidate person probability  
 511 estimates produces a large accuracy increase of about 20%.  
 512 Accuracy using the entire scene only results in 67.26% ac-  
 513 curacy. Combining probabilities  $\text{Pr}_{People}(C|H_1, \dots, H_p)$   
 514 and  $\text{Pr}_{Scene}(C|S)$  achieves accuracy as high as 71.22%,  
 515 which verifies the complementary role of people candidate  
 516 feature and environment feature in a group image.

517 We also compare vertically the accuracy of different ap-  
 518 proaches, to show the role of network feature concatena-  
 519 tion. Using only 7th layer or 6th layer activation from  
 520 the networks  $\text{Nets}_{SDense}$  produces decent results, showing  
 521 both layers' activations can generate high semantic fea-  
 522 tures. Concatenating both layers' activations increases the  
 523 accuracy by 0.5%, indicating the slight information loss of the  
 524 7th fully connected layer.

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

532  $\text{Nets}_{SDistort}$ ,  $\text{Nets}_{SSparse}$ , and  $\text{Nets}_{SDense}$  use different  
 533 patch extraction strategies. Note that we use the same dis-  
 534 torted patch extraction method for person images, as men-  
 535 tioned in Section 3.3.3, while we use three different meth-

540  
541  
542 Table 1: Performance of different approaches using different information.  
543  
544  
545  
546  
547  
548  
549  
550

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> [4]	-	-	-	$46(\text{std: } 2)$

551  
552 ods for scene images. The results for scene images show  
553 the advantage of keeping the aspect ratio of scene images,  
554 and the slight advantage of using dense crops. However, the  
555 final results with People+Scene for the three methods don't  
556 have significant differences, this is due to the combination  
557 with people information.

#### 559 4.2. Convolutional Network feature analysis

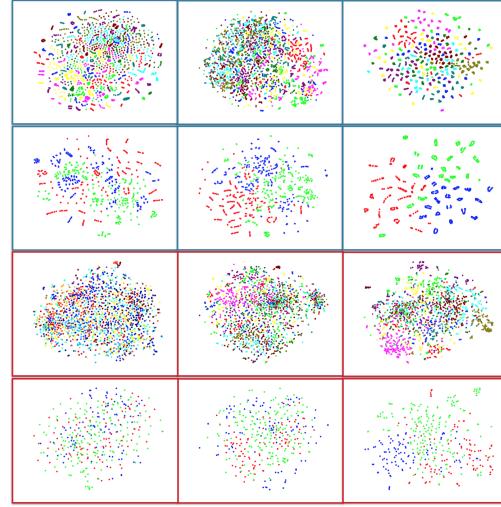
560  
561 **Yufei: Is this section necessary?** We use t-distributed  
562 stochastic neighbor embedding technique [19] to visualize  
563 the power of different layers' features in convolutional  
564 network. Two dimensional embedding of the high dimensional  
565 CNN feature space is extracted, and we plot the data as 2-d  
566 points with different colors indicating different classes they  
567 belong to. Powerful features will pull data of different  
568 semantic classes more apart.

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

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

#### 589 4.3. Urban tribe classes vs. ImageNet classes

590  
591 It's recently being acknowledged that CNN pre-trained  
592 features are generic and can be used for new tasks. In this  
593 section we check the relationship between the new tasks and



594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

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

638 original ImageNet task, using the urban tribe dataset, which  
639 gives us some insight about the features extracted from pre-  
640 trained network.

641 The urban tribe dataset contains groups of people, and  
642 the important features for categorization are mainly human  
643 related features, such as attire, make up, posture and expres-  
644 sions. However, the ILSVRC dataset used for pre-training  
645 contains few images of humans. Instead of examining the  
646 output of the layers directly, we try to find the relationship  
647 between the 1000 classes in ILSVRC dataset and the classes  
648 in two image sets: urban tribe dataset, and candidate person  
649 images extracted from urban tribe dataset.

650 We use the parameters of pre-trained CNN network as  
651 our feature extraction model, and train a softmax layer on  
652 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 ImageNet categories. We use one round of 5-fold cross validation, and use all the images in urban tribe dataset for analysis

We first check the relationship of  $\text{Pr}(l_{\text{urban}})$  and  $\text{Pr}(l_{\text{ImageNet}})$  of candidate person images. We calculate the correlation coefficient of  $\text{Pr}(l_{\text{urban}})$  and  $\text{Pr}(l_{\text{ImageNet}})$  for all 1000 ImageNet classes, denoted as  $R(\text{Pr}(l_{\text{urban}}), \text{Pr}(l_{\text{ImageNet}}))$ . 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  $I(l_{\text{urban}}, l_{\text{ImageNet}})$ .

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  $\text{Pr}(l_{urban})$  and  $\text{Pr}(l_{ImageNet})$ , 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  $l_{ImageNet}$ .

698 The correlation between ImageNet class and urban tribe  
699 class and its relationship with class-wise recognition rate  
700 may indicate that the “generic” features extracted by pre-  
701 trained CNN networks are not so generic. The network is

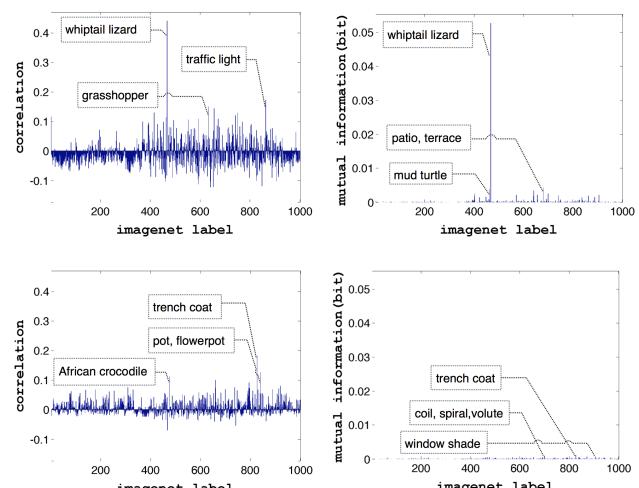


Figure 4: Correlation and mutual information of  $l_{urban}$  and 1000  $l_{ImageNet}$ . The first row is for  $l_{urban}$ =Biker. The second row is for  $l_{urban}$ =Hipster. Top three ImageNet classes are marked.

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.

## 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.

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

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  $l_{ImageNet}$ .

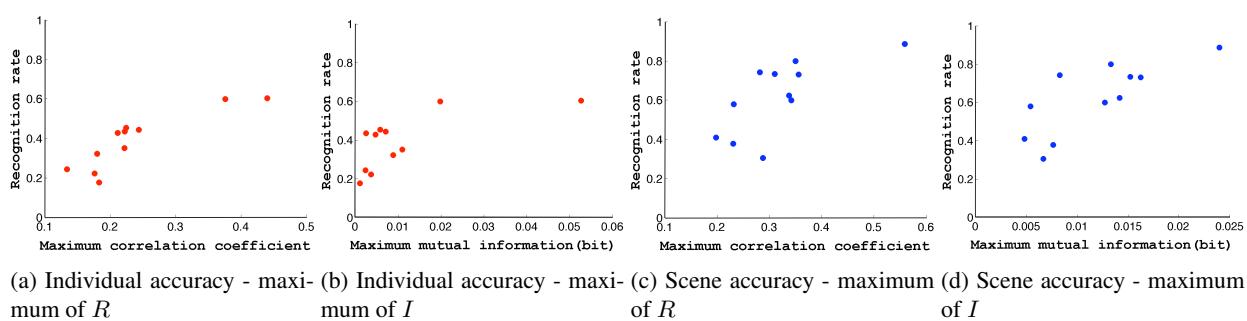


Figure 6: The relationship between class-wise recognition rate and maximum of correlation  $R(\Pr(l_{urban}), \Pr(l_{ImageNet}))$ , class-wise recognition rate and maximum of mutual information  $I(l_{urban}, l_{ImageNet})$

## References

- [1] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, “Network in network,” 2013, arXiv:1312.4400.
- [2] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Improving neural networks by preventing co-adaptation of feature detectors,” *CoRR*, vol. abs/1207.0580, 2012.
- [3] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, “Maxout networks,” 2013, arXiv:1302.4389.
- [4] 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.
- [5] A. C. Murillo, I. S. Kwak, L. Bourdev, D. Kriegman, and S. Belongie, “Urban tribes: Analyzing group photos from a social perspective,” in *CVPR Workshop on Socially Intelligent Surveillance and Monitoring (SISM)*, (Providence, RI), June 2012.
- [6] M. D. Zeiler and R. Fergus, “Visualizing and understanding convolutional networks.,” *CoRR*, vol. abs/1311.2901, 2013.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. R. and Dragomir Anguelov, D. Erhan, V. Vanhoucke, and

864	A. Rabinovich, "Going Deeper with Convolutions," 2014, arXiv:1409.4842.	918
865		919
866		920
867	[12] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated recognition, localiza- tion and detection using convolutional networks," in <i>Inter- national Conference on Learning Representations (ICLR2014)</i> , 2014.	921
868		922
869		923
870		924
871	[13] 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.	925
872		926
873		927
874		928
875	[14] Y. Jia, "Caffe: An open source convolutional architecture for fast feature embedding." <a href="http://caffe.berkeleyvision.org/">http://caffe.berkeleyvision.org/</a> , 2013.	929
876		930
877		931
878	[15] A. C. Gallagher and T. Chen, "Understanding images of groups of people.,," in <i>CVPR</i> , pp. 256–263, IEEE, 2009.	932
879		933
880		934
881	[16] G. Wang, A. C. Gallagher, J. Luo, and D. A. Forsyth, "Seeing people in social context: Recognizing people and social rela- tionships.,," in <i>ECCV</i> (K. Daniilidis, P. Maragos, and N. Par- agios, eds.), vol. 6215, pp. 169–182, Springer, 2010.	935
882		936
883		937
884		938
885	[17] A. Dhall, J. Joshi, I. Radwan, and R. Goecke, "Finding hap- piest moments in a social context.,," in <i>ACCV 2</i> (K. M. Lee, Y. Matsushita, J. M. Rehg, and Z. Hu, eds.), vol. 7725, pp. 613–626, Springer, 2012.	939
886		940
887		941
888		942
889	[18] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.- J. Lin, "Liblinear: A library for large linear classification," <i>Journal of Machine Learning Research</i> , vol. 9, pp. 1871– 1874, 2008.	943
890		944
891		945
892		946
893	[19] G. H. Laurens van der Maaten, "Visualizing data using t-sne," <i>Journal of Machine Learning Research</i> , vol. 9, pp. 2579–2605, November 2008.	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