

# Scene Parsing With Integration of Parametric and Non-Parametric Models

Bing Shuai, *Student Member, IEEE*, Zhen Zuo, *Student Member, IEEE*,  
Gang Wang, *Member, IEEE*, and Bing Wang, *Student Member, IEEE*

**Abstract**—We adopt convolutional neural networks (CNNs) to be our parametric model to learn discriminative features and classifiers for local patch classification. Based on the occurrence frequency distribution of classes, an ensemble of CNNs (CNN-Ensemble) are learned, in which each CNN component focuses on learning different and complementary visual patterns. The local beliefs of pixels are output by CNN-Ensemble. Considering that visually similar pixels are indistinguishable under local context, we leverage the global scene semantics to alleviate the local ambiguity. The global scene constraint is mathematically achieved by adding a global energy term to the labeling energy function, and it is practically estimated in a non-parametric framework. A large margin-based CNN metric learning method is also proposed for better global belief estimation. In the end, the integration of local and global beliefs gives rise to the class likelihood of pixels, based on which maximum marginal inference is performed to generate the label prediction maps. Even without any post-processing, we achieve the state-of-the-art results on the challenging SiftFlow and Barcelona benchmarks.

**Index Terms**—Scene parsing, convolution neural network, CNN-ensemble, global scene constraint, local ambiguity, deep learning.

## I. INTRODUCTION

SCENE parsing (also termed as scene labeling, scene semantic segmentation) builds a bridge towards deeper scene understanding. The goal is to associate each pixel with one semantic class. Generally, “thing” pixels (car, person, etc) in real world images can be quite visually different due to their scale, illumination and pose variation, meanwhile “stuff” pixels are usually visually similar (road, sea, etc) in a local close-up view. Hence, the local classification for pixels is challenging. Besides, the class frequency distribution is highly imbalanced in natural scene images: more than 80% pixels in

Manuscript received August 26, 2015; revised December 3, 2015 and January 11, 2016; accepted February 1, 2016. Date of publication February 23, 2016; date of current version April 7, 2016. This work was supported in part by the Rapid-Rich Object Search (ROSE) Laboratory, Nanyang Technological University, Singapore, in part by the Singapore Ministry of Education under Grant Tier 2 ARC28/14, and in part by the Singapore Agency for Science, Technology and Research within the Science and Engineering Research Council under Grant PSF1321202099. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Guoliang Fan. (*Corresponding author: Gang Wang.*)

The authors are with the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore 639798 (e-mail: bshuai001@ntu.edu.sg; zzu01@ntu.edu.sg; wanggang@ntu.edu.sg; wang0775@ntu.edu.sg).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIP.2016.2533862



Fig. 1. Motivation of our integration model: the parametric model can distinguish visually different pixels very well, but get confused for pixels that are visually similar in local context. However, the local ambiguities can be easily eliminated as long as the correct global scene semantics is revealed. A more consistent labeling result can be achieved by integrating their beliefs. The figure is best viewed in color.

the images belong to only a few number of semantic classes. Thus, the classification model is biased towards frequent classes due to the scarcity of training instances for rare classes. Overall, these issues pose scene parsing as one of the most challenging problems in computer vision.

The recent advance of Convolutional Neural Networks (CNNs) [1], [2] has revolutionized the computer vision community due to their outstanding performance in a wide variety of tasks [1], [3]–[7]. Recently, Farabet et al. [8] and Pinheiro and Collobert [9] has applied CNNs to scene labeling. In this scenario, CNNs are used to model the class likelihood of pixels directly from local image patches. They are able to learn strong features and classifiers to discriminate the local visual subtleties. In general, single CNN fails to produce satisfactory labeling results due to the severely imbalanced class frequency distribution in natural scene images (as exemplified in Figure 3). To address this issue, we propose the **CNN-Ensemble**, which aggregates the predictions from different and complementary CNNs. The CNN component shares the identical network architecture, but it is trained from image patches with disparate statistics,

which are generated from different sampling methods. In the end, the proposed CNN-Ensemble is capable of yielding more reliable labeling results than any single CNN.

Even though powerful, CNN-Ensemble still struggles in differentiating visually similar pixels as it only consider limited context in local classification. As exemplified in Figure 1, the sand pixels are confused with road and sidewalk pixels in a local view. We refer to such problem as *local ambiguity*. This problem has usually been addressed from two perspectives:

- Augmenting the scale of context to represent pixels: [8] considers multi-scale context input, [9] increases the size of context input in a recurrent CNN framework. These methods somehow mitigate the local ambiguity, however they may have a negative effect to small-size objects and may also degrade the efficiency of the system.
- Building a probabilistic graphical model to capture the explicit label dependencies among pixels [10]–[13]. However, the parametric graphical model is usually hard and inefficient to optimize when the higher order potentials are involved, and the low-order potentials suffer from low representation power.

Here in this paper, we propose to utilize the global scene semantics to eliminate the *local ambiguity*. As a simple example in Figure 1, the confusion between ‘road’ and ‘sand’ pixels can be easily removed if the global “coast” scene is revealed. Intuitively, a global scene constraint is implicitly enforced to allow more reliable local classification. Such global constraint is mathematically achieved by adding a global energy term to the labeling energy function. However, due to the extraordinarily huge labeling space, it’s infeasible to model the global energy parametrically. Thus, the global energy is practically modeled in a non-parametric framework by transferring the class dependencies and priors from its global similar exemplar images.

Furthermore, a large margin based metric learning objective is introduced to fine tune the network, thus making the estimation of global belief more accurate. Finally, the class likelihood of pixels are obtained by integrating the local and global beliefs. Based on which, our integration model outputs the label prediction maps. We justify our method on the popular and challenging scene parsing benchmarks: SiftFlow [14] and Barcelona [15] datasets. Even without any post-processing, our integration model is able to achieve very competitive results that are on par with the state-of-the-arts. Overall, the contributions of this paper are summarized as follows:

- 1) We propose the **CNN-Ensemble**, in which each CNN component concentrates on learning distinct visual patterns. The aggregation of single CNNs gives rise to much more powerful model that is able to generate more reliable labeling maps.
- 2) We leverage global scene semantics to remove the *local ambiguity* by transferring class dependencies and priors from similar exemplars.
- 3) We introduce the CNN metric, and show that the learned metrics are beneficial in our non-parametric global belief estimation.

This paper is an extension to the conference paper [16]. The rest of the paper is organized as follows: section II firstly

reviews, discusses and compares our methods with relevant works. Following that, the formulation of our integration model are presented in Section III; Then, details of estimating the local and global beliefs are elaborated in Section IV and V respectively; Section VI demonstrates the experimental setup and reports the results of the proposed methods; Section VII concludes the paper.

## II. RELATED WORK

Scene parsing has attracted more and more attention in recent years. Among all the interesting works, we review and discuss four line of works that are most relevant to ours.

### A. Feature Learning

The first direction exploits extracting better features for classifying pixels/superpixels. Previously, low-level and mid-level hand-crafted features are designed to capture different image statistics. They usually lack discriminative power and suffer from high dimensionality, thus limiting the complexity of the full labeling system. Recently, machine learning techniques are commonly used to learn discriminative features for various computer vision tasks [8], [9], [17], [18]. In accordance, Farabet et al. [8] fed a convolutional neural network with multi-scale raw image patches, and they have presented very interesting results on real-world image scene labeling benchmarks. Furthermore, Pinheiro and Collobert [9] adopted a recurrent CNN to process the large-size image patches. Bulò and Kotschieder [19] learned a more compact random forest by substituting the random split function with a stronger Neural Network. Mostajabi and Gholampour [20] and Mostajabi *et al.* [21] extracted features from different scopes of image regions, and then concatenated them to yield the context-aware representation. Therefore, regional and global context is encoded in the local representation. In their works, the local disambiguation is achieved via augmenting input context. In contrast, we leverage the global scene constraint to mitigate the local ambiguity.

### B. Probabilistic Graphical Models

Another line of works focus on exploring explicit dependency modeling among labels, which is usually formulated as a structure learning problem. Shotton et al. [12] formulated the unary and pairwise features in a 2nd-order sparse Conditional Random Fields (CRF) graphical model. Roy and Todorovic [22], Zhang and Chen [13] and Chen et al. [23] built a fully connected graph to enforce higher order labeling coherence. Kohli et al. [11], Kotschieder et al. [24] and Márquez-Neila et al. [25] modeled the higher order relations by considering patch/superpixel as a clique. He et al. [10] defined a multi-scale CRF that captures different contextual relationships ranging from local to global granularity. Zheng et al. [26] formulates the CRF as a neural network, so its inference equals to applying the same neural network recurrently until some fixed point (convergence) is reached. Recently, Shuai et al. [27], [28] adopt recurrent neural networks (RNNs) to propagate local

contextual information and it shows superiority over PGMs on the applicability to large-scale scene parsing task. Our work is related to this branch of works, but approaches from a different angle. The potentials in these works are usually modeled parametrically, therefore extensive efforts are needed for learning these parameters. Our global energy term can be estimated very efficiently in a non-parametric framework.

### C. Label Transfer Models

Recently, non-parametric label transfer methods [14], [15], [29]–[32] have gained popularity due to their outstanding performance and the scalability to large scale data. They usually estimate the class likelihood of image unit from the globally similar images. In a nutshell, global scene features are firstly utilized to retrieve the relevant images, whose label maps are then leveraged to estimate the class likelihood of image units. The pioneering label transfer work [14] transformed RGB image to SIFT [33] image, which was used to seek correspondences over pixels. Then, an energy function was defined over pixel correspondences, and the label prediction maps are obtained by minimizing the energy. The Superparsing system [15] performed label transfer over image superpixels. Eigen and Fergus [29] learned adaptive weights for each low-level features, and it resulted in better nearest neighbor search. Gould and Zhang [34] and Gould *et al.* [35] built a graph for dense image patch and superpixel to achieve the label transfer. We adopt this framework to evaluate our global energy term. In comparison with these works that are based on hand-crafted features, we used the learned CNN features which are more compact and discriminative. We expect our features to benefit their systems in terms of accuracy and efficiency as well.

### D. Ensemble Models

The ensemble methods [36], [37] have achieved great success in machine learning. The idea is to build a strong predictors by assembling many weak predictors. The assumption is that these weak predictors are complementary to each other when they are trained with different subset of features and training data. Some examples are random forest [38], [39], bagging [40], boosting [41], etc. Random forest has been successfully used in solving image labeling problems. For example, Shotton *et al.* [37] learned an ensemble of decision jungles to output the semantic class labels of pixels. Kotschieder *et al.* [24] constructed a random forest that directly maps the image patches to their corresponding structured output maps. Our CNN-Ensemble is different from these works. First, the individual model in theirs are very weak, whereas each of our single CNN has very strong capability. Second, the data that are used for each individual model training in their works are sampled without discrimination. In contrast, data that are fed to each single CNN are sampled differently, therefore their statistics are different.

Recently, the ensemble models [1], [42], [43] have been pervasively adopted in the large-scale ImageNet classification competition [44]. Specifically, many deep neural networks (their network architecture is identical) are trained, and the

fusion of their decisions give rise to the output. By doing this, the ensemble model is able to enhance its classification accuracy slightly. Our CNN-Ensemble also differs from these works. In their works, each network component is trained with exactly the same data, and the difference of network components mainly originates from nonidentical network initializations. In contrast, every CNN component is trained from entirely different data, which will guide the CNN component to focus on learning different but complementary visual patterns.

## III. FORMULATION

The image labeling task is usually formulated as a discrete energy minimization problem. Specifically in this paper, we consider minimizing the following energy:

$$E(X, Y) = E_I(X, Y) + E_G(X, Y) \quad (1)$$

where  $X = \{X_1, X_2, \dots, X_N\}$  is the observed image and  $X_i$  corresponds to the  $i$ th pixel;  $Y = \{Y_1, Y_2, \dots, Y_N\}$ ,  $Y \in \{1, 2, \dots, |L|\}^N$  denotes a labeling configuration for image  $X$ ;  $E_I(X, Y)$  and  $E_G(X, Y)$  are the local and global energy term respectively.

Individually speaking, the local energy term  $E_I(X, Y)$  measures the likelihood of image  $X$  taking the labeling configuration  $Y$  in the local view. Mathematically, it is expressed as the summation of local unary potential  $\Psi_I(X_i, Y_j)$ :

$$E_I(X, Y) = \sum_{X_i \in X} \Psi_I(X_i, Y_j) \quad (2)$$

where  $\Psi_I(X_i, Y_j) = -\log(P_I(X_i, Y_j))$  is defined as the negative log likelihood of pixel  $X_i$  being labeled as  $Y_j$ ; Hereafter, we call  $P_I(X_i, Y_j)$  the local belief.

The global energy  $E_G(X, Y)$ , on the other hand, reflects the likelihood of image  $X$  taking the labeling configuration  $Y$  in the global view. A naive implementation is to consider all the pixels to be in a single clique, and  $E_I(X, Y)$  evaluates the energy term according to the labeling states. However in this scenario, the labeling states are prohibitively huge ( $|L|^N$ ), it makes the energy evaluation intractable in practice. Inspired by [25], we adopt a non-parametric approach to decompose the global energy to the aggregation of global unary potential  $\Psi_G(X_i, Y_j)$ :

$$E_G(X, Y) = \sum_{X_i \in X} \Psi_G(X_i, Y_j) \quad (3)$$

where  $\Psi_G(X_i, Y_j) = -\log(P_G^{S(X)}(X_i, Y_j))$ ;  $P_G^{S(X)}(X_i, Y_j)$  denotes the likelihood of  $X_i$  being  $Y_j$  in the global scene view, and it is estimated from images in  $\mathcal{S}(X)$  which are expected to share similar global scene semantics with image  $X$ . Likewise, we call  $P_G^{S(X)}(X_i, Y_j)$  the global belief hereafter. As shown, the global energy implicitly captures the pixel dependencies via transferring scene semantics from images in  $\mathcal{S}(X)$ . As an example in Figure 1, the global scene exemplars ( $\mathcal{S}(X)$ ) define a ‘‘coast’’ scene, in which road and sidewalk pixels are invalid and sand pixels are more likely to appear in the bottom regions.

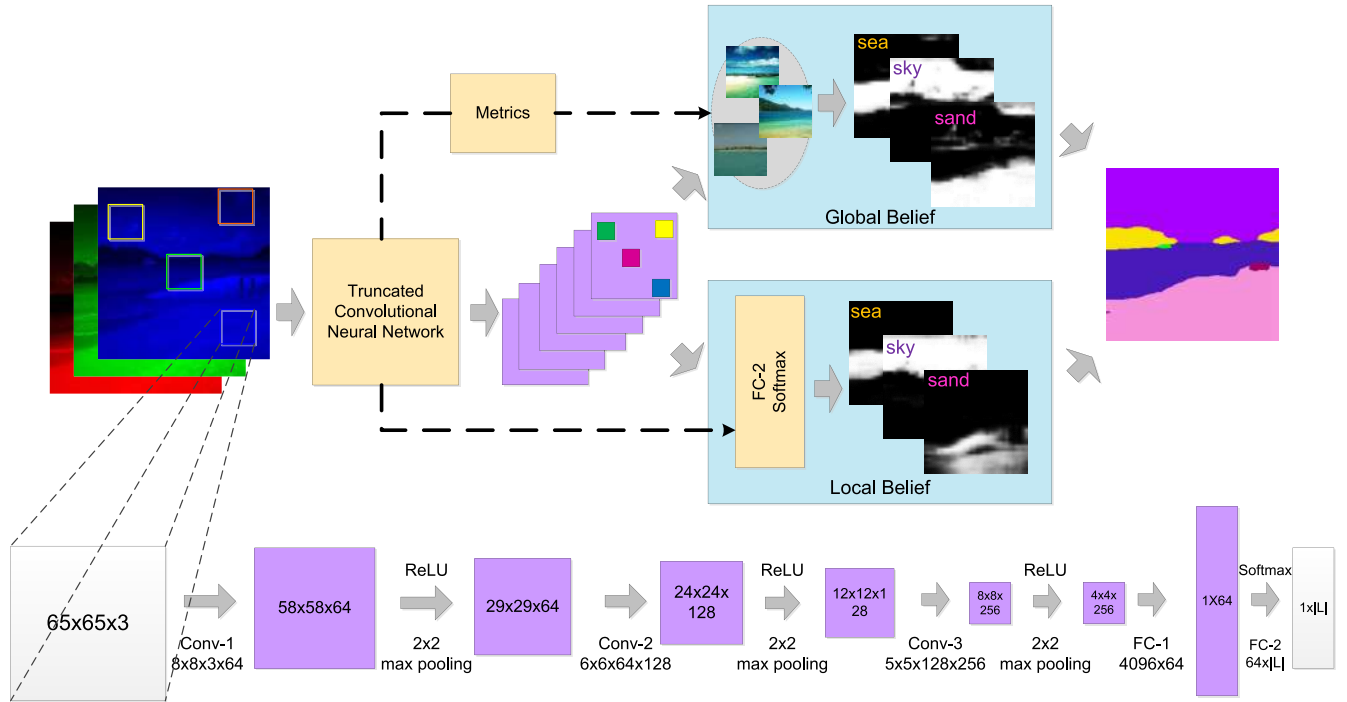


Fig. 2. Framework of our integration model: the parametric CNNs are responsible for emitting the local beliefs, and the non-parametric model is practically used to output the global beliefs. The integration of local and global beliefs gives rise to the un-normalized class likelihood of pixels. The modules painted in yellow represents parametric models (CNN and Metrics).

Finally, by rewriting the energy functions, we generate the following form:

$$E(X, Y) = - \sum_{X_i \in X} \log((P_I(X_i, Y_j) \cdot P_G^{S(X)}(X_i, Y_j))) \quad (4)$$

The above energy is numerically proportional to the integration of beliefs from two sources: (1), Local belief:  $P_I(X_i, Y_j)$  measures the belief for pixel  $X_i$  based on its surrounding local context; (2), Global belief:  $P_G^{S(X)}(X_i, Y_j)$  denotes the belief for  $X_i$  from the global scene view. An intuitive interpretation of Equation 4 is that a global scene constraint (prior) is enforced to the local classification in the form of weighting the local beliefs of pixels with their corresponding global beliefs. Since the estimation of class likelihood for different pixels is independent, the energy minimization (inference) can be done in an efficient pixel-wise manner:  $Y = \bigcup_{i=1:N} Y_i$ ,  $Y_i = \operatorname{argmin}_{1, \dots, |L|} E(X_i, Y_j)$ .

The pipeline of our model is depicted in Figure 2. Systematically, an image is first passed to the truncated CNN, and the corresponding pixel feature maps are generated. Next, the feature maps are fed into two branches: (1), they are independently classified based on the parametric CNNs (CNN-Ensemble), which yield the local beliefs; (2), they are aggregated to produce the global scene envelop, which is used to retrieve the global similar exemplar images. Based on which, the global belief is estimated. Finally, the integration of local and global beliefs yields the un-normalized class likelihood of pixels, based upon which the integration model outputs the label prediction map. We elaborate each module in the following sections.

#### IV. LOCAL BELIEFS

In a local view, the semantic class of each pixel  $X_i$  is determined by its surrounding image region (patch). A parametric model is commonly used to emit its local belief  $P_I(X_i, Y_j)$ . In this paper, this model is parameterized by **CNN-Ensemble** - an ensemble of Convolutional Neural Network (CNNs). The CNN components are trained from entirely different image regions (patches) in terms of the class frequency distribution, thus they capture complementary image statistics. In detail, each CNN component is enforced to focus on discriminating some specific classes by adaptively learning different features and classifiers. By fusing these complementary cues, CNN-Ensemble is expected to output more reliable labeling results.

##### A. Parametric CNNs

The Convolutional Neural Networks (CNNs) learn features and classifiers in an end-to-end trainable system. They are able to learn compact yet discriminative representations, which are easy to be differentiated for the jointly learned classifiers. Specifically in the image labeling task, the parameters of CNNs are optimized based on the image patches, whose labels are associated with the centering pixels. Unlike other CNNs [8], [9] which are fed with multi-scale or large field-of-view patches, we use moderate-size contextual windows to predict the labels for pixels. By doing this, we enforce the network to learn strong representations to discriminate the local visual subtleties. Meanwhile, the locality information is well preserved, which is crucial to differentiate small-size object classes. Moreover, as evidenced by the experiments

later, our CNNs outperform their nets [8], [9] dramatically. In addition, our CNNs are more efficient in terms of inference.

The architecture of our CNNs is demonstrated in Figure 2. It accepts  $65 \times 65$  image patches as valid input. If we assume that the last fully connected layer (FC-2) serves as the functionality of classifiers, the removal of which in the CNN gives rise to the Truncated Convolutional Neural Network. In other words, if we pass an image patch ( $65 \times 65$ ) to the truncated CNN, its output is a representation vector (64 dimension), which summarizes the contextual information of the input patch. In this perspective, truncated CNN can be interpreted as a feature extractor.

### B. Data Sampling

Knowing that the number of training patches are prohibitively huge (thousands of millions), we only use a fraction of them for the network training. Specifically, at the beginning of each epoch, training patches are randomly sampled. By doing this, we are able to decrease the training time dramatically. Meanwhile, the sampling strategy does not harm the performance, as the image patches are highly redundant (patches that are related to neighborhood pixels are usually the same) and the randomness injected into the data sampling during each epoch enables the network to “see” the whole data throughout the whole training process. Here in this paper, we introduce four sampling methods based on the class frequency distribution in natural scene images:

- **Global sampling (GS):** It samples patches randomly from the collection of all the training patches. Then the class frequency distribution in the sampled patches should be very close to that in the original images.
- **Class sampling (CS):** It samples patches in a manner that classes appear equally in the resulting sampled patches. Note that some patches may appear multiple times if they are from extremely rare classes.
- **Hybrid sampling (HS):** It is a compromise between global sampling and class sampling. In detail, it firstly samples patches globally, and then augments the rare-class patches until their occurrence frequencies reach the desired threshold  $\eta$  in the sampled patches.
- **Truncated class sampling (TCS):** It adopts the same sampling procedure as class sampling, but removes all the frequent-class patches.

The threshold  $\eta$  is used to determine whether a class is frequent or rare: the class is considered to be frequent as long as its occurrence frequency in the training data is above  $\eta$ , otherwise, it belongs to a rare class. Obviously, the above sampling methods are expected to yield significantly different sampled patches in the form of class frequency distribution. When these image patches are presented in the network training phase, the CNNs are enforced to learn disparate visual patterns, thus they behave differently. We next discuss how the sampled image patches influence the characteristics of CNNs.

### C. CNN-Ensemble

CNNs, even with the identical architecture, can be functionally different if they are trained from image patches with

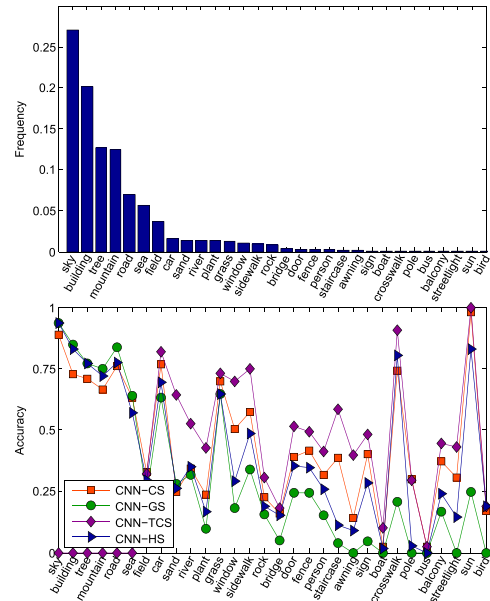


Fig. 3. The first graphic shows an example of the class frequency distribution in natural scene images, through which we can observe that large percentage of pixels belong to very small number of frequent classes. The second figure delineates the class-wise accuracy curves for CNN-M. The experimental statistics are based on the SiftFlow dataset [14].

disparate data distribution. Concretely, to fit the image patches whose statistics are not the same, CNNs adaptively learn different representations and classifiers. In consequence, it leads to the big performance discrepancy on labeling the same images. Let’s denote **CNN-M** the CNN trained with image patches sampled through method M. As shown in Figure 3, CNN-M exhibits significantly different characteristics towards predicting the semantic classes of pixels: **CNN-GS** prioritizes the frequent classes, therefore it optimizes the overall pixel accuracy; **CNN-CS** implicitly normalizes the weights of classes by performing downsampling and over-sampling operations to frequent and rare classes respectively, thus it maximizes the average class accuracy; In comparison with CNN-GS, **CNN-HS** gives slightly higher weights to rare classes, hence it compromises the two above criterions; **CNN-TCS** works exceptionally well on differentiating rare classes.

To produce a satisfactory labeling prediction map, the algorithm is required to perform extraordinarily well to correctly predict frequent classes, and in the mean time should work well towards recognizing rare classes precisely. The former criterion guarantees that the scene semantics of the images are well defined, and the latter one enforces that objects (rare classes) are not missing in the scene. However as manifested in Figure 3, none of the CNN-M satisfies both criterions. Therefore, we propose the **CNN-Ensemble**, which combines the predictions from the CNN-M components. Mathematically, the local belief  $P_I(X_i, Y_j)$  is derived with the following equation:

$$P_I(X_i, Y_j) = \frac{1}{M} \sum_{m=1}^M P_m(X_i, Y_j) \quad (5)$$

where  $M$  is the number of CNN-M components and  $P_m(X_i, Y_j)$  denotes the class likelihood prediction from the  $m$ -th CNN-M. As mentioned, CNN-M components are designed to focus on distinguishing some specific classes. For example, CNN-GS discriminates frequent classes excellently, and CNN-TCS captures the subtleties of infrequent classes. The proposed CNN-Ensemble fuses the complementary predictions from different sources. Our later experiments in Section VI will demonstrate that CNN-Ensemble is able to produce much more reliable labeling results quantitatively, and it performs excellently on both overall pixel accuracy and average per-class accuracy.

## V. GLOBAL BELIEFS

In a global view, the semantic class of each pixel  $X_i$  in image  $X$  is determined by the global scene semantics of  $X$ . In other words, its global class likelihood  $P_G(X_i, Y_j)$  should match the expected scene layout of  $X$ . Even though the pixel  $X_i$  is represented identically, its class belief could vary significantly depending on how the pixel is evaluated. Take the simple example as an illustration in Figure 1, the pixels in the lower part of the image could be ‘sand’, ‘road’ or even ‘rock’ in a local view. We refer to this problem as *local ambiguity*. In contrast, with the awareness of global ‘coast’ scene, it is obvious that ‘sand’ class is preferred in a global view. In this perspective, the global scene prior is a good remedy to alleviate the local ambiguity.

### A. Non-Parametric Global Belief Transfer

As elaborated, the parametric CNNs (CNN-Ensemble) are able to produce good labeling results for the pixels with good local contextual support, they still suffer from the notorious *local ambiguity* problem. Previously, researchers usually addressed this issue by generating contextual aware local features. For example, Farabet et al. [8] fed the network with multi-scale image patches to yield richer contextual aware local features, and likewise Pinheiro and Collobert [9] took the network input as larger image patches. In this paper, the local disambiguation is achieved by enforcing a global scene constraint to local classification. More specifically, a pixel is considered under global context: the class likelihood of pixels should satisfy the scene layout and semantics of the image.

First, we generate the global-level representation for the considered image  $X \in \mathbb{R}^{h \times w \times 3}$ , where  $h, w$  is the height and width of the image  $X$ . The corresponding CNN feature tensor  $F \in \mathbb{R}^{h \times w \times M}$  can be obtained by passing densely sampled image patches in  $X$  to the truncated CNN ( $M = 64$  in our implementation). Next, we introduce the average pooling operator *pool* [1] to aggregate the pixel features, thus giving rise to the global feature  $H$ . In detail, suppose an image is decomposed to regions  $\mathcal{R} = \{R^{(1)}, R^{(2)}, \dots, R^{(J)}\}$ ,<sup>1</sup> the region feature is generated by applying *pool* operator to the constituent pixel features:  $H(R^{(i)}) = \text{pool}(F^i), \forall i \in R^{(i)}$ . The global image representation is defined as the concatenation of region

features  $H = [H(R^{(1)}), H(R^{(2)}), \dots, H(R^{(J)})]$ . As expected, this global image feature  $H$  not only conveys discriminative scene semantics but also encodes scene layout information.

Then, based on  $H$ , the global nearest exemplars  $\mathcal{S}(X)$  are retrieved. Each image in  $\mathcal{S}(X)$  is expected to have the similar scene semantics and layout with image  $X$ . After that, the global class likelihood of pixels (global belief) are transferred from the statistics of pixel features in  $\mathcal{S}(X)$ . Concretely, among all the pixels in  $\mathcal{S}(X)$ , the semantic class of pixel  $X_i$  should match those pixels whose local representations are also close to  $X_i$ . Therefore,  $K$  pixels (from images within  $\mathcal{S}(X)$ ) are firstly retrieved that are similar to  $X_i$  in the local representation space. Then the global belief is generated through a weighted voting based on the  $K$  retrieved pixels. Mathematically, it is derived in the following equation:

$$P_G^{\mathcal{S}(X)}(X_i, Y_j) = \frac{\sum_{k=1}^K \phi(X_i, X_k) \delta(Y(X_k) = Y_j)}{\sum_k \phi(X_i, X_k)} \quad (6)$$

where  $X_k$  is the  $k$ -th nearest neighbor pixel of  $X_i$  among all the pixels in  $\mathcal{S}(X)$ ;  $Y(X_k)$  is the ground truth label for pixel  $X_k$ ;  $\delta(Y(X_k), Y_j)$  is an indicator function;  $\phi(X_i, X_k)$  measures the similarity between  $X_i$  and  $X_k$ , which is defined over spatial and feature space:

$$\phi(X_i, X_j) = \exp(-\alpha \|x_i - x_j\|) \exp(-\gamma \|z_i - z_j\|) \quad (7)$$

where  $x_i = F(X_i)$  denotes the CNN pixel feature for  $X_i$ ,  $z_i$  is the normalized coordinate along the image height axis and  $\alpha, \gamma$  controls the belief exponential falloff.

Meanwhile, as small-size object classes (e.g. ‘bird’ in the sky, ‘boat’ in the sea, etc) make negligible contribution to the global scene semantics, they may not appear in its nearest global exemplar images  $\mathcal{S}(X)$ . Thus, the global belief will be highly skewed to frequent classes, which potentially harms the final class likelihood of rare classes according to the integration rule in Equation 4. To address this issue, we introduce an auxiliary pixel transfer set  $\mathcal{A}(X)$ . In detail, we first find the rare-class pixels whose quantities are below  $K$  (same value as in Equation 6) in  $\mathcal{S}(X)$ , and then augment the corresponding rare-class pixels until their quantities reach  $K$ . More specifically,  $\mathcal{A}(X)$  is a set of rare-class pixels, and it is derived by randomly sampling the desired quantity of pixels from images outside  $\mathcal{S}(X)$ . Thus, the final transfer set is the combination of  $\mathcal{S}(X)$  and  $\mathcal{A}(X)$ , and the quantity of every rare-class pixel is at least  $K$ . Thus, the global belief is expected to preserve the salient objects in the scene. Our non-parametric global belief estimation is reminiscent of popular label transfer works [14], [15], [29]–[31], two differences need to be highlighted:

- Instead of adopting hand-engineered low-level local and global features, we use more discriminative and compact features learned from CNN for label transfer.
- Our non-parametric model works as global scene constraints for local pixel features. Generally, small size retrieval images are sufficient to define the scene semantic and layout. However, previous works have to seek large retrieval set to cover all the possible semantic classes.

<sup>1</sup> $J$  is the number of regions. In our experiments, the image is divided into rectangular regions in a 2-layer spatial pyramid fashion [45].

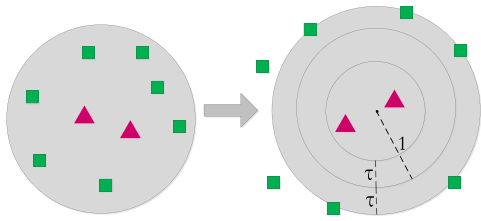


Fig. 4. Graphical illustration of the effect of our large margin based metrics. Due to the highly imbalanced data distribution, the nearest neighbors of the testing feature (triangle) are dominated by imposter classes (rectangle). After the metric transformation, the imposters stay far away from it. Thus, their contribution is significantly attenuated in the global belief estimation.

### B. CNN-Metric

As shown in Equation 6, the estimation of global belief  $P_G^{S(X)}(X_i, Y_j)$  is highly dependent on the distance metric between two pixel features. However, our features are learned by optimizing pixel/patch classification accuracy, while do not take distance metric into consideration. Therefore we propose to learn a large-margin based metric to mitigate the inaccurate class likelihood estimation for rare classes (Figure 4). In detail, the Mahalanobis metric  $M = W^T W$  is learned by minimizing the loss function, which is formally written as:

$$L = \frac{\lambda}{2} \|W\|^2 + \frac{1}{2N} \sum_{i,j} g(x_i, x_j)$$

$$g(x_i, x_j) = \max(0, 1 - \ell_{i,j}(\tau - \|Wx_i - Wx_j\|^2)) \quad (8)$$

where  $\ell_{i,j}$  indicates whether two features have the same semantic label or not, and  $\ell_{i,j} = 1$  if  $X_i$  and  $X_j$  are from the same class, or  $\ell_{i,j} = -1$  otherwise;  $\tau (> 1)$  is the margin and  $\lambda$  controls the effect of regularization;  $x_i = F(X_i)$  is the feature representation for  $X_i$  and  $N$  is the number of features. The objective function would enforce the pixel features from the same semantic class to be close and stay within the ball with radius  $1 - \tau$ , and enforce data from different classes to be far away from each other by at least  $1 + \tau$ . The graphical illustration of the metrics is depicted in Figure 4.

Instead of simply learning a metric based on the extracted CNN features, we further replace the softmax layer with our metric learning layer, so that the feature extraction parameters can also be adapted. We replace the softmax layer of previous CNN (CNN-softmax) with a fully connected layer parameterized by  $W$  (or more layers to learn non-linear metrics [46]) and fix the biases to be zero, which serves as a Mahalanobis metric ( $M = W^T W$ ). We call the new network CNN-metric. These two networks do not share any parameters except that the feature extraction parameters of CNN-metric are initiated from the corresponding layers of CNN-softmax. The errors are back propagated through the chain rule, and  $\frac{\partial L}{\partial W}$   $\frac{\partial L}{\partial x_i}$  for the last layer are given in Equation 9.

$$\frac{\partial L}{\partial W} = \lambda W + \frac{1}{N} \sum_{i,j} \zeta_{ij}$$

$$\zeta_{ij} = g'(c) \ell(i, j) (Wx_i - Wx_j) (x_i - x_j)^T$$

$$\frac{\partial L}{\partial x_i} = \frac{1}{N} \sum_{i,j} g'(c) (W^T \ell(i, j) (Wx_i - Wx_j))$$

$$c = 1 - \ell(i, j)(\tau - \|Wx_i - Wx_j\|^2)$$

$$x_i = F(X_i)$$

$$g'(c) = \begin{cases} 0, & c \leq 0 \\ 1, & c > 0 \end{cases} \quad (9)$$

We adopt the Stochastic Gradient Descent (SGD) to optimize the CNN-Metric. Considering that the quantity of patches is prohibitively huge, we also sample a fraction of patches during each epoch in the **Class Sampling** manner (Section IV). Other sampling methods are not appropriate in this scenario, as the imbalanced data distribution will result in the scarcity of training examples for some class pairs, which is expected to skew the learned metric mapping. Furthermore, due to the infeasibility of feeding the sampled patches to the network in a single propagation, the training data are divided to several batches, among which the proposed metric constraint is enforced to any class pairs.

## VI. EXPERIMENTS

### A. Evaluation Benchmarks

We evaluate our approach on two popular and challenging scene parsing benchmarks:

- SiftFlow [14]: It has 2688 images generally captured from 8 typical natural scenes. Every image has  $256 \times 256$  pixels, which belongs to one of 33 semantic classes. We use the training/testing (2488/200 images) split provided by [14] to conduct our experiments.
- Barcelona [15]: It consists of 14871 training and 279 testing images. The size of the images varies across different instances, and each pixel is labelled as one of the 170 semantic classes. Note that the class frequency distribution is more imbalanced than that in the SiftFlow dataset. Meanwhile, the scene categories of training images range from indoor to outdoors, whereas the testing images are only captured from the barcelona street. These issues pose Barcelona as an extremely challenging dataset.

To quantitatively evaluate our methods, we report two types of scores: the percentage of all correctly classified pixels - Global Pixel Accuracy (**GPA**), and Average per-Class Accuracy (**ACA**).

### B. Local Labeling Results

We first present the implementation details of training the parametric CNNs. We use Stochastic Gradient Descent (SGD) with momentum to train the CNN-M; During each training epoch, around  $5 \times 10^5$  pixels are randomly sampled from training pixel pools; The learning rate is initialized to be 0.01, and it is decreased by 10 times after 20 epoches; The momentum is fixed to 0.9, and the batch size is set as 100. We train our CNN-M based on MatConvnet [47] toolbox.<sup>2</sup> The reported results are based on the model trained in 35 epoches. Each image is preprocessed by first subtracting the mean and then performing contrast normalization by dividing its variance. Besides, the threshold  $\eta$  that discriminates rare classes are

<sup>2</sup>The code is publicly available under the homepage of authors.



Fig. 5. Visualization of the learned convolution filters for the first layer of CNN-CS, CNN-TCS, CNN-GS and CNN-HS respectively. They are trained on the Barcelona datasets. Note the visual difference among the learned filters of different CNN-M. The figure is best viewed in color.

TABLE I  
QUANTITATIVE PERFORMANCE OF DIFFERENT CNNs. DETAILS OF EACH METHOD ARE ELABORATED IN THE TEXT

	SiftFlow		Barcelona	
	GPA	ACA	GPA	ACA
Multiscale ConvNet [8]	67.9%	45.9%	37.8%	12.1%
Recurrent CNN (67×67) [9]	65.5%	20.8%	N/A	N/A
CNN-asymmetric [21]	42.4%	38.4%	20.0%	13.3%
CNN-GS	75.4%	30.2%	68.5%	11.4%
CNN-CS	70.9%	42.6%	24.7%	18.4%
CNN-TCS	7.91%	38.7%	6.33%	16.8%
CNN-HS	74.7%	39.4%	61.0%	16.7%
CNN-Ensemble	75.3%	<b>44.8%</b>	61.3%	19.5%
Ensemble CNN-GS	<b>77.1%</b>	32.0%	<b>69.7%</b>	11.5%
Ensemble CNN-CS	72.8%	43.7%	26.6%	<b>20.0%</b>
Ensemble CNN-TCS	8.33%	40.0%	7.53%	18.3%
Ensemble CNN-HS	76.4%	40.2%	63.2%	17.5%

empirically set to 5% and 1% for the SiftFlow and Barcelona datasets respectively. The learned convolutional filters for the first layer of CNN-M on the Barcelona dataset are shown in Figure 5.

Next, we evaluate the performance of CNN-M. Specifically, CNNs output their local belief maps, based on which maximum marginal inference is performed to output the label prediction maps. Table I lists the quantitative results. As shown, our CNN-M (GS, HS, CS) achieves much better results than Multiscale ConvNet and Recurrent CNN. In terms of the individual performance of CNN-M, **CNN-GS** achieves the best accuracy on global pixel accuracy (GPA), whereas its performance on average class accuracy (ACA) is not satisfactory. In contrast, **CNN-CS** claims the best ACA among all the single CNNs, which indicates that it predicts the semantic classes of pixels in a more equal manner. A favorable performance compromise is achieved by **CNN-HS**, which works considerably well on both GPA and ACA. **CNN-TCS** performs extremely poor on GPA as it ignores the frequent classes, it however achieves very competitive ACA. Importantly, CNN-TCS captures image statistics that are significantly disparate from other networks (CNN-GS, HS, CS): it performs the best to correctly recognize the rare-class pixels. Quantitatively, the CNN-Ensemble that excludes CNN-TCS only achieves 40.3% and 17.8% in terms of ACA on the SiftFlow and Barcelona dataset respectively, and their performance are boosted to 44.8% and 19.5% after including CNN-TCS. By fusing the complementary decisions from different CNN-M, **CNN-Ensemble** achieves the best performance tradeoff.

Furthermore, we train **Ensemble CNN-M** to compare their performance behaviour with CNN-Ensemble. In detail, the CNN-M in Ensemble CNN-M is trained with the image patches sampled from the identical method M. As manifested in Table I, the Ensemble CNN-M significantly improves the local labeling performance over CNN-M. However, they optimize only one evaluation criterion (either GPA or ACA), and this phenomenon is more obvious in the severely imbalanced Barcelona dataset. In contrast, the CNN-Ensemble performs competitively excellently on both GPA and ACA. It's also important to note that CNN-HS behaves similarly with CNN-Ensemble and it produces very promising results on both GPA and ACA. However, the performance of Ensemble CNN-HS is significantly inferior to CNN-Ensemble in terms of ACA.

In the end, we discuss the issue of imbalanced class frequency distribution in scene parsing. In order to boost the recognition rates for infrequent classes, we train another **CNN-asymmetric** as in [21], whose log-loss is modulated by the inverse frequency of each class, thus the rare classes are effectively given more attention. In this case, the training data can be generated as in [21] by collecting all image patches (or via global sampling strategy). The corresponding result is reported in Table I, which shows that CNN-asymmetric fails to achieve the desirable results as in object segmentation benchmarks [21]. This phenomenon can be explained from the following perspective. The class frequency distribution in the object segmentation task is not as imbalanced as in the scene parsing task.<sup>3</sup> If the inverse frequency is used to scale the log-loss in this scenario, the scaled losses w.r.t the frequent classes are negligible. Consequently, the CNN-asymmetric performs poorly on GPA on the scene parsing benchmarks. In contrast, the sampling based CNNs (e.g. CNN-CS, CNN-HS, CNN-GS) achieve much better performance tradeoff, which positively elucidates the effectiveness of the proposed sampling strategy to address the class imbalance issue. More importantly, by feeding the networks with different sampled patches during the training phase, we are able to train a number of complementary CNNs, and combine them to produce much stronger local prediction model (CNN-Ensemble).

### C. Evaluation of Global Features

In this section, we demonstrate that the pooling operation of pixel features is capable of generating semantically consistent global features. To achieve this goal, we calculate the KNN matching score  $p$  - the average genuine matching percentage among their  $K$  nearest neighbors. It is Mathematically derived in the following equation:  $p = \frac{\sum_i^N \sum_k^K \delta(i, NN(i, k))}{NK}$ , where  $N$  is the number of test images,  $NN(i, k)$  stands for the  $k$ -th nearest neighbor for image  $I_i$ , and  $\delta(i, j)$  outputs value 1 if  $i$  and  $j$  are a genuine match, or 0 otherwise. A genuine matching image pair means that they belong to the identical semantic class. We test the global features on the

<sup>3</sup>Statistically, the frequency ratio between the most frequent and rare classes on the PASCAL VOC 2011 [21] and SiftFlow datasets (see Figure 3) are approximately 240 and  $3.5 \times 10^4$  respectively.



TABLE II

AVERAGE GENUINE MATCHING PERCENTAGE AMONG THEIR K-NEAREST NEIGHBORS FOR DIFFERENT GLOBAL FEATURES. GT IS THE SEMANTIC FEATURE POOLED FROM GROUND TRUTH LABEL MAPS

	Dim	K=1	K=5	K=10
GIST [48]	512	74.0%	70.7%	68.3%
SIFT-SPM [45]	2100	76.5%	71.3%	69.1%
Global Feature (CNN-GS)	<b>320</b>	<b>91.5%</b>	<b>88.3%</b>	86.5%
Global Feature (CNN-HS)	<b>320</b>	88.5%	85.8%	85.2%
Global Feature (CNN-CS)	<b>320</b>	90.5%	84.7%	84.1%
Global Feature (CNN-TCS)	<b>320</b>	79.5%	73.9%	72.8%
Global Feature (CNN-Ensemble)	1280	<b>91.5%</b>	87.4%	<b>86.7%</b>
GT	165D	94.0%	91.0%	89.5%

SiftFlow dataset, as it provides the global scene label for each image.

Four global features are compared in our experiment: GIST [48] is a global summary of scene images that captures scene structure and layout; SIFT-SPM (GT) [45] is pooled from low-level local SIFT [33] (ground truth label map [14]) in a 3(2)-layer spatial pyramid. They are commonly used in scene classification and non-parametric label transfer framework. GT is the ideal global semantic feature. As mentioned in Section V, our global feature is pooled from the output of truncated CNN-M in a 2-layer spatial pyramid fashion. Euclidean distance is used to retrieve nearest neighbors for non-histogram features (GIST, Ours), and histogram intersection similarity measurement is applied for the rest histogram features (SIFT-SPM and GT).

The quantitative matching scores for different global features are listed in Table II. As demonstrated, our global feature is more likely to group semantically relevant images together. Meanwhile, among all of the global features (CNN-M), CNN-GS performs the best and CNN-TCS works the worst, as the scene semantics of outdoor images are mostly determined by the appearance of frequent classes. It’s worth noting that the concatenation of global features from different CNN-M fails to outperform CNN-GS. Hence, only global feature (CNN-GS) is used subsequently to retrieve the nearest exemplars. It’s also interesting to observe that the retrieval performance based on GT features are imperfect, which implies that different scenes can have very similar building blocks. For example, ‘inside city’ and ‘street’ scenes are dominated by sky and building pixels. We believe that the quality of nearest neighbor retrieval directly determines the correctness of global belief. Therefore, our global feature is expected to benefit other label transfer works as well. Some qualitative examples are shown in Figure 6, in which the retrieved nearest exemplar images have very similar scene layout.

#### D. Global Labeling Results

We first present the implementation details for our non-parametric model. The non-overlapping patches are adopted as label transfer units, within which labels are assumed to be identical. Specifically, as our CNN-M has three pooling layers, the feature extractor (truncated CNN-M) can be regarded as sliding the images with a stride of 8. In other words,

TABLE III

QUANTITATIVE PERFORMANCE OF DIFFERENT METHODS ON THE SIFTFLOW DATASETS

	Global (GPA)	Class (ACA)
SuperParsing [15]	76.9%	29.4%
Liu et al. [14]	74.8%	N/A
Gould et al. [35]	78.4%	25.7%
Singh et al. [30]	79.2%	33.8%
Tighe et al. [49]	78.6%	39.2%
Yang et al. [32]	79.8%	48.7%
Tung et al. [50]	79.9%	49.3%
George et al. [51]	<b>81.7%</b>	<b>50.1%</b>
Raw Multiscale ConvNet [8]	67.9%	45.9%
Raw Multiscale ConvNet [8] + Cover	72.3%	50.8%
Raw Multiscale ConvNet [8] + Cover	78.5%	29.4%
Plain CNN (133×133) [9]	76.5%	30.0%
Recurrent CNN (133×133) [9]	77.7%	29.8%
Gatta et al. [52]	78.7%	32.1%
Long et al. [53] (ImageNet Pretrain)	85.2%	51.7%
Local Labeling (CNN-Ensemble)	75.3%	44.8%
Global Labeling	78.7%	36.2%
Global Labeling (Metric)	78.8%	39.6%
Integration model	81.0%	44.6%
Integration model (Metric)	<b>81.2%</b>	<b>45.5%</b>

the dimension of the output feature tensor  $F$  is  $\frac{1}{8}$  of original image size: one feature in  $F$  corresponds to a  $8 \times 8$  image patch. In our experiments, we only estimate the class likelihood for each feature in  $F$ , implicitly assuming that class labels within the  $8 \times 8$  regions are the same.<sup>4</sup>

The global belief is calculated by Equation 6, in which  $|\mathcal{S}(X)|$  (size of nearest exemplar images) and  $K$  (size of nearest pixel/patch neighbors) are empirically set to be 5 and 200 respectively. As manifested in Table II, 5 images are sufficient to correctly define the scene semantics. However, as many images are not fully annotated in the Barcelona dataset, a larger retrieval image set is used ( $|\mathcal{S}(X)| = 100$ ).  $\alpha$  and  $\gamma$  in Equation 7 are empirically set to 15 and 5 respectively. To fine-tune the CNN-metric,  $1 \times 10^4$  patches are sampled for each class in each epoch (which results in  $3.3 \times 10^5$  and  $1.7 \times 10^6$  training patches on the SiftFlow and Barcelona dataset respectively), and 2000 patches are used to be a training batch.  $\lambda$  and  $\tau$  in Equation 8 makes marginal difference to the performance, and they are fixed to 0.01 and 3 respectively. The learning rate is initialized to  $10^{-3}$  and decays exponentially with the rate of 0.9. The reported results are obtained under the models learned in 20 epoches.

Similarly, the global labeling results are obtained through performing maximum marginal inference over global beliefs. The quantitative results on the SiftFlow and Barcelona dataset are listed in Table III and IV respectively. From which, we clearly observe that our simple non-parametric model achieves very promising results that are comparable or even better than most complicated label transfer counterparts. For example, we reach 78.8% (39.6%) in terms of GPA (ACA) on the SiftFlow dataset, whereas SuperParsing [15] and graph transfer [35] only achieves 76.9% (29.4%) and

<sup>4</sup>This is not the optimal setting, as it oversmooths the label prediction maps. However, it is fast and easy to implement, and we don’t expect that the global belief to preserve the boundary information as it simply reflects the global scene prior. Hence, it’s a reasonable compromise.

TABLE IV  
QUANTITATIVE PERFORMANCE OF DIFFERENT  
METHODS ON THE BARCELONA DATASET

	Global (GPA)	Class (ACA)
SuperParsing [15]	66.9%	7.6%
Raw Multiscale ConvNet [8]	37.8%	12.1%
Raw Multiscale ConvNet [8] + Cover	46.4%	<b>12.5%</b>
Raw Multiscale ConvNet [8] + Cover	<b>67.8%</b>	9.5%
Local Labeling (CNN-Ensemble)	61.3%	19.5%
Global Labeling	68.1%	13.1%
Global Labeling (Metric)	68.7%	14.1%
Integration Model	69.7%	16.8%
Integration Model (Metric)	<b>70.3%</b>	<b>17.2%</b>

78.4% (25.7%) respectively. Moreover, our global labeling alone already achieves state-of-the-art results on the challenging Barcelona dataset. We attribute the performance superiority to the highly discriminative CNN features which we work with in the non-parametric transfer model.

Furthermore, as evidenced by Table III and IV, the learned metric is capable of improving the quality of global beliefs for rare classes. As illustrated in Figure 4, the learned metric is expected to shrink the distance of pixel features between identical classes and enlarge the distance between disparate classes. Consequently, it benefits the estimation of global beliefs for infrequent classes. In our experiments, we do observe the desirable average class accuracy (ACA) boost by incorporating metric tuning on both datasets.

#### E. Integration Labeling Results

The integration model applies Equation 4 to integrate the local and global beliefs, which yields the un-normalized class likelihood of pixels. Based on which, the final prediction map is produced. Table III, IV list the quantitative performance on the SiftFlow and Barcelona dataset respectively. As shown, our integration model is able to take advantage of both models, therefore it outputs more reliable label maps quantitatively and qualitatively. On one hand, *it dramatically boosts the global pixel accuracy (GPA) compared to the local labeling (CNN-Ensemble): 5.9% and 9.0% GPA improvement on the SiftFlow and Barcelona dataset respectively.* These results elucidate that a higher quality labeling prediction map can be obtained by enforcing a global scene constraint to the local classification. On the other hand, *it also tremendously enhances the average class accuracy (ACA) compared to the global labeling: 9.3% and 4.1% ACA improvement on the Siftflow and Barcelona dataset respectively.* These results indicate that some object classes are ignorable in a global view (e.g. a small bird flies in broad sky), and the integration of local cues is akin to delineating objects in a global scene image. A number of qualitative examples are demonstrated in Figure 6.

In comparison with other representation networks [8], [9], our integration model outperforms them by a large margin. Thus, enforcing a global scene constraint to local classification is a promising solution to alleviate *local ambiguities*. Furthermore, we compare our integration model with state-of-the-art counterparts. As listed in Table III, our method achieves

TABLE V  
PER-CLASS ACCURACY COMPARISON OF DIFFERENT MODELS  
ON THE SIFTFLOW DATASET. THE STATISTICS OF CLASS  
FREQUENCY IS OBTAINED IN TEST IMAGES

	Frequency	CNN-Ensemble	Integration Model
sky	27.1%	93.5%	<b>96.7%</b>
building	20.2%	80.6%	<b>88.0%</b>
tree	12.6%	74.8%	<b>84.7%</b>
mountain	12.4%	71.5%	<b>80.3%</b>
road	6.93%	76.8%	<b>85.8%</b>
sea	5.6%	59.7%	<b>75.0%</b>
field	3.64%	33.1%	<b>37.6%</b>
car	1.6%	80.5%	78.7%
sand	1.41%	33.6%	<b>37.5%</b>
river	1.37%	44.0%	<b>50.1%</b>
plant	1.33%	27.5%	7.56%
grass	1.22%	75.2%	72.6%
window	1.07%	45.1%	33.8%
sidewalk	0.89%	60.2%	53.7%
rock	0.85%	23.8%	13.0%
door	0.26%	36.8%	<b>40.2%</b>
fence	0.24%	44.5%	<b>44.6%</b>
person	0.23%	30.9%	<b>44.2%</b>
staircase	0.18%	45.4%	24.6%
awning	0.11%	7.54%	<b>14.5%</b>
sign	0.11%	33.4%	<b>50.3%</b>
boat	0.06%	2.22%	<b>3.05%</b>
crosswalk	0.05%	85.4%	74.1%
bridge	0.04%	18.8%	<b>20.4%</b>
pole	0.04%	2.43%	<b>7.1%</b>
balcony	0.03%	34.4%	29.9%
bus	0.03%	0.05%	0.05%
streetlight	0.02%	15.1%	2.36%
sun	0.01%	93.3%	86.0%
bird	$\approx 0$	13.2%	<b>28.3%</b>
GPA	-	75.3%	<b>81.2%</b>
ACA	-	44.8%	<b>45.5%</b>

very competitive results on the SiftFlow dataset. It's important to note that [51] uses 20 types of low-level features and further augments it with Fisher Vector (FV) [54] descriptor (based on the SIFT feature) to represent each image super-pixel. In contrast, we only use very compact (64-dimension) features learned from CNNs. In addition, [53] adopts the very deep CNNs [42] pretrained on the large-scale ImageNet dataset [44] to generate the local features, whereas we utilize much shallower CNNs, which are trained with image patches only from the target dataset. Meanwhile, Table IV clearly manifests that our method achieves the new best results on the more challenging Barcelona dataset, which significantly outperforms the previous state-of-the-arts.

#### F. Analysis of Local and Global Beliefs

In this section, we first compare the characteristics of local and global beliefs. By looking into the local/global labeling results quantitatively and qualitatively, we notice that the global labeling is more likely to output globally consistent label maps, while the local labeling focuses on discriminating local regions. Specifically, the global belief prioritizes the background classes, whereas the local belief works excellently on differentiating locally distinct object classes (especially small size classes). Hence, the global belief defines the global scene prior, while local belief preserves the locality information. Moreover, the reciprocity of local and

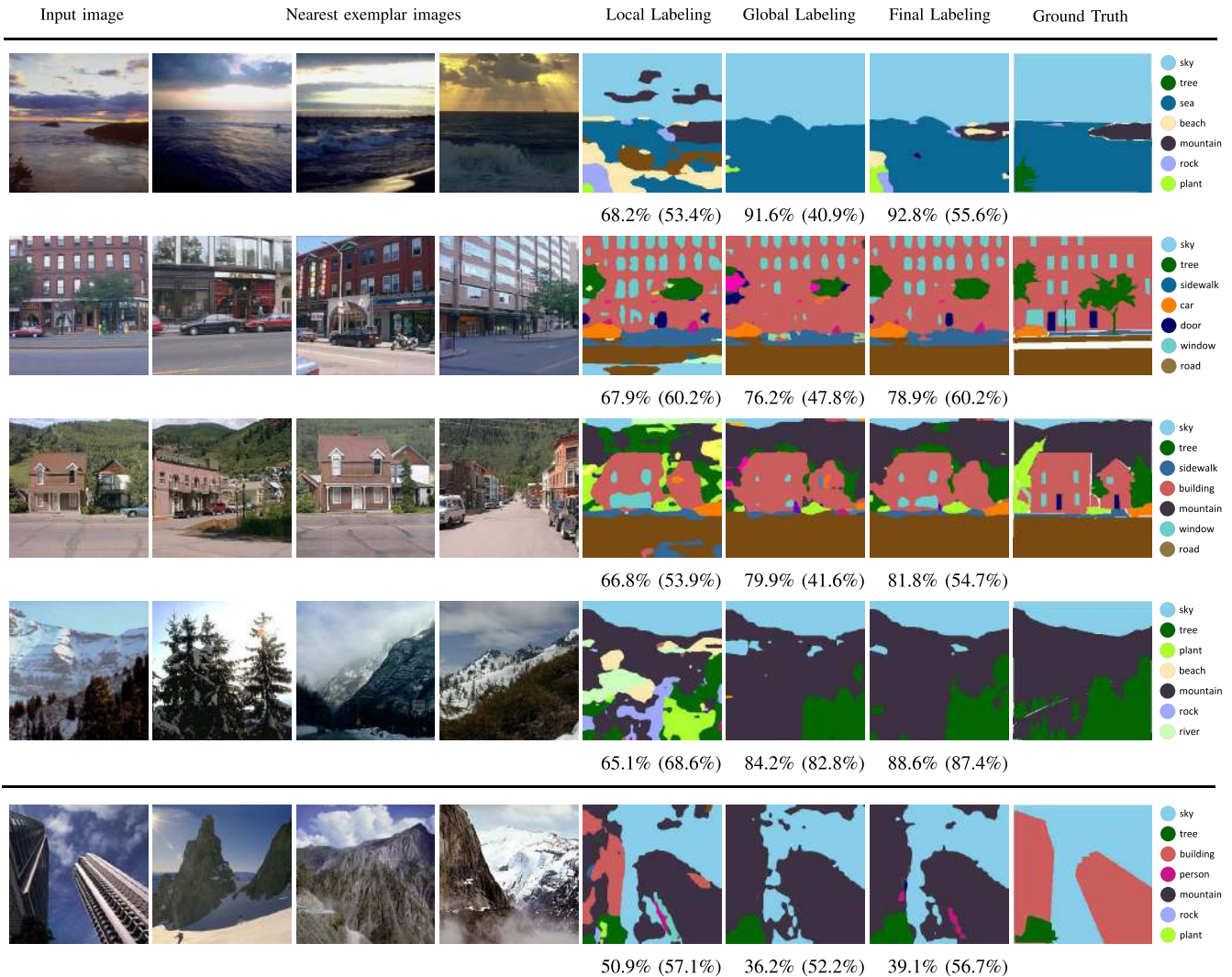


Fig. 6. Qualitative labeling examples on the SiftFlow dataset (best viewed in color). In each row, we show input images, top three nearest exemplar images, local labeling maps (output by the parametric model - CNN-Ensemble), global labeling maps (output by the non-parametric model), final label prediction maps (output by the integration model) and their ground truth maps respectively. The numbers outside and inside parentheses are global pixel accuracy (GPA) and average class accuracy (ACA) respectively. The last row shows an example, where the *local ambiguity* cannot be removed when their aggregated global features fail to reveal the true scene semantics.

global beliefs necessitates the competitive performance of our integration model. A glimpse of these properties are depicted in the qualitative labeling examples of Figure 6. Meanwhile, It is interesting to see that the global labeling outperforms local labeling (CNN-GS), which illuminates the significance of global context for local classification.<sup>5</sup>

We further investigate the per-class accuracy changes of the parametric model after it integrates with the non-parametric model. Table V shows the quantitative results on the SiftFlow dataset: the integration model boosts the accuracy significantly for frequent classes, while slightly washes away some rare “object” classes. In detail, the global belief is more helpful for classes which are more stable in positions, and large-size classes are preferred because the target classes to be included in the nearest exemplars. Overall, our integration model is able to achieve very competitive average class accuracy (ACA), and

<sup>5</sup>Under this situation, these two models are working on the same pixel features, which are generated by truncated CNN-GS.

in the same time dramatically improve the qualitative labeling results. As evidenced by Table I and IV, similar results are also observed on more challenging Barcelona datasets.

## VII. CONCLUSION

In this paper, we first present a very effective parametric model - **CNN-Ensemble** - for local classification. The CNN components in the CNN-Ensemble are trained from image patches which are very different in the form of class frequency distribution. Therefore, each CNN component learns nonidentical visual patterns, and their decision fusion gives rise to more accurate local beliefs. Then, we alleviate the notorious *local ambiguity* problem by introducing a global scene constraint, which is mathematically achieved by adding a global energy term to the labeling energy function, and it is practically estimated in a non-parametric framework. Furthermore, A large margin based CNN metric learning method is also proposed for better global belief estimation. The final

class likelihood of pixels are obtained by integrating local and global cues. The outstanding quantitative and qualitative results on the challenging SiftFlow and Barcelona datasets illuminate the effectiveness of our methods.

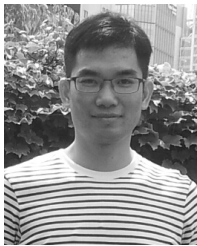
#### ACKNOWLEDGEMENTS

The authors would also like to thank NVIDIA for their generous donation of GPU.

#### REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [2] 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, Nov. 1998.
- [3] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Learning and transferring mid-level image representations using convolutional neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2014, pp. 1717–1724.
- [4] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf, "DeepFace: Closing the gap to human-level performance in face verification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2014, pp. 1701–1708.
- [5] N. Zhang, M. Paluri, M. Ranzato, T. Darrell, and L. Bourdev. (2013). "PANDA: Pose aligned networks for deep attribute modeling." [Online]. Available: <http://arxiv.org/abs/1311.5591>.
- [6] J. Tompson, A. Jain, Y. Lecun, and C. Bregler, "Joint training of a convolutional network and a graphical model for human pose estimation," in *Proc. NIPS*, 2014, pp. 1799–1807.
- [7] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2014, pp. 818–833.
- [8] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 8, pp. 1915–1929, Aug. 2013.
- [9] P. Pinheiro and R. Collobert, "Recurrent convolutional neural networks for scene labeling," in *Proc. 31st Int. Conf. Mach. Learn.*, 2014, pp. 82–90.
- [10] X. He, R. S. Zemel, and M. A. Carreira-Perpinan, "Multiscale conditional random fields for image labeling," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, vol. 2, Jun./Jul. 2004, pp. II-695–II-702.
- [11] P. Kohli, L. Ladický, and P. H. S. Torr, "Robust higher order potentials for enforcing label consistency," *Int. J. Comput. Vis.*, vol. 82, no. 3, pp. 302–324, May 2009.
- [12] J. Shotton, J. Winn, C. Rother, and A. Criminisi, "TexonBoost: Joint appearance, shape and context modeling for multi-class object recognition and segmentation," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2006, pp. 1–15.
- [13] Y. Zhang and T. Chen, "Efficient inference for fully-connected CRFs with stationarity," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2012, pp. 582–589.
- [14] C. Liu, J. Yuen, and A. Torralba, "Nonparametric scene parsing: Label transfer via dense scene alignment," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2009, pp. 1972–1979.
- [15] J. Tighe and S. Lazebnik, "SuperParsing: Scalable nonparametric image parsing with superpixels," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2010, pp. 352–365.
- [16] B. Shuai, G. Wang, Z. Zuo, B. Wang, and L. Zhao, "Integrating parametric and non-parametric models for scene labeling," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 4249–4258.
- [17] L. Wang, T. Liu, G. Wang, K. L. Chan, and Q. Yang, "Video tracking using learned hierarchical features," *IEEE Trans. Image Process.*, vol. 24, no. 4, pp. 1424–1435, Apr. 2015.
- [18] Z. Zuo, G. Wang, B. Shuai, L. Zhao, and Q. Yang, "Exemplar based deep discriminative and shareable feature learning for scene image classification," *Pattern Recognit.*, vol. 48, no. 10, pp. 3004–3015, 2015.
- [19] S. Bulò and P. Kotschieder, "Neural decision forests for semantic image labelling," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2014, pp. 81–88.
- [20] M. Mostajabi and I. Gholampour, "A robust multilevel segment description for multi-class object recognition," *Mach. Vis. Appl.*, vol. 26, no. 1, pp. 15–30, 2015.
- [21] M. Mostajabi, P. Yadollahpour, and G. Shakhnarovich, "Feedforward semantic segmentation with zoom-out features," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 3376–3385.
- [22] A. Roy and S. Todorovic, "Scene labeling using beam search under mutex constraints," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2014, pp. 1178–1185.
- [23] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. (2014). "Semantic image segmentation with deep convolutional nets and fully connected CRFs." [Online]. Available: <http://arxiv.org/abs/1412.7062>.
- [24] P. Kotschieder, S. R. Bulo, H. Bischof, and M. Pelillo, "Structured class-labels in random forests for semantic image labelling," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Nov. 2011, pp. 2190–2197.
- [25] P. Márquez-Neila, P. Kohli, C. Rother, and L. Baumela, "Non-parametric higher-order random fields for image segmentation," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2014, pp. 269–284.
- [26] S. Zheng *et al.*, "Conditional random fields as recurrent neural networks," in *Proc. IEEE Conf. Comput. Vis.*, Dec. 2015, pp. 1529–1537.
- [27] B. Shuai, Z. Zuo, and G. Wang, "Quaddirectional 2D-recurrent neural networks for image labeling," *IEEE Signal Process. Lett.*, vol. 22, no. 11, pp. 1990–1994, Nov. 2015.
- [28] B. Shuai, Z. Zuo, G. Wang, and B. Wang. (2015). "DAG-recurrent neural networks for scene labeling." [Online]. Available: <http://arxiv.org/abs/1509.00552>.
- [29] D. Eigen and R. Fergus, "Nonparametric image parsing using adaptive neighbor sets," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2012, pp. 2799–2806.
- [30] G. Singh and J. Kosecka, "Nonparametric scene parsing with adaptive feature relevance and semantic context," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 3151–3157.
- [31] F. Tung and J. J. Little, "CollageParsing: Nonparametric scene parsing by adaptive overlapping windows," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2014, pp. 511–525.
- [32] J. Yang, B. Price, S. Cohen, and M.-H. Yang, "Context driven scene parsing with attention to rare classes," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2014, pp. 3294–3301.
- [33] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, 2004.
- [34] S. Gould and Y. Zhang, "PatchMatchGraph: Building a graph of dense patch correspondences for label transfer," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2012, pp. 439–452.
- [35] S. Gould, J. Zhao, X. He, and Y. Zhang, "Superpixel graph label transfer with learned distance metric," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2014, pp. 632–647.
- [36] Z.-H. Zhou, *Ensemble Methods: Foundations and Algorithms*. Boca Raton, FL, USA: CRC Press, 2012.
- [37] J. Shotton, T. Sharp, P. Kohli, S. Nowozin, J. Winn, and A. Criminisi, "Decision jungles: Compact and rich models for classification," in *Proc. Adv. Neural Inf. Process. Syst.*, 2013, pp. 234–242.
- [38] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.
- [39] A. Criminisi and J. Shotton, *Decision Forests for Computer Vision and Medical Image Analysis* (Advances in Computer Vision and Pattern Recognition). London, U.K.: Springer, 2013.
- [40] L. Breiman, "Bagging predictors," *Mach. Learn.*, vol. 24, no. 2, pp. 123–140, 1996.
- [41] Y. Freund, R. Schapire, and N. Abe, "A short introduction to boosting," *J.-Jpn. Soc. Artif. Intell.*, vol. 14, nos. 771–780, p. 1612, 1999.
- [42] K. Simonyan and A. Zisserman. (2014). "Very deep convolutional networks for large-scale image recognition." [Online]. Available: <http://arxiv.org/abs/1409.1556>.
- [43] C. Szegedy *et al.* (2014). "Going deeper with convolutions." [Online]. Available: <http://arxiv.org/abs/1409.4842>.
- [44] O. Russakovsky *et al.*, "ImageNet large scale visual recognition challenge," *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, 2015.
- [45] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 2, Jun. 2006, pp. 2169–2178.
- [46] J. Hu, J. Lu, and Y.-P. Tan, "Discriminative deep metric learning for face verification in the wild," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2014, pp. 1875–1882.
- [47] A. Vedaldi and K. Lenc, "MatConvNet: Convolutional neural networks for MATLAB," in *Proc. 23rd Annu. ACM Conf. Multimedia Conf.*, Oct. 2015, pp. 689–692.

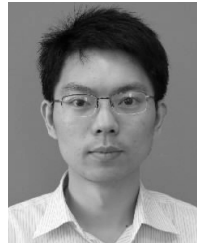
- [48] A. Oliva and A. Torralba, "Building the gist of a scene: The role of global image features in recognition," *Prog. Brain Res.*, vol. 155, pp. 23–36, Oct. 2006.
- [49] J. Tighe and S. Lazebnik, "Finding things: Image parsing with regions and per-exemplar detectors," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 3001–3008.
- [50] F. Tung and J. J. Little, "Scene parsing by nonparametric label transfer of content-adaptive windows," *Comput. Vis. Image Understand.*, vol. 143, pp. 191–200, Feb. 2015.
- [51] M. George, "Image parsing with a wide range of classes and scene-level context," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2015, pp. 3622–3630.
- [52] C. Gatta, A. Romero, and J. van de Veijer, "Unrolling loopy top-down semantic feedback in convolutional deep networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2014, pp. 504–511.
- [53] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 3431–3440.
- [54] F. Perronnin, Y. Liu, J. Sánchez, and H. Poirier, "Large-scale image retrieval with compressed Fisher vectors," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2010, pp. 3384–3391.



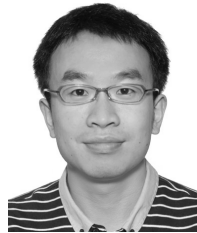
**Bing Shuai** received the B.E. degree from Chongqing University, China, in 2010, and the M.S. degree from Xiamen University, China, in 2013. He is currently pursuing the Ph.D. degree with the Rapid-Rich Object Search Laboratory, School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore. His research interests include computer vision and machine learning.



**Zhen Zuo** received the B.S. degree from the Huazhong University of Science and Technology, China, in 2011. She is currently pursuing the Ph.D. degree with the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore. Her research interests include computer vision and machine learning.



**Gang Wang** received the B.S. degree in electrical engineering from the Harbin Institute of Technology, in 2005, and the Ph.D. degree in electrical and computer engineering from the University of Illinois at Urbana-Champaign (UIUC), in 2010. He is currently an Assistant Professor with the School of Electrical and Electronic Engineering, Nanyang Technological University (NTU), and a Research Scientist with the Advanced Digital Science Center. His research interests include computer vision and machine learning. Particularly, he is focusing on object recognition, scene analysis, large scale machine learning, and deep learning. During his Ph.D. studies, he was a recipient of the prestigious Harriett and Robert Perry Fellowship (2009-2010) and the CS/AI Award (2009) at UIUC. His research is currently sponsored by NTU, the Ministry of Education, the Media Development Authority, SPRING Singapore, and the Agency for Science Technology and Research. Beyond pursuing scientific impact, he is also interested in solving real-world problems and commercializing research technologies. He serves as a Founder and Consultant for several (startup) companies.



**Bing Wang** received the B.E. degree from Anhui University, China, in 2011, and the M.E. degree from Nanyang Technological University, Singapore, in 2012. He is currently pursuing the Ph.D. degree with the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore. His research interests include computer vision, pattern recognition, and machine learning.