

Hybrid CNN and Dictionary-Based Models for Scene Recognition and Domain Adaptation

Guo-Sen Xie, Xu-Yao Zhang, Shuicheng Yan, *Senior Member, IEEE*, Cheng-Lin Liu, *Fellow, IEEE*

Abstract—Convolutional neural network (CNN) has achieved state-of-the-art performance in many different visual tasks. Learned from a large-scale training dataset, CNN features are much more discriminative and accurate than the hand-crafted features. Moreover, CNN features are also transferable among different domains. On the other hand, traditional dictionary-based features (such as BoW and SPM) contain much more local discriminative and structural information, which is implicitly embedded in the images. To further improve the performance, in this paper, we propose to combine CNN with dictionary-based models for scene recognition and visual domain adaptation. Specifically, based on the well-tuned CNN models (e.g., AlexNet and VGG Net), two dictionary-based representations are further constructed, namely mid-level local representation (MLR) and convolutional Fisher vector representation (CFV). In MLR, an efficient two-stage clustering method, i.e., weighted spatial and feature space spectral clustering on the parts of a single image followed by clustering all representative parts of all images, is used to generate a class-mixture or a class-specific part dictionary. After that, the part dictionary is used to operate with the multi-scale image inputs for generating mid-level representation. In CFV, a multi-scale and scale-proportional GMM training strategy is utilized to generate Fisher vectors based on the last convolutional layer of CNN. By integrating the complementary information of MLR, CFV and the CNN features of the fully connected layer, the state-of-the-art performance can be achieved on scene recognition and domain adaptation problems. An interested finding is that our proposed hybrid representation (from VGG net trained on ImageNet) is also complementary with GoogLeNet and/or VGG-11 (trained on Place205) greatly.

Index Terms—Convolutional neural networks, Scene recognition, Domain adaptation, Dictionary, Part learning, Fisher vector.

I. INTRODUCTION

WITH the development of deep learning, convolutional neural network (CNN) [1], [2] has been successfully applied to various fields, such as object recognition [2], [3], [4], [5], [6], image detection [7], [8], [9], image segmentation [10], [11], [12], [13], image retrieval [14], and so on. State-of-the-art performance achieved by CNN in these fields identify the powerful feature representation ability of CNN for different visual tasks.

G.-S. Xie, X.-Y. Zhang and C.-L. Liu are with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, No. 95 Zhongguancun East Road, Beijing 100190, P.R. China. C.-L. Liu is also with the Research Center for Brain-Inspired Intelligence, Institute of Automation, and the CAS Center for Excellence in Brain Science and Intelligence Technology, Chinese Academy of Sciences. E-mail: {guosen.xie, xyz, liucl}@nlpr.ia.ac.cn.

S. Yan is with the Department of Electrical and Computer Engineering, National University of Singapore, 117583, Singapore. Email: eleyan.s@nus.edu.sg.

The traditional Bag of Visual Words (BoW) models [15], [16], [17], [18], [19], [20], which had been quite popular before 2012 in the research community of object recognition, is now being gradually neglected since the CNN model gained the champion in the large-scale competition of ImageNet classification [2] in 2012 (ILSVRC-12). Due to the strong representation ability of CNN trained on a large dataset, e.g., ImageNet, some works [21], [22] advocate directly using CNN features for classification, and have gained much better performance than traditional methods. There also exist some works [23], [24] that have tried to combine traditional models with the CNN model to get better results. Traditional features and CNN features can complement each other, and better performance can be obtained by combining them.

The images from the ImageNet database for training CNN are mostly object-oriented, thus it seems that the trained CNN model is more suitable for object recognition than for other tasks, e.g., scene recognition [25]. Zhou et al. [26] collected a large-scale place database to train the CNN (PlaceNet), with the same architecture as [22], and found that based on the features of PlaceNet, better performance can be obtained than the original CNN features [22]. Recently, Mircea et al. [27] found that stronger CNN architectures can approach and outperform PlaceNet even if trained on ImageNet data. Yosinski et al. [28] verified that CNN features in the lower layer are more “general” while the features of higher layers are more “specific”. Therefore, how to explore the underlined representation and discriminative ability of CNN, no matter what kind of database it is trained on, is still an interesting problem. In this paper, we focus on scene recognition [25] and visual domain adaptation [29]. Unlike in generic object categorization, scene images have many discriminative part regions, which is beneficial for distinguishing the categories. During the past few years, many works [30], [31], [32], [33], [34], [35] have been devoted to discover discriminative parts for scene recognition. The above methods first train a detector or classifier and then detect the part regions, which are not scalable for large databases. To discover discriminative parts more efficiently, there have also been some methods that are based on bottom cues to discover discriminative proposals, e.g., [36], [37], [38], [39]. Visual domain adaptation (DA) can be classified into two types, unsupervised and semi-supervised adaptation, according to whether the target domain samples are used or not for training the classifier. Pioneer works of DA include [29], [40], [41], [42], [43], [44], which are based on learning with regularization on the manifold. Recent works on DA [45], [46], [47], [48] have also adopted the strategy of adding maximum mean discrepancy (MMD) constraint on the

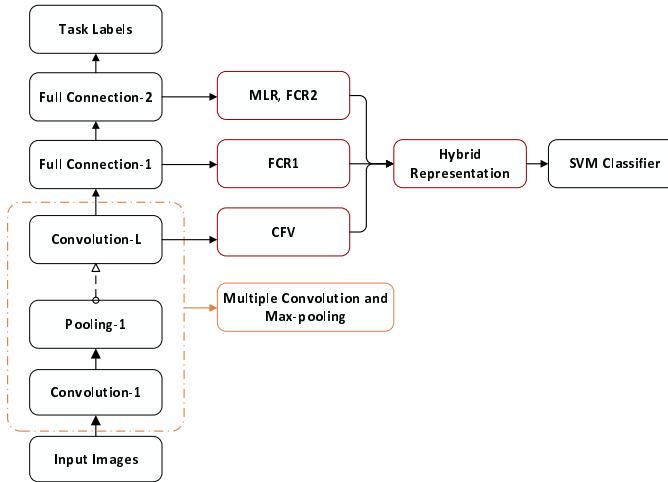


Fig. 1. The proposed hybrid representation based on the trained CNN. CFV is calculated based on the last convolutional layer (Convolution-L); FCR1 indicates the global representation based on the first fully connected layer (Full Connection-1); similarly, FCR2 is calculated based on the second fully connected layer (Full Connection-2); and our proposed MLR is based on the second fully connected layer too. Operations in the dashed box are multiple convolution, max-pooling and other normalization, determined by the used CNN architecture.

fully connected layer of neural networks during the training process.

For both scene recognition and domain adaptation, very few works have utilized local discriminative information implicitly contained in the images in the context of CNN. Especially for domain adaptation, no one knows whether the local part information can improve the domain transfer performance or not. Herein, we propose to utilize the local part information to improve the discriminative ability of CNN features. Specifically, we propose to combine CNN features with two dictionary-based models. The first one is the mid-level local discriminative representation (MLR). MLR aims to discover local discriminative information contained in the images, which are beneficial for afterward classification. To construct MLR, we utilize selective search [36] to generate initial parts for each image, followed by our proposed two-stage clustering to filter out redundancy parts and generate part dictionary, finally the locality-constrained linear coding (LLC) [17] and spatial pyramid matching (SPM) [49] are used for generating the MLR. (See Fig. 2 for illustration.) On the other hand, Fisher vectors of the last convolutional layer of CNN (CFV) are generated to further boost the performance. As we know, Fisher vector [18] consists of first order and second order differences between the descriptors and the GMM centers, which gives CFV the same representation ability to well distinguish different categories. As for GMM training before Fisher vector coding, multi-scale and scale-proportional descriptors are sampled. As for Fisher vector coding, we use the same strategy as [50], which is denoted as Multi-scale Pyramid Pooling. Another commonly used feature of CNN is the fully connected layer representation (FCR), which contains the global information of input images. By combining these several representations, i.e., mid-level local representation (MLR), convolutional Fisher vector representation (CFV), and the global representations of

the last two fully connected layers of CNN (FCR) together, we can obtain our hybrid representation (See Fig. 1 for the whole flowchart). The advantages of our hybrid representation lie in 1) having more discriminative ability for classification, and 2) being complementary to each other. Experimental results on both scene recognition and domain adaptation validate the strong domain transfer ability from other large database of our hybrid representation. Moreover, we also find that MLR is both complementary with CFV and FCR.

The remainder of this paper is organized as follows. In Section II, we give some related works. In Section III, we illustrate the proposed pipeline elaborately. Section IV shows that the hybrid representation can achieve the best results in scene recognition and visual domain adaptation. In Section V, we conclude this paper and discuss the future work.

II. RELATED WORK

In this section, we introduce the related works to our hybrid representation. Gong et al. [24] proposed to calculate the VLAD [51] representation based on the FCR of the trained CNN, of which the performance is better than the original FCR and other traditional methods. Let $X = [x_1, x_2, \dots, x_N] \in \mathbb{R}^{d \times N}$ be the single-scale activates of the image I based on the FCR of the trained CNN, and $C = [c_1, c_2, \dots, c_M] \in \mathbb{R}^{d \times M}$ be the codebook (dictionary) learned by the K-means algorithm based on all the training images. Then, the VLAD coding of the image I can be represented as

$$v = \left[\sum_{NN(x_i)=c_1} x_i - c_1, \dots, \sum_{NN(x_i)=c_M} x_i - c_M \right], \quad (1)$$

where $NN(x_i) = c_j$ denotes the set of x_i whose nearest neighbor is c_j . To improve the performance, Gong et al. [24] calculated multi-scale representations of Eq. 1 and concatenated all the scales. Motivated by [24], the CFV [18] can be further used to improve the performance of image representation [27], [50].

Suppose we are given the multi-scale activates $\{X^{(s)} | X^{(s)} \in \mathbb{R}^{d \times N^{(s)}}\}$, $X^{(s)} = [x_1^{(s)}, x_2^{(s)}, \dots, x_{N^{(s)}}^{(s)}]$, $s = 1, 2, \dots, S$ from all the training images. Let $u_\lambda = \sum_{t=1}^M \omega_t u_t(x)$ denote a Gaussian Mixture Model (GMM), where $\lambda = \{\lambda_i = \{\omega_i, \mu_i, \sigma_i\}, i = 1, 2, \dots, M\}$ represents the parameters of the GMM, and λ can be optimized by the Maximum Likelihood (ML) estimation based on $\{X^{(s)}, s = 1, 2, \dots, S\}$. Denote the multi-scale activates of the image I as $\chi = \{\chi_s = [x_1^{(s)}, x_2^{(s)}, \dots, x_{I^{(s)}}^{(s)}], s = 1, 2, \dots, S\}$. The gradients of the $u_\lambda(\chi)$ w.r.t. the i -th Gaussian can be represented as follows

$$g_i = \frac{1}{|\chi|} \sum_{s=1}^S \sum_{n=1}^{|\chi_s|} \nabla_{\lambda_i} \log u_\lambda(x_n^{(s)}). \quad (2)$$

The Fisher vector of the image I is obtained by concatenating all the gradients w.r.t. those M Gaussians. In [50], while calculating gradients w.r.t. each Gaussian, Yoo et al. adopted Multi-scale Pyramid Pooling, where first the GMM parameters are generated based on all the descriptors from different scaled

training images, and then scale-specific normalization and max-pooling are implemented.

Recently, Liu et al. [52] also proposed to learn a part dictionary based on the LC-KSVD method [53] directly. In this work, each part element is given a label which is the same as that of the image where it is located on. LC-KSVD is a popular dictionary learning method, but it is very time-consuming and therefore cannot address large-scale problems, e.g., the part dictionary learning of the SUN database [54]. Another drawback is that many local parts usually have no explicit semantic (label) information at all, so unsupervised part dictionary learning may be more reasonable than the supervised counterpart.

III. THE PROPOSED HYBRID METHOD

In this section, we illustrate the whole pipeline of the construction of MLR, CFV and FCR, followed by the details of each part.

A. The Whole System

To explore the representation ability of deep CNN features, we propose to combine the three kinds of representations based on CNN features, i.e., our proposed MLR, CFV, and FCR. We illustrate the whole pipeline in Fig. 1. In the figure, given the trained CNN models with (without) fine-tuning, e.g. AlexNet [2] and VGG Net [3], we further extract the multi-scale CFVs, single-scale FCRs and our local discriminative MLRs for each image. Then the concatenated hybrid representations are fed into the linear SVM classifier [55]. As validated by our experiments, the hybrid representations are very powerful features due to their complementary components, i.e., CFV which contains one or two order gradient information, FCR which contains global information of images, and MLR which contains local discriminative (structural) information. The whole procedure of our hybrid representation is illustrated in Algorithm 1.

B. MLR: Integrating Local Discrimination and CNN Features

CNN has become a strong tool to learn invariant features for various visual applications. Nevertheless, during the training process of CNN, the input images are all global ones, and no good strategy has been found to incorporate the local information into the training process of CNN yet. In this paper, we propose a strategy to generate the mid-level local representation (MLR), which is based on part dictionary clustering and multi-scale mid-level representation generating. MLR is constructed based on the local discriminative part dictionary, which leads to its local discrimination. In this subsection, we illustrate the details of our method for constructing the mid-level representation (MLR) by utilizing the local discriminative ability of images.

At the beginning, we utilize selective search [36] to get the discriminative part proposals for every image, with the constraints of the pixel number of proposals within $[60 \times 60, 160 \times 160]$, and width/height (height/width) smaller than 3, which is used to catch the local information of the images.

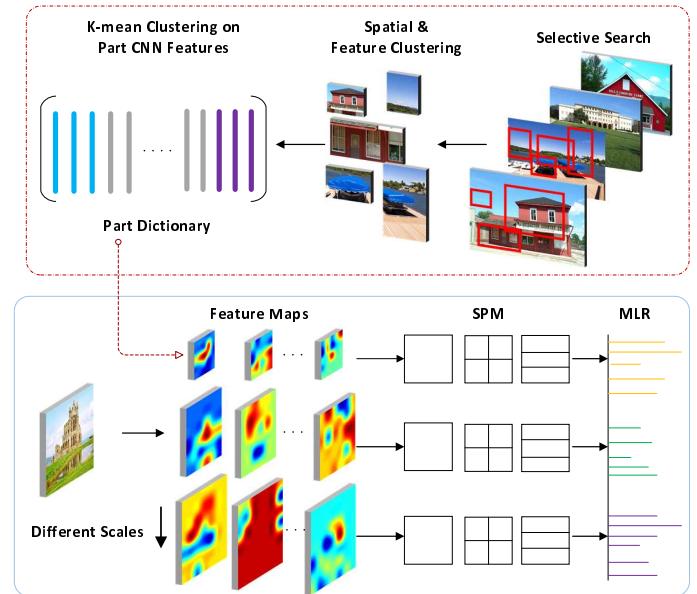


Fig. 2. The process of calculating MLR based on the trained part dictionary. The operations in the red dashed box are the part dictionary generating process based on two-stage clustering, and operations in the below are the MLR generating process which is based on part dictionary and SPM. Best viewed in color.

Furthermore, the spectral clustering in both the spatial and the feature space is conducted with the bounding box information of each image. Specifically, suppose the selected bounding boxes (by selective search) are $B = [B_1, B_2, \dots, B_{n_I}] \in \mathbb{R}^{4 \times n_I}$ with each B_i denoting the coordinates of top-left and bottom-right on the image I , and the corresponding last fully connected layer activates of the trained CNN for B are denoted as $F = [f_1, f_2, \dots, f_{n_I}] \in \mathbb{R}^{d \times n_I}$ ($d = 4096$ for current popular networks). Herein, we use the activates after ReLU [2]. Under our constraints, while extracting B , the number of n_I is usually less than 500, so the forward propagation for extracting F will be acceptable. We construct the final similarity graph $G = (V, E)$ based on both B and F , and the weights on the edges are as follows:

$$W = \lambda_B W_B + \lambda_F W_F, \quad (3)$$

where λ_B and λ_F are the weighting parameters and $\lambda_B + \lambda_F = 1$ is used in our paper. Specifically, the elements of W_B and W_F are denoted as

$$W_B(i, j) = \frac{|B_i| \cap |B_j|}{|B_i| \cup |B_j|}, \quad (4)$$

where $|\cdot|$ indicates calculating the area of the input box.

$$W_F(i, j) = \exp\left(-\frac{\|f_i - f_j\|_2^2}{2\sigma^2}\right), \quad (5)$$

for $i, j = 1, 2, \dots, n_I$.

After spectral clustering on the graph G , we obtain Q clusters $\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_Q$. Then we sort the Q clusters in descending order according to the number of bounding boxes and select the top T ($T \leq Q$) clusters. We randomly select one bounding box from each $\mathcal{C}_i, i = 1, 2, \dots, T$, thus totally T bounding boxes are reserved. We further implement context padding [7]

on the selected T bounding boxes and do forward propagation (based on CNN) again to get the new feature representation based on the context padded boxes. In this way, we can get T representations $P = [P_1, P_2, \dots, P_T]$ for each image I , and we denote them as the prototypes of the corresponding image.

Finally, for all the prototypes of the training images, we do K-means clustering [56] to get the class-specific part dictionary and the class-mixture part dictionary, denoted as $D_{cs} \in \mathbb{R}^{d \times K}$ and $D_{cm} \in \mathbb{R}^{d \times K}$ respectively. Here the elements of D_{cs} are obtained by first clustering the prototypes from each class and then concatenating them together. While D_{cm} is based on clustering all the prototypes without considering the class information. The part dictionary contains local discriminative and multi-scale information of the training images, due to the selective search algorithm (SSA) and constraints while running SSA.

With the part dictionary D_{cs} (D_{cm}) learned, we can consider it as a group of local discriminative filter banks. Motivated by Object-Bank [57], given an input image at a single scale, we sample square regions at multiple scales densely, i.e., sampled squares with size 128×128 , 92×92 and 64×64 and step size 32 pixels for all three scales. Then we calculate the activates of the last fully connected layer of the trained CNN for all these squares under different scales, which generates three scaled activate tensors. This can also be seen as the local feature extraction on square regions. After that, for each scaled activate tensor of each image, we use the D_{cs} (D_{cm}) to operate with the activates on each location of the tensors, resulting in K new feature maps for each scaled activate tensor, which can be seen as another kind of convolutional operation based on the part dictionary [58]. For every K feature maps under different scales, we further apply spatial pyramid matching (SPM) [49] to divide the map region into spatial cells in three levels, i.e., 1×1 , 2×2 , 3×1 , and do max-pooling on each cell. The final MLR is the concatenation of all the max-pooled features, with the dimension being $3 \times K \times (1 \times 1 + 2 \times 2 + 3 \times 1)$.

Many methods can be adopted as the strategy of operation between the part dictionary D_{cs} (D_{cm}) and the activates of local square regions, e.g., inner production, sparse coding [16], locality-constrained linear coding (LLC) [17], and auto-encoder based coding [59]. There are detailed coding speed comparisons in [59]. Here we utilize the LLC for feature coding, due to its relative fast speed and locality-preserve property.

Specifically, denote the activating tensor of the image I under one scale by $A \in \mathbb{R}^{d \times h \times h}$, and $A_{(\cdot, i, j)} \in \mathbb{R}^d$, ($i, j = 1, 2, \dots, h$) is the activating vector on the location (i, j) of the $h \times h$ input tensor. To obtain the values $v^* \in \mathbb{R}^K$ on the location (i, j) of K new maps, we only need to solve the following LLC problem:

$$v^* = \arg \min_v \|A_{(\cdot, i, j)} - Dv\|_2^2 + \lambda \|dist \odot v\|_2^2, \quad (6)$$

where D can be taken as D_{cs} or D_{cm} , \odot denotes the element-wise multiplication, and $dist = [\exp(\|v - d_1\|^2/\tau), \exp(\|v - d_2\|^2/\tau), \dots, \exp(\|v - d_K\|^2/\tau)] \in \mathbb{R}^K$ is the adaptive vector between v and each dictionary element, which preserves the locality between v and the dictionary D . In this paper, we

Algorithm 1 Extracting of Hybrid Representation

Input: Training (source) images: $\{I_i\}_{i=1}^m$ and test (target) images: $\{I_i\}_{i=m+1}^n$. A CNN model with or without fine-tuning.
Output: Hybrid Representation $\{H_i\}_{i=1}^n$ for each training and test images.
Procedure:

- 1: **for** $i = 1 \rightarrow n$ **do**
- 2: **Selective search** on image I_i , obtain part set $B = [B_1, B_2, \dots, B_{n_i}] \in \mathbb{R}^{4 \times n_i}$.
- 3: **Graph construction** based on Eq. 3, 4, 5, obtain graph G with weight as W .
- 4: **Spectral clustering** on G , obtain Q clusters.
- 5: **Sort** the Q clusters by the number of boxes contained in each cluster.
- 6: **Select** the top T clusters, and randomly take one box from each of the T cluster.
- 7: **K-means clustering** on the selected parts, generate class-specific or class-mixture dictionary.
- 8: **LLC coding** (Eq. 6) and **SPM**, generate the **MLR_i**
- 9: **Fisher vector coding** (Eq. 2) on the last convolutional layer of CNN, obtain **CFV_i**.
- 10: **Fully connected representations**, **FCR1_i** and **FCR2_i**
- 11: $H_i = [\text{MLR}_i, \text{CFV}_i, \text{FCR1}_i, \text{FCR2}_i]$
- 12: **end for**

utilize the approximated LLC [17] for the fast calculation of v^* .

The flowchart of generating the MLR given the part dictionary D_{cs} or D_{cm} is shown in Fig. 2.

C. CFV: Convolutional Fisher Vector

In this part, we briefly describe the construction of the Fisher vector [18] based on the last convolutional layer of CNN. Like [18], [24], [27], and [50], we also use multiple scales as input while constructing CFVs.

Moreover, considering the different scale information, we also calculate CFVs for each scale followed by $L2$ normalization and max-pooling. This strategy is known as Multi-scale Pyramid Pooling [50]. The difference between our method and [50] is that for sampling descriptors for GMM [60] training before Fisher vector coding, we also consider scale information (i.e., the number of sampled descriptors is proportional to that of the total descriptors under different scales in each image). Here the descriptors from the last convolutional layer are without ReLU [2] throughout our paper.

To improve the performance, power and $L2$ normalization are further applied to the max-pooled representation, which generates the final CFV representation.

D. FCR: CNN Features from Well-Tuned Networks

Given the well-trained CNN model (with or without fine-tuning, such as AlexNet and VGG Net), we can use it for other different visual tasks, which is termed as parameter transfer learning. The fully connected representation (FCR) of the resized input images are extracted by forward propagating them until the fully connected layers of the network. We denote FCR1 and FCR2 (Fig. 1) as the penultimate and the last fully connected layer representation with ReLU [2], respectively. Under this definition, our MLR is constructed based on FCR2

and local discriminative part (spectral clustering) on each image.

Before being fed into SVM classifier for training, L_2 normalization is also applied to the FCRs.

IV. EXPERIMENTS

In this section, we show the classification accuracy of our hybrid representation in two applications, i.e., scene recognition and visual domain adaptation, compared with state-of-the-art models, including traditional and CNN based ones. We first introduce the datasets and experimental settings, then report experimental results on each dataset, after that we give an analysis of the key parameters followed by the analysis of the complementary ability of the proposed hybrid representation with the representations from other Nets, and finally we analysis the time complexity of calculating our representation.

A. Datasets and Experimental Settings

For scene recognition, we utilize three trained CNN models, i.e., AlexNet [2], VGG-19 net [3], GoogLeNet [4], and VGG-11 [76] net based on their caffe implementations [61]. Specifically, AlexNet and VGG-19 net are trained on ImageNet database, and these two nets are used to evaluate our hybrid representation system. GoogLeNet and VGG-11 net are trained on Place205 database, which are used to validate the complementary ability of our hybrid representation w.r.t. the representations from other nets. For domain adaptation, we use the AlexNet only, so that we can compare with other CNN based methods fairly. In our experiments, we resize all the images into resolution of 256×256 before the following operations, such as selective search and Fisher vector extraction. After obtaining the hybrid representation, we use linear SVM to train the classifier in all our experiments. The databases and experimental details are as follows:

1) **MIT Indoor-67 database**: MIT Indoor-67 database [25] is a popular indoor scene database, including 15,620 images of 67 indoor scenes. It is difficult to distinguish different classes, because the categories are all indoor scenes, and inter-class variance between different classes is very little. We follow the same training-test partition as [25] by using approximately 80 images from each class for training and 20 for testing. The average class accuracy is reported in our paper. Sample images of Indoor-67 can be found in Fig. 3 (a). It can be seen from Fig. 3 (a) that the categories of movie theater, meeting room and classroom are very similar.

2) **SUN-397 database**: SUN-397 [54] is a large-scale scene recognition database, containing 397 categories and with scenes varying from abbey to zoo, including both indoor and outdoor scenes. At least 100 images are contained in each category. We use the publicly available training-test partitions, and report the average accuracy and standard errors for all the partitions. We use 50 images from each class for training and 50 for testing. Sample images can be found in Fig. 3 (b).

3) **Office database**: Office dataset [29] contains three sub-datasets from different domains, i.e., Amazon, Webcam, and Dslr. Each sub-dataset is from a separate domain: images from Amazon are collected from online catalogs of amazon.com,

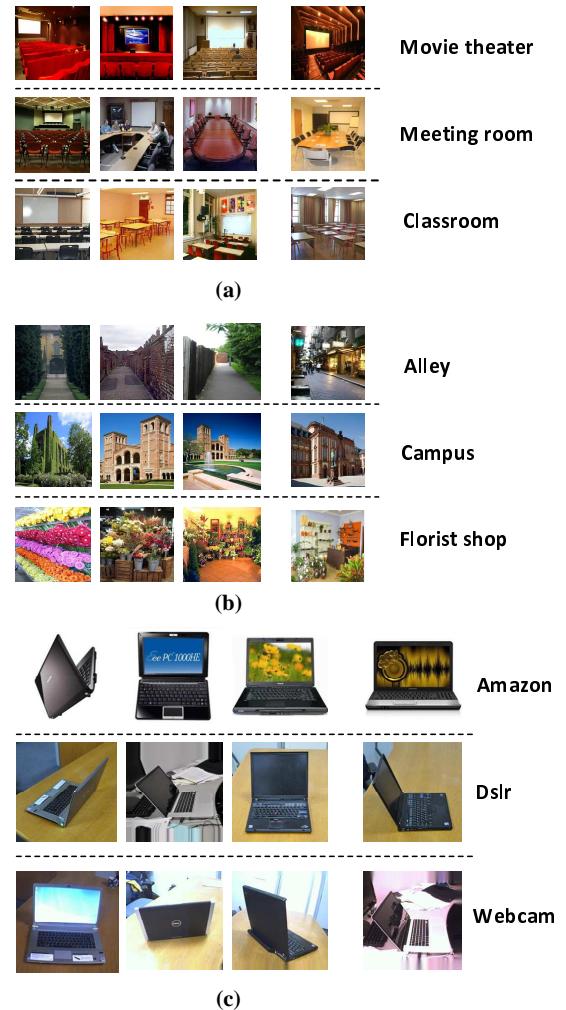


Fig. 3. The sample images from different databases. (a) Images from MIT Indoor-67 database. It can be seen that categories “movie theater”, “meeting room” and “classroom” are very similar and difficult to distinguish. (b) Images from SUN-397 database. It includes outdoor scenes, e.g. “alley” and “campus”, and indoor scenes, e.g. “florist shop”. (c) Images from “destop” category in Office database. It contains 3 sub-sets, i.e., Amazon, Dslr, and Webcam. Best viewed in color.

and images in Dslr and Webcam are obtained in the daily office environment by a digital SLR camera and a webcam with high and low resolutions, respectively. There are totally 31 categories which are common for these three domains, and the number of images per category per domain ranges from 8 to 100, and totally 4,652 images are included. Fig. 3 (c) shows the domain shift between different domains.

4) **Details for domain adaptation**: In domain adaptation, if the training data (source domain) with labels and the test data (target domain) without labels are given, it is called unsupervised domain adaptation; if a source domain with labels and a target domain with a small amount of labeled data are given, then the problem is denoted as semi-supervised domain adaptation. We will perform both unsupervised and semi-supervised domain adaptation on all the 31 categories in our experiments. We adopt the standard experimental setup presented in [29]. We use 20 source examples per category when Amazon is taken as the source domain, and 8 images

per category when Webcam or Dslr is taken as the source domain [29]. For the semi-supervised adaptation, three more labeled target examples per category are added into the source domain. We also try another setting presented in [42], [46], where we use all the source domain examples with labels for unsupervised adaptation and three more target domain samples per class for semi-supervised adaptation.

We evaluate our method across five random training-test partitions for each of the three domain transfer tasks commonly used for evaluation, i.e., Amazon \rightarrow Webcam ($A \rightarrow W$), Dslr \rightarrow Webcam ($D \rightarrow W$) and Webcam \rightarrow Dslr ($W \rightarrow D$), and report average accuracy and standard errors for each setting.

Actually, given the CNN model, calculating the hybrid representation based on Algorithm 1 is already a domain (parameter) transfer process. As for domain adaptation experiments, we run Algorithm 1 to generate hybrid representation for both source and target data, which is used for further domain transfer by training classifier on the source representation, and testing on the target one. Note that in domain adaptation, after obtaining the prototype boxes based on the first stage spectral clustering, we first carry out box filtering based on variance thresholding of the prototype boxes, followed by the second stage K-means clustering. The purpose of this trick is to filter out the box regions with low variance, which may be the surrounding background regions in the Office database. Specifically, given the prototype box I_b , we calculate the variance of its gray scale image as follows:

$$VAR = var(gray(I_b)). \quad (7)$$

Then if $VAR < 125$, we delete the prototype. Fig. 4 illustrates the deleted prototype boxes based on our strategy.

5) **Parameter settings:** In Algorithm 1, there are several key parameters, we give a detailed description of these parameters. As for graph constructing, the parameter (λ_B, λ_F) in Eq. 3 is set as $(1, 0)$ for Indoor-67 database, and $(0.5, 0.5)$ for SUN397 and Office databases¹. σ in Eq. 5 is set as 1. The number G of clusters for spectral clustering is set as 10, and the number of top ranking T clusters is set as 5, which is equal to the number of boxes selected from each image. As for Fisher vector coding, the Gaussian components of GMM training is fixed as 64, and we use multiple scales in our paper. Specifically, given the CNN input size of $L \times L$ for the image I , the used five scales are $L \times \sqrt{2}^{[0 \ 1 \ 2 \ 3 \ 4]} = [L \ \sqrt{2}L \ 2L \ 2\sqrt{2}L \ 4L]$ ². As for FCR1 and FCR2, we extract features with the whole image as input for Indoor-67 and SUN-397, and extract features with the central cropped sub-region as input for Office. The SVM parameter C is fixed as 1 in all our experiments.

6) **Abbreviation in the experiments:** To make some key definition clear, we list these abbreviation in Table I.

¹ Without considering the feature space information on Indoor-67 ($\lambda_F = 0$) is based on the observations: the layout in Images of Indoor-67 is almost dense, while the layouts in images of SUN-397 and Office are sparse, i.e., some sub-regions are the same, e.g. the sky in the campus category of SUN-397, and the background of ruler category in Office dataset.

²Note that for Indoor-67 experiment under VGG-19 net, we use the 10 scales the same as [27], i.e., $L \times \sqrt{2}^{[-6 \ -5 \ -4 \ -3 \ -2 \ -1 \ 0 \ 1 \ 2 \ 3]}$ to reproduce their results.

TABLE I
ABBREVIATION IN THIS PAPER.

Abbreviation	Description
FCR1-c	FCR1 based on central crop as input
FCR1-w	FCR1 based on whole image as input
FCR2-c	FCR2 based on central crop as input
FCR2-w	FCR2 based on whole image as input
CM	Construct MLR based on class-mixture dictionary
CS	Construct MLR based on class-specific dictionary
G_P205	The average pooled representation before loss3 of GoogLeNet (trained on Place205)
VGG-11_P205	FCR1-w of VGG-11 trained on Place205
VGG-19_Hybrid	FCR1-w+CFV+MLR based on VGG-19 trained on ImageNet

TABLE II
COMPARISON OF CLASSIFICATION RATE ON MIT INDOOR-67 DATABASE.

Traditional Methods	Accuracy (%)	
ROI [25]	26.05	
MM-Scene [62]	28.00	
DPM [35]	30.40	
CENTRIST [63]	36.90	
Object Bank [57]	37.60	
RBOB [64]	37.93	
Discriminative Patches [31]	38.10	
Hybrid parts [65]	39.80	
LPR-LIN [66]	44.84	
BOP [30]	46.10	
MI-SVM [67]	46.40	
Hybrid parts+GIST-color+SP [65]	47.20	
ISPR [33]	50.10	
MMDL [68]	50.15	
Discriminative Parts [69]	51.40	
DSFL [23]	52.24	
Discriminative List Group [70]	55.58	
IFV [30]	60.77	
IFV+BOP [30]	63.10	
Mode-Seeking [32]	64.03	
Mode-Seeking + IFV [32]	66.87	
ISPR+IFV [33]	68.50	
CNN based Methods	Accuracy (%)	
FCR2 (placeNet) [26]	68.24	
Order-less Pooling [24]	68.90	
CNNaug-SVM [21]	69.00	
FCR2 HybridNet [26]	70.80	
URDL+CNNaug [52]	71.90	
MPP-FCR2 (7 scales) [50]	75.67	
DSFL+CNN [23]	76.23	
MPP+DSFL [50]	80.78	
CFV (VGG-19) [27]	81.00	
Our hybrid methods (CS)	Accuracy (%)	
	AlexNet	VGG-19
FCR1-c	59.71	72.96
FCR1-w	61.63	71.48
FCR2-c	59.48	70.23
FCR2-w	61.01	68.71
FCR2-w (fine-tuning)	64.78	—
MLR	69.33	77.51
CFV	68.68	81.05
FCR1-w+FCR2-w	62.43	70.59
MLR+FCR2-w	70.20	77.94
CFV+FCR2-w	71.27	80.75
MLR+FCR1-w	70.66	78.81
CFV+FCR1-w	70.94	81.60
MLR+CFV	72.80	81.47
MLR+CFV+FCR1-w	74.09	82.24
MLR+CFV+FCR2-w	73.17	81.57
CFV+FCR1-w+FCR2-w	70.72	79.70
MLR+FCR1-w+FCR2-w	69.57	78.34
MLR+CFV+FCR1-w+FCR2-w	72.98	80.91

B. MIT Indoor-67 Experiments

In this Subsection, we report and analyze the experimental results on the MIT Indoor-67 database. We first show the changing tendency of the classification rate under different part dictionary sizes, in the class-mixture and the class-specific manner respectively, without fine-tuning the CNN. The part

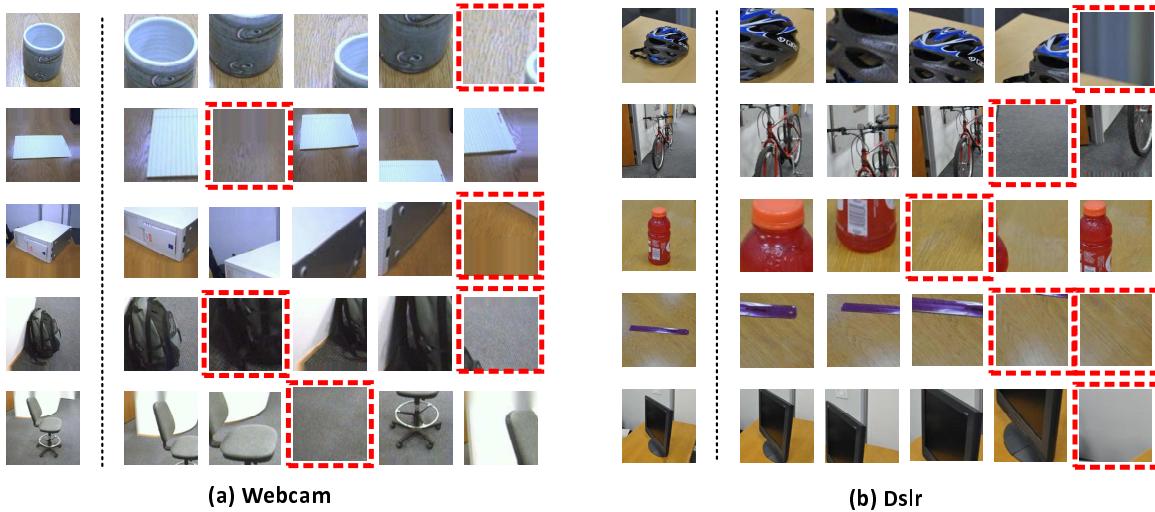


Fig. 4. Examples of original images (first column in both (a) Webcam and (b) Dslr domains) and their corresponding prototype boxes (2nd-6th columns) based on our proposed first stage spectral clustering. The boxes are resized to fit CNN input. Red dashed boxes are removed ones based on Eq. 7. Best viewed in color.

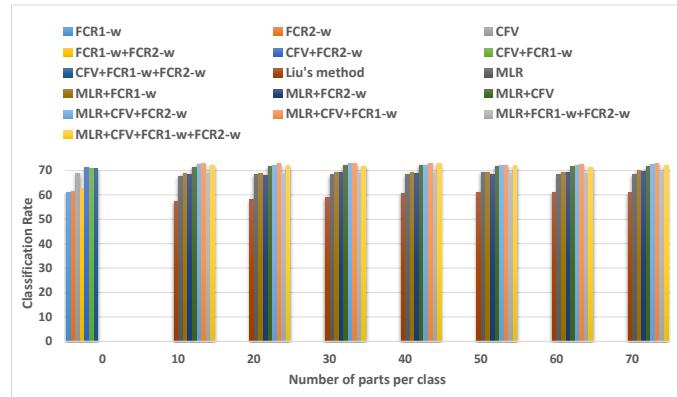


Fig. 5. Changing tendency of classification rate with different part element numbers of class-mixture dictionary D_{cm} on Indoor-67 database, based on AlexNet. Best viewed in color.

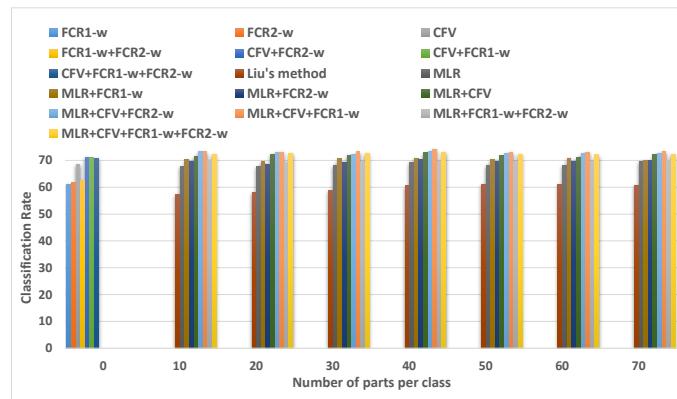


Fig. 6. Changing tendency of classification rate with different part element numbers of class-specific dictionary D_{cs} on Indoor-67 database, based on AlexNet. Best viewed in color.

dictionary size can be seen as the number of representative parts per category multiplying the total category number (see Fig. 5-7 for details).

In both Fig. 5 and Fig. 6, “0” on the x-axis means that representations FCR1-w, FCR2-w, CFV, FCR1-w+FCR2-w,

CFV+FCR2-w, CFV+FCR1-w and CFV+FCR1-w+FCR2-w are not generated based on the part dictionary. The detailed meaning of FCR1-w and FCR2-w are listed in Table I. From Fig. 5 and 6, it can be seen that our part dictionary is much better than that of Liu et al. [52] in the classification rate under different part dictionary sizes, and the hybrid representations combined with MLR are much better than their counterpart representations. From Fig. 7, a conclusion can be drawn that class-specific dictionaries are always better than class-mixture ones. Moreover, learning a class-specific dictionary is much faster than learning a class-mixture one based on K-means, thus the best choice of the second stage clustering is the class-specific part dictionary.

To better evaluate the representation ability of our proposed MLR, we also compare the classification rate with or without fine-tuning based on AlexNet [2]. We use the AlexNet with the architecture of Caffe’s implementation [61], which contains five convolutional layers, two fully connected layers and one output layer with a node number equal to the number of categories (i.e., the output number of nodes is 67 for Indoor-67 database). We first fine-tune AlexNet with all the global training images (resized to 256×256) of Indoor-67. The initialization of the parameters of the first seven layers is the same as the model trained on ImageNet database, and the parameters of the last layer are randomly initialized with Gaussian distribution. The learning rates of the first seven layers and the last fully-connected layer are initialized as 0.001 and 0.01 respectively, and reduced to one tenth of the current rates after fixed iterations (10,000 in our experiments). By setting a larger learning rate for the last layer, due to the random initialization of this layer, we hope the parameters of this layer can be updated to be more suitable for the new task-specific database. For the previous layers of AlexNet, we hope the parameters change as little as possible to preserve the already learned texture and shape information during the learning based on ImageNet database. Then based on our fine-tuned task-specific model, we further fine-tune the AlexNet based on the prototype bounding boxes (with labels the same

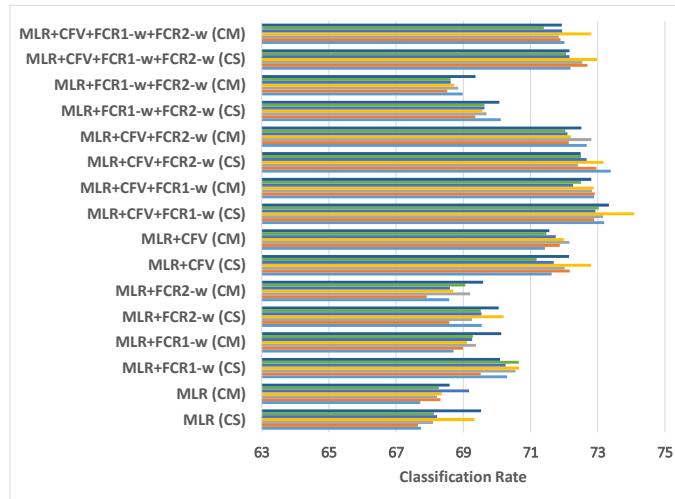


Fig. 7. Comparison of changing tendency of classification rate under different part element numbers and different part dictionaries, i.e., Class-Mixture (CM) and Class-Specific (CS) on Indoor-67 database, based on AlexNet. Take MLR(CS) as an example. The accuracies from bottom to top are corresponding to the part element number that varies from 10×67 to 70×67 . Best viewed in color.

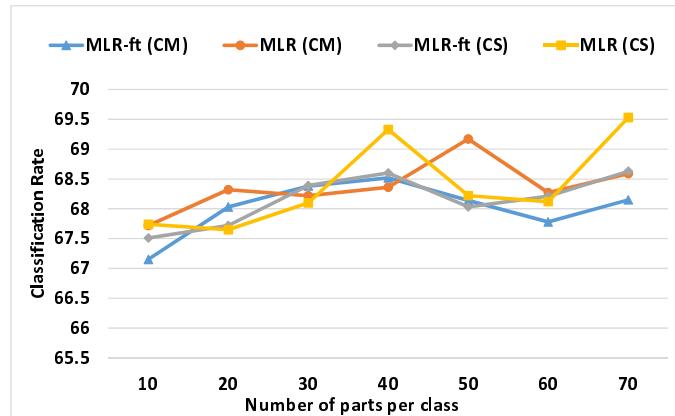


Fig. 8. Comparison of classification rate (with or without fine-tuning) under different part element numbers and different dictionaries on Indoor-67 database, based on AlexNet. Best viewed in color.

as the images where they are located on), which are extracted based on the first stage spatial and feature spectral clustering.

An observation is that the classification rate of MLR based on fine-tuned AlexNet is comparable with their counterpart representation calculated based on AlexNet without fine-tuning. See Fig. 8 for details. Therefore, we do not fine-tune our used CNN models in the later experiments.

Finally, we list the classification accuracies on Indoor-67 and compare our hybrid methods with other state-of-the-art methods, both under AlexNet and VGG-19 networks. Table II gives the detailed accuracies of the traditional models, the CNN based models, and our hybrid representation based methods. Here the dictionary size for constructing MLR is 67×40 and 67×30 for AlexNet and VGG-19, respectively. It can be seen that our hybrid representation achieves the result of 82.24% based on VGG-19 network, which has been the best result on Indoor-67 based on the net trained on ImageNet database. This justifies the effectiveness of combining CNN and dictionary-based features in improving the classification

TABLE III
COMPARISONS OF CLASSIFICATION RATE ON SUN-397 DATABASE.

Traditional Methods	Accuracy (%)
S-manifold [71]	28.90
OTC [72]	34.56
contextBow+semantic [73]	35.60
Xiao et al. [54]	38.00
FV (SIFT + Local Color Statistic) [74]	47.20
OTC + HOG2x2 [72]	49.60
CNN based Methods	Accuracy (%)
Decaf [22]	40.94
MOP-CNN [24]	51.98
Hybrid-CNN [26]	53.86±0.21
Place-CNN [26]	54.32±0.14
Place-CNN ft [26]	56.2
Our hybrid methods (CS)	Accuracy (%)
	AlexNet
FCR2-w	46.27±0.37
FCR1-w	46.42±0.47
MLR	53.84±0.16
CFV	52.30±0.09
FCR1-w+FCR2-w	47.19±0.53
MLR+FCR2-w	55.22±0.34
CFV+FCR2-w	54.65±0.40
MLR+FCR1-w	55.17±0.33
CFV+FCR1-w	54.34±0.22
MLR+CFV	56.66±0.17
MLR+CFV+FCR1-w	57.15±0.26
MLR+CFV+FCR2-w	57.31±0.12
CFV+FCR1-w+FCR2-w	53.80±0.53
MLR+FCR1-w+FCR2-w	54.84±0.32
MLR+CFV+FCR1-w+FCR2-w	56.83±0.13
Our hybrid methods (CS)	VGG-19
	Accuracy (%)

accuracy.

C. SUN-397 Experiments

In this subsection, we report and analyze the classification results on the SUN-397 database. It can be seen from Table III that our hybrid representation can also achieve the best results under both AlexNet and VGG-19 net configurations based on the net trained on ImageNet database. In our experiments on SUN-397 database, to improve the computation efficiency, we fix the part number of each class as 10 while generating the class-mixture part dictionary D_{cm} or the class-specific one D_{cs} (i.e., the total number of elements for D_{cm} and D_{cs} is both 3790).

Note that only under AlexNet, our best result is 57.31%, which is already better than the current best result of 56.20% by fine-tuning the PlaceNet trained on the large place database [26]. On the other hand, we fix the part dictionary size as 3970 in our experiments on this database. If we enlarge the dictionary size, better results can be further obtained. On this database, the performance of the class-specific dictionary is also better than the class-mixture counterpart (see Fig. 9 for details).

D. Domain Adaptation Experiments

In this part, we report and analyze the domain adaptation (DA) experiments, under both unsupervised setting and semi-supervised setting with two different source-target partitions. The datasets used here are the Amazon, Webcam and Dslr datasets. We only conduct experiments based on AlexNet in DA experiments, which almost all yield best results under the four settings. If our representation is constructed under

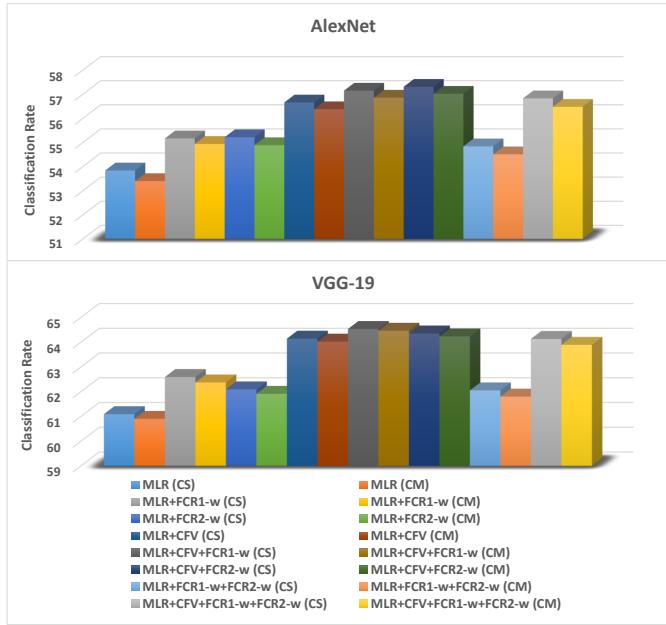


Fig. 9. Comparison of classification rate under different part dictionaries (dictionary size is fixed as 3970), i.e., Class-Mixture (CM) and Class-Specific (CS) on SUN-397 database, based on AlexNet and VGG-19 Net. Best viewed in color.

more strong CNN models, better improvements can be further obtained.

Table IV-VII show the DA performances of our method compared with other methods. It can be seen that our results are much better than other state-of-the-art DA approaches under $D \rightarrow W$ and $W \rightarrow D$. For the domain transfer $A \rightarrow W$, only comparative results are obtained. Here, the reason may be that the domain biases between domains A and W are very large, which can greatly reduce the importance of local information. However, we can construct the hybrid representation based on the trained models of DDC [45] and DAN [46] which are better than our results based on AlexNet. Another strategy to enhance the representation ability of our representation is to use stronger CNN models, e.g., VGG-19 Net, to construct the hybrid representation. Surprisingly, we get 100% domain transfer accuracy for transferring from Webcam to Dslr, under the setting of using all Webcam images and three more Dslr samples per class for training. This indicates that relatively small domain shifts can be totally avoided by our local representation (MLR) combined with global representations (FCR1-c+FCR2-c) of CNN features.

E. Parameter Analysis

In this part, we take Indoor-67 as an example to evaluate the influence of the key parameters (λ_B, λ_F) (in Eq. 3) w.r.t. the final accuracy of the generated MLR, under AlexNet. Specifically, we evaluate eight combinations which are listed in the first column of Table VIII. The results of MLR in Table VI-II are obtained under the class-mixture dictionary with size of $10 \times \text{Category number}$, and CFV and FCRs are obtained under AlexNet, the same as the original experimental settings. It can be seen from Table VII that the best accuracy of MLR is achieved under $(\lambda_B, \lambda_F) = (1, 0)$, while the best accura-

TABLE IV
CLASSIFICATION ACCURACY OF 31-CATEGORY OFFICE DATASET UNDER UNSUPERVISED ADAPTATION SETTINGS. 20 SAMPLES PER CLASS ARE USED WITH AMAZON AS THE SOURCE DOMAIN, AND 8 SAMPLES PER CLASS WITH WEBCAM AND DSLR AS THE SOURCE DOMAIN RESPECTIVELY.

Methods	$A \rightarrow W$	$D \rightarrow W$	$W \rightarrow D$
GFK (PLS,PCA) [41]	15.0 ± 0.4	44.6 ± 0.3	49.7 ± 0.5
SA [43]	15.3	50.1	56.9
DaNBNN [75]	23.3 ± 2.7	67.2 ± 1.9	67.4 ± 3.0
DaNN [47]	35.0 ± 0.2	70.5 ± 0.0	74.3 ± 0.0
Dlid [44]	26.1	68.9	84.9
Decaf [22]	52.20 ± 1.70	91.50 ± 1.50	—
DDC [45]	59.40 ± 0.80	92.50 ± 0.30	91.70 ± 0.80
FCR2-c	54.21 ± 1.59	91.55 ± 0.64	90.44 ± 0.66
FCR1-c	53.46 ± 0.94	91.35 ± 0.61	91.24 ± 1.25
MLR	40.43 ± 3.03	89.79 ± 1.32	89.68 ± 1.21
CFV	21.06 ± 2.72	76.50 ± 1.32	80.08 ± 1.13
FCR1-c+FCR2-c	55.17 ± 1.91	92.53 ± 0.78	92.21 ± 1.06
MLR+FCR2-c	52.28 ± 2.61	93.36 ± 1.58	93.49 ± 0.96
CFV+FCR2-c	49.36 ± 1.92	91.80 ± 0.91	92.77 ± 0.97
MLR+FCR1-c	49.46 ± 3.03	92.53 ± 1.41	92.89 ± 1.17
CFV+FCR1-c	45.74 ± 1.51	91.27 ± 0.78	92.05 ± 1.18
MLR+CFV	36.15 ± 3.72	90.79 ± 1.00	90.20 ± 1.22
MLR+CFV+FCR1-c	45.21 ± 3.43	92.98 ± 1.19	93.29 ± 1.45
MLR+CFV+FCR2-c	49.33 ± 3.53	93.18 ± 1.16	93.57 ± 1.40
CFV+FCR1-c+FCR2-c	52.25 ± 1.98	92.83 ± 1.02	93.41 ± 1.27
MLR+FCR1-c+FCR2-c	54.29 ± 2.96	93.66 ± 1.51	93.98 ± 0.95
MLR+CFV+FCR1-c+FCR2-c	51.45 ± 3.33	93.46 ± 1.26	94.22 ± 1.36

TABLE V
CLASSIFICATION ACCURACY OF 31-CATEGORY OFFICE DATASET UNDER SEMI-SUPERVISED ADAPTATION SETTINGS. 20 SAMPLES PER CLASS ARE USED WITH AMAZON AS SOURCE DOMAIN, AND 8 SAMPLES PER CLASS WITH WEBCAM AND DSLR AS THE SOURCE DOMAIN RESPECTIVELY.

Methods	$A \rightarrow W$	$D \rightarrow W$	$W \rightarrow D$
GFK (PLS,PCA) [41]	46.4 ± 0.5	61.3 ± 0.4	66.3 ± 0.4
SA [43]	45.0	64.8	69.9
DaNBNN [75]	52.8 ± 3.7	76.6 ± 1.7	76.2 ± 2.5
DaNN [47]	53.6 ± 0.2	71.2 ± 0.0	83.5 ± 0.0
Dlid [44]	51.9	78.2	89.9
Decaf S+T [22]	80.7 ± 2.3	94.8 ± 1.2	—
DDC [45]	84.1 ± 0.6	95.4 ± 0.4	96.3 ± 0.3
DAN [46]	85.6 ± 0.3	95.8 ± 0.2	96.7 ± 0.2
FCR2-c	80.68 ± 1.50	94.81 ± 0.93	94.02 ± 1.41
FCR1-c	80.06 ± 1.23	95.01 ± 0.57	95.21 ± 0.67
MLR	79.32 ± 1.09	94.25 ± 0.80	94.27 ± 1.15
CFV	73.05 ± 1.86	87.52 ± 1.42	88.59 ± 1.88
FCR1-c+FCR2-c	81.40 ± 1.63	95.27 ± 0.87	95.51 ± 1.08
MLR+FCR2-c	83.96 ± 1.66	95.95 ± 0.65	96.84 ± 0.44
CFV+FCR2-c	83.93 ± 1.94	95.21 ± 1.07	95.85 ± 0.75
MLR+FCR1-c	82.74 ± 1.37	95.78 ± 0.54	96.74 ± 0.64
CFV+FCR1-c	82.56 ± 1.79	94.62 ± 0.97	96.05 ± 0.82
MLR+CFV	81.42 ± 1.22	94.79 ± 0.87	95.06 ± 1.34
MLR+CFV+FCR1-c	84.16 ± 0.96	95.61 ± 0.78	96.44 ± 0.85
MLR+CFV+FCR2-c	84.79 ± 1.62	95.98 ± 0.50	96.89 ± 0.37
CFV+FCR1-c+FCR2-c	83.50 ± 2.06	95.24 ± 0.64	96.54 ± 0.84
MLR+FCR1-c+FCR2-c	83.59 ± 1.45	96.15 ± 0.48	96.94 ± 0.64
MLR+CFV+FCR1-c+FCR2-c	84.64 ± 1.57	95.95 ± 0.54	97.14 ± 0.37

cy (72.94%) of hybrid representation (MLR+CFV+FCR1-w) is achieved under $(\lambda_B, \lambda_F) = (0.7, 0.3)$, which is only slightly better than the result (72.89%) under $(\lambda_B, \lambda_F) = (1, 0)$. It can be concluded that (1) MLR is slightly sensitive to the values of (λ_B, λ_F) , the best result is obtained under $\lambda_F = 0$ which further validates the conclusion that feature space information is not necessary during clustering proposals of each image for dense layout database of Indoor-67. (2) The final best hybrid representation (MLR+CFV+FCR1-w) is actually robust to the values of (λ_B, λ_F) , which is very important since we cannot try too many combinations of (λ_B, λ_F) in the real applications. Therefore, the hybrid representation of MLR+CFV+FCRs (therein, we set $(\lambda_B, \lambda_F) = (1, 0)$,

TABLE VI

CLASSIFICATION ACCURACY OF 31-CATEGORY OFFICE DATASET UNDER UNSUPERVISED ADAPTATION SETTINGS. ALL SOURCE DOMAIN DATA ARE USED FOR TRAINING.

Methods	$A \rightarrow W$	$D \rightarrow W$	$W \rightarrow D$
GFK [46]	21.4±0.2	69.1±0.3	65.0±0.2
CNN [46]	59.40±0.5	94.40±0.3	98.80±0.2
LapCNN [46]	60.40 ± 0.30	94.70±0.50	99.10±0.20
DDC [46]	60.50 ± 0.70	94.80±0.50	98.50±0.40
DAN [46]	64.50 ± 0.40	95.20±0.30	98.60±0.20
FCR2-c	61.38	94.72	98.80
FCR1-c	59.50	95.72	99.40
MLR	50.31	93.08	97.79
CFV	25.91	82.52	92.17
FCR1-c+FCR2-c	63.02	95.85	99.60
MLR+FCR2-c	61.51	96.10	99.60
CFV+FCR2-c	55.22	96.35	99.40
MLR+FCR1-c	57.74	95.85	99.20
CFV+FCR1-c	51.57	95.47	99.60
MLR+CFV	44.78	94.72	97.99
MLR+CFV+FCR1-c	53.33	96.48	99.00
MLR+CFV+FCR2-c	56.60	96.48	99.40
CFV+FCR1-c+FCR2-c	58.87	96.73	99.60
MLR+FCR1-c+FCR2-c	63.52	96.73	99.80
MLR+CFV+FCR1-c+FCR2-c	59.50	96.48	99.80

TABLE VII

CLASSIFICATION ACCURACY OF 31-CATEGORY OFFICE DATASET UNDER SEMI-UNSUPERVISED ADAPTATION SETTINGS. ALL SOURCE DOMAIN DATA AND THREE MORE TARGET SAMPLES PER CLASS ARE ADDED FOR TRAINING.

Methods	$A \rightarrow W$	$D \rightarrow W$	$W \rightarrow D$
FCR2-c	76.27±0.59	96.18±0.48	98.86±0.37
FCR1-c	77.75±0.87	96.87±0.66	99.36±0.22
MLR	77.69±1.39	96.21±0.72	99.01±0.60
CFV	69.86±1.58	91.99±1.38	94.72±1.15
FCR1-c+FCR2-c	80.26±1.42	97.12±0.31	99.56±0.11
MLR+FCR2-c	83.65±1.53	97.89±0.31	99.80±0.21
CFV+FCR2-c	83.13±1.19	97.61±0.56	99.65±0.22
MLR+FCR1-c	82.28±1.73	97.64±0.42	99.95±0.11
CFV+FCR1-c	81.40±1.90	97.38±0.52	99.80±0.21
MLR+CFV	80.34±1.95	97.04±0.86	98.86±0.81
MLR+CFV+FCR1-c	83.56±1.82	97.81±0.57	99.95±0.11
MLR+CFV+FCR2-c	84.30±1.74	98.12±0.53	99.90±0.22
CFV+FCR1-c+FCR2-c	83.33±1.78	97.75±0.23	99.95±0.11
MLR+FCR1-c+FCR2-c	83.65±1.63	98.06±0.31	100.0±0.00
MLR+CFV+FCR1-c+FCR2-c	85.24±1.98	98.09±0.42	100.0±0.00

which makes MLR have the best performance) should be still an effective final features to train classifiers and conduct predictions.

F. Complementarity with Other Nets

In this part, we combine our proposed hybrid representation (from VGG-19) with the representations (from GoogLeNet [4] and/or VGG-11 [76] trained on Place205 database[26]) to verify their complementarity to each other. For GoogLeNet, given the input resized image of size 224×224 , the last convolutional layer is first extracted and followed by average pooling to output a 1024 dimensional vector, which is further L_2 normalized. For VGG-11, the first fully connected layer representations are extracted and L_2 normalized. The results on both Indoor-67 and SUN-397 database are listed in Table IX. It can be concluded that our proposed hybrid representation can greatly enhance the performance when combined with GoogLeNet and/or VGG-11 representations. As far as we know, the performances of 85.97% and 70.69% are the best results on Indoor-67 and SUN-397 databases

TABLE IX

COMPLEMENTARY CLASSIFICATION RATE OF VGG-19, VGG-11 AND GOOGLENET.

Methods	Indoor-67 (%)	SUN-397 (%)
G_P205	76.84	63.39 ± 0.22
VGG-11_P205	82.74	66.71 ± 0.08
VGG-19_Hybrid	82.24	64.53 ± 0.24
VGG-11_P205+G_P205	83.42	68.01 ± 0.24
VGG-19_Hybrid+VGG-11_P205	85.59	69.55 ± 0.13
VGG-19_Hybrid+G_P205	84.91	69.48 ± 0.17
VGG-19_Hybrid+VGG-11_P205+G_P205	85.97	70.69±0.15

respectively, which are much better than current state-of-the-art performance.

G. Time Complexity Analysis

As our representations consist of three parts, i.e., MLR, CFV and FCR, the total time consumptions will be the sum of constructing MLR, CFV and FCR. Specifically, time consumption of MLR is mainly due to the two stage clustering (T_1) and and approximated LLC coding [17] for generating different feature maps (T_2). And the time consumption of CFV is the GMM training (T_3) and Fisher vector coding (T_4). Finally, time consumption of FCR is the time of forward propagation of CNN (T_5). So the total time consumption for constructing our representation is $T = \sum_{i=1}^5 T_i$. After obtaining the representations of all the training and testing images, we train linear SVM classifier. The time consumption for training linear SVM and testing should also be considered for the total time complexity. In our experiments, T is approximately two hours for SUN-397 database and less than one hour for Indoor-67 and Office databases.

V. CONCLUSION AND FUTURE WORKS

In this paper, we propose a hybrid representation method for scene recognition and domain adaptation by integrating the powerful CNN features with the traditional well-studied dictionary-based features. An efficient two-stage part dictionary generating method is used to exploit local discriminative and structural information. Based on the generated class-specific and class-mixture part dictionaries, we can get a novel mid-level local representation (denoted as MLR) containing rich local discriminative information, which is not considered during the CNN training process. Moreover, the Fisher vector representation of convolutional layer (CFV) is further used to boost the accuracy. The global representation of the fully connected layer (denoted as FCR) of CNN is validated to be a powerful representation. By combining the complementary information in MLR, CFV, and FCR, a hybrid feature representation can be generated which is much more accurate than the traditional CNN features. Experiments on scene recognition and domain adaptation demonstrate the excellent performance of our hybrid models. In future, to further boost the accuracy, we will jointly train the CNN model and MLR (CFV) under a unified framework. Another interesting future work is to construct our hybrid representations based on the networks trained on place database.

TABLE VIII

THE CLASSIFICATION RATE OF MLR AND ITS COMBINATIONS WITH OTHER REPRESENTATIONS WITH CLASS-MIXTURE DICTIONARY SIZE OF $10 \times (\text{NUMBER OF CATEGORIES})$, UNDER ALEXNET ON INDOOR-67.

Methods / (λ_B, λ_F)	(1,0)	(0.9,0.1)	(0.8,0.2)	(0.7,0.3)	(0.6,0.4)	(0.5,0.5)	(0.4,0.6)	(0.3,0.7)	(0.2,0.8)	(0.1,0.9)	(0,1)
MLR	67.72	66.72	66.54	66.65	65.85	65.64	66.67	65.27	66.18	66.66	66.55
MLR+FCR2-w	68.58	68.45	68.38	68.43	67.61	68.53	68.20	68.07	67.07	68.42	68.41
MLR+FCR1-w	68.71	68.87	68.93	68.56	68.70	68.79	68.78	67.99	68.13	68.26	68.84
MLR+CFV	71.43	71.13	71.09	71.01	71.18	71.25	70.77	71.23	70.62	70.41	70.46
MLR+CFV+FCR1-w	72.89	72.56	72.48	72.94	71.81	72.77	72.10	72.91	72.24	71.87	72.04
MLR+CFV+FCR2-w	72.67	72.25	71.80	72.21	71.52	71.96	71.36	71.51	71.90	72.46	71.86
MLR+FCR1-w+FCR2-w	68.98	68.62	68.17	67.70	68.09	68.90	68.14	67.77	68.18	68.25	68.47
MLR+CFV+FCR1-w+FCR2-w	72.01	71.91	71.64	71.72	71.38	71.91	71.45	71.17	71.21	71.67	71.23

ACKNOWLEDGMENT

This work was supported by the National Basic Research Program of China (973 Program) Grant 2012CB316302, the Strategic Priority Research Program of the CAS (Grant X-DA06040102) and National Natural Science Foundation of China (NSFC) Grant 61403380.

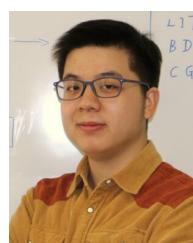
REFERENCES

- [1] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998. [1](#)
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012, pp. 1097–1105. [1, 3, 4, 5, 7](#)
- [3] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015. [1, 3, 5](#)
- [4] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015. [1, 5, 10](#)
- [5] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015. [1](#)
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," *arXiv:1502.01852*, 2015. [1](#)
- [7] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR*, 2014, pp. 580–587. [1, 3](#)
- [8] P. Sermanet, K. Kavukcuoglu, S. Chintala, and Y. LeCun, "Pedestrian detection with unsupervised multi-stage feature learning," in *CVPR*, 2013, pp. 3626–3633. [1](#)
- [9] W. Ouyang, P. Luo, X. Zeng, S. Qiu, Y. Tian, H. Li, S. Yang, Z. Wang, L. Chen-change, and X. Tang, "Deepid-net: multi-stage and deformable deep convolutional neural networks for object detection," in *CVPR*, 2015. [1](#)
- [10] K. Kang and X. Wang, "Fully convolutional neural networks for crowd segmentation," *arXiv:1411.4464*, 2014. [1](#)
- [11] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *CVPR*, 2015. [1](#)
- [12] P. O. Pinheiro and R. Collobert, "From image-level to pixel-level labeling with convolutional networks," in *CVPR*, 2015. [1](#)
- [13] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected crfs," in *ICLR*, 2015. [1](#)
- [14] F. Zhao, Y. Huang, L. Wang, and T. Tan, "Deep semantic ranking based hashing for multi-label image retrieval," in *CVPR*, 2015. [1](#)
- [15] F. Jurie and B. Triggs, "Creating efficient codebooks for visual recognition," in *ICCV*, 2005, pp. 604–610. [1](#)
- [16] J. Yang, K. Yu, Y. Gong, and T. Huang, "Linear spatial pyramid matching using sparse coding for image classification," in *CVPR*, 2009, pp. 1794–1801. [1, 4](#)
- [17] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong, "Locality-constrained linear coding for image classification," in *CVPR*, 2010, pp. 3360–3367. [1, 2, 4, 10](#)
- [18] F. Perronnin and C. Dance, "Fisher kernels on visual vocabularies for image categorization," in *CVPR*, 2007, pp. 1–8. [1, 2, 4](#)
- [19] F. Perronnin, J. Sánchez, and T. Mensink, "Improving the fisher kernel for large-scale image classification," in *ECCV*, 2010, pp. 143–156. [1](#)
- [20] X. Zhou, K. Yu, T. Zhang, and T. S. Huang, "Image classification using super-vector coding of local image descriptors," in *ECCV*, 2010, pp. 141–154. [1](#)
- [21] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "Cnn features off-the-shelf: an astounding baseline for recognition," in *CVPRW*, 2014, pp. 512–519. [1, 6](#)
- [22] 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," *arXiv:1310.1531*, 2013. [1, 8, 9](#)
- [23] Z. Zuo, G. Wang, B. Shuai, L. Zhao, Q. Yang, and X. Jiang, "Learning discriminative and shareable features for scene classification," in *ECCV*, 2014, pp. 552–568. [1, 6](#)
- [24] Y. Gong, L. Wang, R. Guo, and S. Lazebnik, "Multi-scale orderless pooling of deep convolutional activation features," in *ECCV*, 2014, pp. 392–407. [1, 2, 4, 6, 8](#)
- [25] A. Quattoni and A. Torralba, "Recognizing indoor scenes," in *CVPR*, 2009, pp. 413–420. [1, 5, 6](#)
- [26] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, "Learning deep features for scene recognition using places database," in *NIPS*, 2014, pp. 487–495. [1, 6, 8, 10](#)
- [27] M. Cimpoi, S. Maji, and A. Vedaldi, "Deep convolutional filter banks for texture recognition and segmentation," in *CVPR*, 2015. [1, 2, 4, 6](#)
- [28] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" in *NIPS*, 2014, pp. 3320–3328. [1](#)
- [29] K. Saenko, B. Kulis, M. Fritz, and T. Darrell, "Adapting visual category models to new domains," in *ECCV*, 2010, pp. 213–226. [1, 5, 6](#)
- [30] M. Juneja, A. Vedaldi, C. Jawahar, and A. Zisserman, "Blocks that shout: Distinctive parts for scene classification," in *CVPR*, 2013, pp. 923–930. [1, 6](#)
- [31] S. Singh, A. Gupta, and A. A. Efros, "Unsupervised discovery of mid-level discriminative patches," in *ECCV*, 2012, pp. 73–86. [1, 6](#)
- [32] C. Doersch, A. Gupta, and A. A. Efros, "Mid-level visual element discovery as discriminative mode seeking," in *NIPS*, 2013, pp. 494–502. [1, 6](#)
- [33] D. Lin, C. Lu, R. Liao, and J. Jia, "Learning important spatial pooling regions for scene classification," in *CVPR*, 2014, pp. 3726–3733. [1, 6](#)
- [34] C. Doersch, S. Singh, A. Gupta, J. Sivic, and A. Efros, "What makes paris look like paris?" *ACM Transactions on Graphics*, vol. 31, no. 4, 2012. [1](#)
- [35] M. Pandey and S. Lazebnik, "Scene recognition and weakly supervised object localization with deformable part-based models," in *ICCV*, 2011, pp. 1307–1314. [1, 6](#)
- [36] J. R. Uijlings, K. E. van de Sande, T. Gevers, and A. W. Smeulders, "Selective search for object recognition," *IJCV*, vol. 104, no. 2, pp. 154–171, 2013. [1, 2, 3](#)
- [37] M.-M. Cheng, Z. Zhang, W.-Y. Lin, and P. Torr, "Bing: Binarized normed gradients for objectness estimation at 300fps," in *CVPR*, 2014, pp. 3286–3293. [1](#)
- [38] P. Arbelaez, J. Pont-Tuset, J. Barron, F. Marques, and J. Malik, "Multiscale combinatorial grouping," in *CVPR*, 2014, pp. 328–335. [1](#)
- [39] C. L. Zitnick and P. Dollár, "Edge boxes: Locating object proposals from edges," in *ECCV*, 2014, pp. 391–405. [1](#)
- [40] R. Gopalan, R. Li, and R. Chellappa, "Domain adaptation for object recognition: An unsupervised approach," in *ICCV*, 2011, pp. 999–1006. [1](#)
- [41] B. Gong, Y. Shi, F. Sha, and K. Grauman, "Geodesic flow kernel for unsupervised domain adaptation," in *CVPR*, 2012, pp. 2066–2073. [1, 9](#)
- [42] B. Gong, K. Grauman, and F. Sha, "Connecting the dots with landmarks: Discriminatively learning domain-invariant features for unsupervised domain adaptation," in *ICML*, 2013, pp. 222–230. [1, 6](#)
- [43] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars, "Unsupervised visual domain adaptation using subspace alignment," in *ICCV*, 2013, pp. 2960–2967. [1, 9](#)

- [44] S. Chopra, S. Balakrishnan, and R. Gopalan, "Dlid: Deep learning for domain adaptation by interpolating between domains," in *ICML workshop*, vol. 2, 2013, p. 5. 1, 9
- [45] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell, "Deep domain confusion: Maximizing for domain invariance," *arXiv:1412.3474*, 2014. 1, 9
- [46] L. Mingsheng, Y. Cao, J. Wang, and M. Jordan, "Learning transferable features with deep adaptation networks," 2015. 1, 6, 9, 10
- [47] M. Ghifary, W. B. Kleijn, and M. Zhang, "Domain adaptive neural networks for object recognition," in *PRICAI*, 2014, pp. 898–904. 1, 9
- [48] K. M. Borgwardt, A. Gretton, M. J. Rasch, H.-P. Kriegel, B. Schölkopf, and A. J. Smola, "Integrating structured biological data by kernel maximum mean discrepancy," *Bioinformatics*, vol. 22, no. 14, pp. e49–e57, 2006. 1
- [49] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *CVPR*, vol. 2. IEEE, 2006, pp. 2169–2178. 2, 4
- [50] D. Yoo, S. Park, J.-Y. Lee, and I. S. Kweon, "Fisher kernel for deep neural activations," in *CVPRW*, 2015. 2, 4, 6
- [51] H. Jégou, M. Douze, C. Schmid, and P. Pérez, "Aggregating local descriptors into a compact image representation," in *CVPR*, 2010, pp. 3304–3311. 2
- [52] B. Liu, J. Liu, J. Wang, and H. Lu, "Learning a representative and discriminative part model with deep convolutional features for scene recognition," in *ACCV*, 2014. 3, 6, 7
- [53] Z. Jiang, Z. Lin, and L. S. Davis, "Learning a discriminative dictionary for sparse coding via label consistent k-svd," in *CVPR*, 2011, pp. 1697–1704. 3
- [54] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, "Sun database: Large-scale scene recognition from abbey to zoo," in *CVPR*, 2010, pp. 3485–3492. 3, 5, 8
- [55] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "Liblinear: A library for large linear classification," *JMLR*, vol. 9, pp. 1871–1874, 2008. 3
- [56] T. Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu, "An efficient k-means clustering algorithm: Analysis and implementation," *TPAMI*, vol. 24, no. 7, pp. 881–892, 2002. 4
- [57] L.-J. Li, H. Su, L. Fei-Fei, and E. P. Xing, "Object bank: A high-level image representation for scene classification & semantic feature sparsification," in *NIPS*, 2010, pp. 1378–1386. 4, 6
- [58] K. Kavukcuoglu, P. Sermanet, Y.-L. Boureau, K. Gregor, M. Mathieu, and Y. L. Cun, "Learning convolutional feature hierarchies for visual recognition," in *NIPS*, 2010, pp. 1090–1098. 4
- [59] G.-S. Xie, X.-Y. Zhang, and C.-L. Liu, "Efficient feature coding based on auto-encoder network for image classification," in *ACCV*, 2014. 4
- [60] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, "Speaker verification using adapted gaussian mixture models," *Digital signal processing*, vol. 10, no. 1, pp. 19–41, 2000. 4
- [61] Y. Jia, "Caffe: An open source convolutional architecture for fast feature embedding," <http://caffe.berkeleyvision.org>, 2013. 5, 7
- [62] J. Zhu, L.-J. Li, L. Fei-Fei, and E. P. Xing, "Large margin learning of upstream scene understanding models," in *NIPS*, 2010, pp. 2586–2594. 6
- [63] J. Wu and J. M. Rehg, "Centrist: A visual descriptor for scene categorization," *TPAMI*, vol. 33, no. 8, pp. 1489–1501, 2011. 6
- [64] S. N. Parizi, J. G. Oberlin, and P. F. Felzenswalb, "Reconfigurable models for scene recognition," in *CVPR*, 2012, pp. 2775–2782. 6
- [65] Y. Zheng, Y.-G. Jiang, and X. Xue, "Learning hybrid part filters for scene recognition," in *ECCV*, 2012, pp. 172–185. 6
- [66] F. Sadeghi and M. F. Tappen, "Latent pyramidal regions for recognizing scenes," in *ECCV*, 2012, pp. 228–241. 6
- [67] Q. Li, J. Wu, and Z. Tu, "Harvesting mid-level visual concepts from large-scale internet images," in *CVPR*, 2013, pp. 851–858. 6
- [68] X. Wang, B. Wang, X. Bai, W. Liu, and Z. Tu, "Max-margin multiple-instance dictionary learning," in *ICML*, 2013, pp. 846–854. 6
- [69] J. Sun and J. Ponce, "Learning discriminative part detectors for image classification and cosegmentation," in *ICCV*, 2013, pp. 3400–3407. 6
- [70] C. Xu, C. Lu, J. Gao, W. Zheng, T. Wang, and S. Yan, "Discriminative analysis for symmetric positive definite matrices on lie groups," *IEEE Transactions on Circuits and Systems for Video Technology*, 2015. 6
- [71] R. Kwitt, N. Vasconcelos, and N. Rasiwasia, "Scene recognition on the semantic manifold," in *ECCV*, 2012, pp. 359–372. 8
- [72] R. Margolin, L. Zelnik-Manor, and A. Tal, "Otc: A novel local descriptor for scene classification," in *ECCV*, 2014, pp. 377–391. 8
- [73] Y. Su and F. Jurie, "Improving image classification using semantic attributes," *IJCV*, vol. 100, no. 1, pp. 59–77, 2012. 8
- [74] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek, "Image classification with the fisher vector: Theory and practice," *IJCV*, vol. 105, no. 3, pp. 222–245, 2013. 8
- [75] T. Tommasi and B. Caputo, "Frustratingly easy nbnn domain adaptation," in *ICCV*, 2013, pp. 897–904. 9
- [76] L. Wang, S. Guo, W. Huang, and Y. Qiao, "Places205-vggnet models for scene recognition," *arXiv:1508.01667*, 2015. 5, 10



Guo-Sen Xie is currently pursuing his Ph.D. degree in pattern recognition and intelligent systems, at the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China. His research interests include machine learning, deep learning, and their applications to object recognition and DNA sequence analysis.



Xu-Yao Zhang received the BS degree in computational mathematics from Wuhan University, Wuhan, China, in 2008, and the PhD degree in pattern recognition and intelligent systems from Institute of Automation, Chinese Academy of Sciences, Beijing, China, in 2013. From July 2013, he has been an Assistant Professor at the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, China. His research interests include machine learning, pattern recognition, handwriting recognition, and deep learning.



Shuicheng Yan is currently an Associate Professor at the Department of Electrical and Computer Engineering at National University of Singapore, and the founding lead of the Learning and Vision Research Group (<http://www.lv-nus.org>). Dr. Yan's research areas include machine learning, computer vision and multimedia, and he has authored/co-authored hundreds of technical papers over a wide range of research topics, with Google Scholar citation > 19,000 times and H-index 51. He is ISI Highly-cited Researcher, 2014 and IAPR Fellow 2014. He has been serving as an associate editor of IEEE TKDE, TCSVT and ACM Transactions on Intelligent Systems and Technology (ACM TIST). He received the Best Paper Awards from ACM MM'13 (Best Paper and Best Student Paper), ACM MM12 (Best Demo), PCM'11, ACM MM10, ICME10 and ICIMCS'09, the runner-up prize of ILSVRC'13, the winner prize of ILSVRC14 detection task, the winner prizes of the classification task in PASCAL VOC 2010–2012, the winner prize of the segmentation task in PASCAL VOC 2012, the honourable mention prize of the detection task in PASCAL VOC'10, 2010 TCSVT Best Associate Editor (BAE) Award, 2010 Young Faculty Research Award, 2011 Singapore Young Scientist Award, and 2012 NUS Young Researcher Award.



Cheng-Lin Liu received the B.S. degree in Electronic Engineering from Wuhan University, Wuhan, China, the M.E. degree in Electronic Engineering from Beijing Polytechnic University (current Beijing University of Technology), Beijing, China, the Ph.D. degree in Pattern Recognition and Intelligent Systems from the Institute of Automation of Chinese Academy of Sciences, Beijing, China, in 1989, 1992 and 1995, respectively. He was a postdoctoral fellow at Korea Advanced Institute of Science and Technology (KAIST) and later at Tokyo University of Agriculture and Technology from March 1996 to March 1999. From 1999 to 2004, he was a research staff member and later a senior researcher at the Central Research Laboratory, Hitachi, Ltd., Tokyo, Japan. From 2005, he has been a Professor at the National Laboratory of Pattern Recognition (NLPR), Institute of Automation of Chinese Academy of Sciences, Beijing, China, and is now the Director of the laboratory. His research interests include pattern recognition, image processing, neural networks, machine learning, and the applications to character recognition and document analysis. He has published over 200 technical papers at prestigious international journals and conferences. He is a Fellow of the IAPR and the IEEE.