

# Semantic Image Segmentation via Deep Parsing Network

Ziwei Liu\* Xiaoxiao Li\* Ping Luo Chen Change Loy Xiaoou Tang  
 Department of Information Engineering, The Chinese University of Hong Kong  
 {lz013, lx015, pluo, ccloy, xtang}@ie.cuhk.edu.hk

## Abstract

Con

wise functions including long range dependencies, high order potentials and semantic label contexts. For example, Krahenbuehler et al. attained accurate segmentation boundary by inferring on a fully connected graph. Ineet et al. extended by defining both high order and long range terms between pixels. Globa et al. investigate semantic contexts between labels were also investigated by. Although they accomplished promising results, they modeled the unary terms as SVM or Adaboost, whose learning capacity becomes a bottleneck. The learning and inference of complex pairwise terms are often expensive.

In the second group, people earned a strong unary classifier by leveraging the recent advances of deep learning, such as the Convolutional Neural Network (CNN). With deep models, these works demonstrated encouraging results using simple definition of the pairwise function or even ignore it. For instance, Long et al. transferred fully connected layers of CNN into convolutional layers, allowing accurate per pixel classification possible using the contemporary CNN architectures that were pre-trained on ImageNet. Chen et al. improved by feeding the outputs of CNN into a MRF with simple pairwise potentials, but it treated CNN and MRF as separated components. A recent advance was obtained by, which jointly trained CNN and MRF by passing the error of MRF inference backward into CNN, but iterative inference of MRF such as the mean field algorithm, MF, is required for each training image during backpropagation (BP). Zheng et al. further showed that the procedure of MF inference can be represented as a Recurrent Neural Network (RNN), but their computational costs are still high. We found that directly combining CNN and MRF as above is inefficient, because CNN typically has millions of parameters while MRF infers thousands of latent variables, and even worse, incorporating complex pairwise terms into MRF becomes impractical, limiting the performance of the entire system.

## 1. Introduction

Markov Random Field (MRF) or Conditional Random Field (CRF) has achieved great successes in semantic image segmentation, which is one of the most challenging problems in computer vision. Existing works such as can be generally categorized into two groups based on their definitions of the unary and pairwise terms of MRF.

In the first group, researchers improved labeling accuracy by exploring rich information to define the pair

\*indicates shared first authorship.

Project page: <http://personal.ie.cuhk.edu.hk/~lz013/projects/DPN.html>. For more technical details, please contact the corresponding author Ping Luo (pia.uo.hi.g.ai.co).

This work proposes a novel Deep Parsing Network (DPN), which is able to jointly train CNN and complex pairwise terms. DPN has several appealing properties.

DPN solves MRF with a single feed forward pass, reducing computational cost and meanwhile maintaining high performance. Specifically DPN encodes unary terms by extending the GG16 network, GG16 pretrained on ImageNet, while additional layers are carefully designed to encode complex pairwise terms. Learning of these terms is transferred into deterministic end-to-end computation by BP instead of embedding MF into BP as we did. Although MF can be represented by RNN, it needs to recurrently compute the forward pass so as to achieve good performance and thus is time-consuming. Each forward pass contains hundred thousands of weights. DPN approximates MF by using only one iteration. This is made possible by jointly learning strong unary terms and rich pairwise information. Pairwise terms determine the graphical structure. In previous works, if the former is changed, so is the latter as well as its inference procedure. But with DPN, modifying the complexity of pairwise terms over a range of pixels and contexts is as simple as modifying the receptive fields of convolutions, without varying BP. DPN is able to represent multiple types of pairwise terms, including any previous works as its special cases. DPN approximates MF with convolutional and pooling operations, which can be speeded up by low-rank approximation and easily parallelized in a Graphics Processing Unit (GPU).

Our contributions are summarized as below. A novel DPN is proposed to jointly train GG16 and rich pairwise information. Compared to existing deep models, DPN can approximate MF with only one iteration, reducing computational cost but still maintaining high performance. We discuss that DPN represents multiple types of MRFs, including any previous works such as RNN and DeepLab as its special cases. Extensive experiments investigate which component of DPN is crucial to achieve high performance. A single DPN model achieves a new state-of-the-art accuracy of ... on the PASCAL3D+ test set. We analyze the time complexity of DPN on GPU.

## 2. Our Approach

DPN learns MRF by extending GG16 to encode unary terms and additional layers are carefully designed for pairwise terms.

**Overview** MRF is an undirected graph where each node represents a pixel in an image  $I$  and each edge represents relation between pixels. Each node is associated with a binary latent variable  $y_u^i \in \{0, 1\}$  indicating whether a pixel  $i$  has label  $u$ . We have  $\forall u \in L = \{1, 2, \dots, l\}$  representing a set of  $l$  labels. The energy function of MRF

is written as

$$E(y) = \sum_{i \in V} \psi_i(y_i^u) + \sum_{(i,j) \in E} \psi_{i,j}(y_i^u, y_j^v),$$

where  $V$  and  $E$  denote a set of latent variables nodes and edges respectively.  $\psi_i(y_i^u)$  is the unary term measuring the cost of assigning label  $u$  to the  $i$ th pixel. For instance, if pixel  $i$  belongs to the first category other than the second one, we should have  $\psi_i^1 < \psi_i^2$ . Moreover,  $\psi_{i,j}(y_i^u, y_j^v)$  is the pairwise term that measures the penalty of assigning labels  $u, v$  to pixels  $i, j$  respectively.

Intuitively the unary terms represent per pixel class assignments, while the pairwise terms represent a set of smoothness constraints. The unary term in Eqn. 1 is typically defined as

$$\psi_i(y_i^u) = -\ln p(y_i^u = 1 | I)$$

where  $p(y_i^u = 1 | I)$  indicates the probability of the presence of label  $u$  at pixel  $i$ , denoted by GG16. To simplify discussions, we abbreviate it as  $p_i^u$ . The smoothness term can be formulated as

$$\psi_{i,j}(y_i^u, y_j^v) = (u, v) d(i, j),$$

where the first term earns the penalty of global co-occurrence between any pair of labels, the output  $(u, v)$  of  $(u, v)$  is large if  $u$  and  $v$  should not coexist, while the second term captures the distances between pixels.  $d(i, j) = \omega_1 \|I_i - I_j\|^2 + \omega_2 \|[x_i, y_i] - [x_j, y_j]\|^2$ . Here  $I_i$  indicates a feature vector such as RGB features extracted from the  $i$ th pixel,  $x, y$  denote coordinates of pixel positions, and  $\omega_1, \omega_2$  are the constant weights. Eqn. 2 implies that if two pixels are close and co-similar, they are encouraged to have labels that are compatible. It has been adopted by most of the recent deep models for semantic image segmentation.

However, Eqn. 2 has two main drawbacks. First, its first term captures the co-occurrence frequency of two labels in the training data, but neglects the spatial context between objects. For example, person may appear beside table, but not at its bottom. This spatial context is a mixture of patterns, as different object configurations may appear in different images. Second, it defines only the pairwise relations between pixels, missing their high-order interactions.

To resolve these issues, we define the smoothness term by leveraging rich information between pixels, which is one of the advantages of DPN over existing deep models. We have

$$\psi_{i,j}(y_i^u, y_j^v) = \sum_{k=1}^K \lambda_k \sum_{z \in N_j} \psi_k(i, u, j, v, z) d(j, z) p_z^v.$$

The first term in Eqn. 3 earns a mixture of local label contexts, penalizing label assignment in a local region.

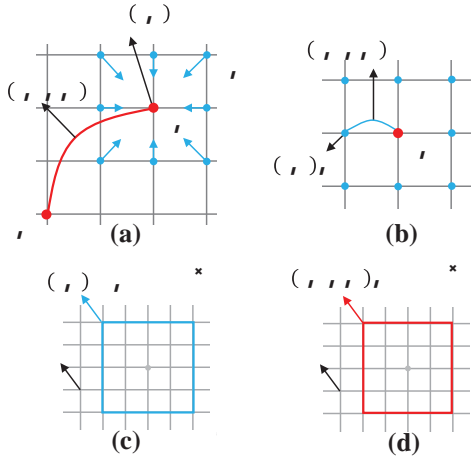


Figure 1. a) Illustration of the pairwise terms in DPN. b) explains the above contexts. c) and d) show that each +ve update of DPN corresponds to combinations.

where  $K$  is the number of components in mixture and  $\lambda_k$  is an indicator determining which component is activated. We define  $\lambda_k \in \{0, 1\}$  and  $\sum_{k=1}^K \lambda_k = 1$ . An intuitive illustration is given in Fig. 1, b, where the dots in red and blue represent a center pixel  $i$  and its neighboring pixels  $j$ ,  $j \in \mathcal{N}_i$  and  $(i, u)$  indicates assigning label  $u$  to pixel  $i$ . Here  $(i, u, j, v)$  outputs labeling cost between  $(i, u)$  and  $(j, v)$  with respect to their relative positions. For instance, if  $u, v$  represent person and table, the earned penalties of positions  $j$  that are at the bottom of center  $i$  should be large. The second term basically codes a **triple penalty** which involves pixels  $i, j$  and  $j$ 's neighbors implying that if  $(i, u)$  and  $(j, v)$  are compatible then  $(i, u)$  should be also compatible with  $j$ 's nearby pixels  $(z, v), \forall z \in \mathcal{N}_j$  as shown in Fig. 1, a.

Learning parameters, weights of GG16 and costs of above contexts in Eqn. (1) is to initialize the distances between ground truth label map and  $\mathbf{y}$ , which needs to be inferred subject to the smoothness constraints.

**Inference Overview** Inference of Eqn. (1) can be obtained by the standard MF algorithm, which estimates the joint distribution of MRE  $P(\mathbf{y}) = \frac{1}{Z} \exp\{-E(\mathbf{y})\}$ , by using a fully factorized proposal distribution,  $Q(\mathbf{y}) = \prod_{i \in \mathbf{V}} \prod_{u \in \mathbf{L}} q_i^u$ , where each  $q_i^u$  is a variable we need to estimate, indicating the predicted probability of assigning label  $u$  to pixel  $i$ . To simplify the discussion, we denote  $(y_i^u)$  and  $(y_i^u, y_j^v)$  as  $y_i^u$  and  $y_{ij}^{uv}$  respectively.  $Q(\mathbf{y})$  is typically optimized by initializing a free energy function (4) of MRE

$$F(Q) = \sum_{i \in \mathbf{V}} \sum_{u \in \mathbf{L}} q_i^u \log q_i^u + \sum_{i, j \in \mathbf{V}} \sum_{u \in \mathbf{L}} \sum_{v \in \mathbf{L}} q_i^u q_j^v y_{ij}^{uv} + \sum_{i \in \mathbf{V}} \sum_{u \in \mathbf{L}} q_i^u \log q_i^u.$$

Specifically, the first term in Eqn. (4) characterizes the cost of each pixel's predictions, while the second term characterizes the consistencies of predictions between pixels. The last term is the entropy, measuring the confidences of predictions. To estimate  $q_i^u$ , we differentiate Eqn. (4) with respect to it and equate the resulting expression to zero. We then have a closed-form expression

$$q_i^u \propto \exp \left\{ - \left( \sum_{j \in \mathbf{N}_i} \sum_{v \in \mathbf{L}} q_j^v y_{ij}^{uv} \right) \right\},$$

such that the predictions for each pixel is independently attained by repeating Eqn. (5), which implies whether pixel  $i$  has label  $u$  is proportional to the estimated probabilities of all its neighboring pixels weighted by their corresponding smoothness penalties. Substituting Eqn. (5) into (4), we have

$$q_i^u \propto \exp \left\{ - \sum_{k=1}^K \lambda_k \sum_{v \in \mathbf{L}} \sum_{j \in \mathbf{N}_i} d(j, z) q_j^v q_z^v \right\},$$

where each  $q_i^u$  is initialized by the corresponding  $p_i^u$  in Eqn. (1), which is the unary prediction of GG16. Eqn. (6) satisfies the smoothness constraints.

In the following, DPN approximates one iteration of Eqn. (5) by decomposing it into two steps. Let  $Q^v$  be a predicted label map of the  $v$ th category. In the first step as shown in Fig. 1, c, we calculate the triple penalty term in (5) by applying a  $m \times m$  filter on each position  $j$ , where each element of this filter equals  $d(j, z) q_j^v$ , resulting in  $Q^v$ . Apparently, this step smooths the prediction of pixel  $j$  with respect to the distances between it and its neighborhood. In the second step as illustrated in Fig. 1, d, the labeling contexts can be obtained by combining  $Q^v$  with a  $n \times n$  filter, each element of which equals  $y_{k(i, u, j, v)}$ , penalizing the triple relations as shown in Fig. 1, a.

### 3. Deep Parsing Network

This section describes the implementation of Eqn. (6) in a Deep Parsing Network (DPN). DPN extends GG16 as unary term and additional layers are designed to approximate one iteration of MF inference as the pairwise term. The hyperparameters of GG16 and DPN are compared in Table (1).

**VGG16** As listed in Table (1), a, the first row represents the number of layer and  $\beta$  in the second row represents the number of the receptive field and the number of combination, respectively. For instance, (1, 1) in the combination layer implies that the receptive field of each filter is  $1 \times 1$ , and it is applied on every single pixel of an input feature map, while  $e_j$  in the max pooling layer indicates each feature



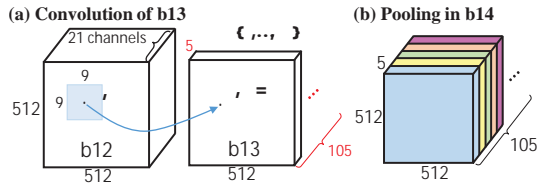


Figure 1. (a) and (b) illustrates the convolutions of  $b_{13}$  and the poolings in  $b_{14}$ .

can be formulated as

$$b_{(j,v)}^{12} = \text{lin}(k_{(j,v)} * b_{(j,v)}^{11}),$$

where  $\text{lin}(x) = ax + b$  representing the linear activation function,  $*$  is the convolutional operator and  $k_{(j,v)}$  is a  $5 \times 5$  kernel at position  $j$  of channel  $v$ . We have  $k_{(j,1)} = k_{(j,2)} = \dots = k_{(j,21)}$  shared across all channels.  $b_{(j,v)}^{11}$  indicates a local patch in  $b_{(j,v)}^{10}$ , while  $b_{(j,v)}^{12}$  is the corresponding output of  $b_{(j,v)}^{11}$ . Since  $b_{(j,v)}^{11}$  has stride one, the result of  $k_j * b_{(j,v)}^{11}$  is a scalar. In unary  $b_{(j,v)}^{11}$  has different kernels and produces 5 output feature maps.

Eqn. 11 represents the triple penalty of Eqn. 10. Recall that each output feature map of  $b_{(j,v)}^{11}$  indicates a probabilistic appearance of a specific object appearing in the image. As a result, Eqn. 11 suggests that the probability of object  $v$  presented at position  $j$  is updated by weighted averaging over the probabilities at its nearby positions. Thus, as shown in Fig. 1c,  $b_{(j,v)}^{11}$  corresponds to a patch of  $Q^v$  centered at  $j$  which has values  $p_z^v \forall z \in \mathcal{N}_j^{50 \times 50}$ . Similarly,  $k_{(j,v)}$  is initialized by  $d(j, z)p_z^v$  implying each kernel captures dissimilarities between positions. These kernels are fixed during BP other than learned as in conventional CNN.

• **b13** As shown in Table 1b and Fig. 1a,  $b_{13}$  is a convolutional layer that generates 5 feature maps by using 5 kernels of size  $5 \times 5$ . For example, the value of  $(i, u = 1)$  is attained by applying a  $5 \times 5$  kernel at positions  $\{(j, v = 1, \dots, 21)\}$ . In other words,  $b_{13}$  earns a kernel for each category to penalize the probabilistic appearance of  $b_{12}$  corresponding to the local label contexts in Eqn. 11 by assuming  $K = 5$  and  $n = 9$  as shown in Fig. 1d.

• **b14** As illustrated in Table 1 and Fig. 1b,  $b_{14}$  is a block in pooling layer that pools over every  $5 \times 5$  region with one stride across every input channels, leading to 5 output channels.  $b_{14}$  activates the contextual pattern with the smallest penalty.

Each kernel in  $b_{14}$  actually represents a distance metric between pixels in a specific region. In OC, the patterns of a the training images in a specific region are heterogeneous because of various object shapes. Therefore, we initialize each kernel with Euclidean distance. Nevertheless, Eqn. 11 is a more general form than the triple penalty in Eqn. 10. Kernels in  $b_{14}$  can be automatically learned from data if the patterns in a specific region are homogeneous such as face or human images which have more regular shapes than images in OC.

• **b15** This layer combines both the unary and smoothness terms by summing the outputs of  $b_{13}$  and  $b_{14}$  in an element-wise manner similar to Eqn. 12.

$$b_{(i,u)}^{15} = \frac{\exp(\ln(b_{(i,u)}^{11}) - b_{(i,u)}^{14})}{\sum_{u=1}^{21} \exp(\ln(b_{(i,u)}^{11}) - b_{(i,u)}^{14})},$$

where probability of assigning label  $u$  to pixel  $i$  is normalized over all the labels.

**Relation to Previous Deep Models** Many existing deep models such as AlexNet employed Eqn. 11 as the pairwise terms which are the special cases of Eqn. 10. To see this, let  $K=1$  and  $j=i$  the right hand side of Eqn. 10 reduces to

$$\begin{aligned} & \exp\left\{-\frac{u}{i} - \lambda_1 \frac{1}{v} \frac{d(i, z)p_i^v p_z^v}{L} \right\} \\ & = \exp\left\{-\frac{u}{i} - \frac{(u, v)}{v} \frac{d(i, z)p_z^v}{L} \right\}, \end{aligned}$$

where  $(u, v)$  and  $d(i, z)$  represent the global label co-occurrence and pairwise pixel similarity of Eqn. 10 respectively. This is because  $\lambda_1$  is a constant,  $d(i, i) = 0$  and  $(i, u, i, v) = (u, v)$ . Eqn. 11 is the corresponding MF update equation of Eqn. 10.

### 3.3. Learning Algorithms

**Learning** The first ten groups of DPN are initialized by GG16, while the last four groups can be initialized randomly. DPN is then fine-tuned in an incremental manner with four stages. During fine-tuning, at these stages so far the pixel-wise softmax loss is used, but updating different sets of parameters.

First, we add a loss function to  $b_{10}$  and fine-tune the weights from  $b_{10}$  to  $b_{11}$  without the last four groups in order to learn the unary terms. Second, to learn the triple relations, we stack  $b_{13}$  on top of  $b_{11}$  and update its parameters.  $\omega_1, \omega_2$  in the distance measure, but the weights of the preceding groups,  $b_{10} \sim b_{11}$  are fixed. Third,  $b_{14}$  and  $b_{15}$  are stacked onto  $b_{13}$  and similarly their weights are updated with all the preceding parameters fixed so as to learn the local label contexts. Finally, all the parameters are jointly fine-tuned.

**Implementation** DPN transforms Eqn. 10 into convolutions and poolings in the groups from  $b_{10}$  to  $b_{15}$ , such that filtering at each pixel can be performed in a parallel manner. Assume we have  $f$  input and  $f$  output feature maps  $N \times N$  pixels kernels with  $s \times s$  receptive field and a mini batch with  $M$  samples.  $b_{10}$  takes a total  $f \cdot N^2 \cdot s^2 \cdot M$  operations,  $b_{11}$  takes  $f \cdot f \cdot N^2 \cdot s^2 \cdot M$  operations, while both  $b_{13}$  and  $b_{14}$  require  $f \cdot N^2 \cdot M$  operations.

We use the released GG16 code, which is publicly available at [http://www.robots.ox.ac.uk/~vgg/research/very\\_deep/](http://www.robots.ox.ac.uk/~vgg/research/very_deep/)



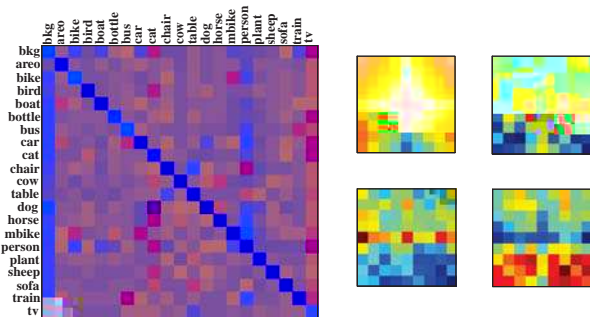


Figure 4: Visualization of (a) learned compatibility, (b) learned contextual information. (Best viewed in color)

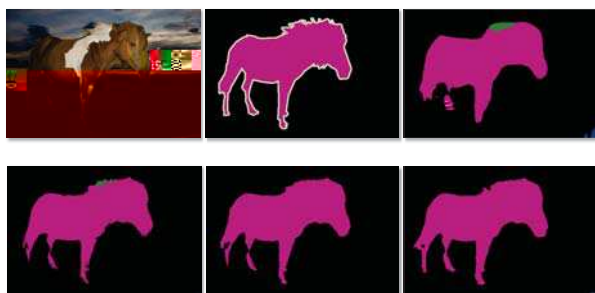


Figure 5: Step by step visualization of DPN. (Best viewed in color)

how likely a common object would present when a row object is presented. Blue represents high possibility, and red represents low possibility. For example, when horse is presented, person is more likely to present than the other objects. A sofa chair is compatible with table and bag is compatible with all the objects. (b) visualizes some contextual patterns, where  $A \rightarrow B$  indicates that when A is presented where B is more likely to present. For example, bag is around train, motorbike is below person, and person is sitting on chair.

**Incremental Learning** As discussed in Sec. 4.1, DPN is trained in an incremental manner. The right hand side of Table 1 demonstrates that each stage leads to performance gain compared to its previous stage. For instance, triple penalty improves unary term by 2.1 percent, while pairwise contexts improve triple penalty by 1.2 percent. More importantly, joint fine-tuning all the components, unary terms and pairwise terms in DPN achieves another gain of 1.2 percent. A step by step visualization is provided in Fig. 5.

We also compare incremental learning with joint learning, which fine-tunes all the components of DPN at the same time. The training curves of the two are plotted in Fig. 6(a), showing that the former leads to higher and more stable accuracies with respect to different iterations, while the latter may get stuck at local minima. This difference

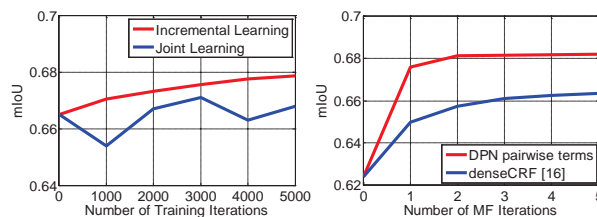


Figure 6: Ablation study of (a) training strategy, (b) required MF iterations. (Best viewed in color)

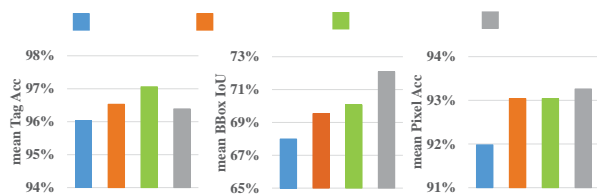


Figure 7: Stage wise analysis of (a) mean tagging accuracy, (b) mean localization accuracy, (c) mean boundary accuracy.

is easy to understand because incremental learning only introduces new parameters until all existing parameters have been fine-tuned.

**One-iteration MF DPN** approximates one iteration of MF. Fig. 6(b) illustrates that DPN reaches a good accuracy with one MF iteration. A CRF [16] with dense pairwise edges needs more than 4 iterations to converge. It also has a large gap compared to DPN. Note that the existing deep models such as [23, 24] required 100 iterations to converge as well.

**Different Components Modeling Different Information** We further evaluate DPN using three metrics. The results are given in Fig. 7. For example, (a) illustrates that the tagging accuracy can be improved in the third stage, as it captures pairwise co-occurrence with a mixture of contextual patterns. However, TA is decreased a little after the final stage. Since joint tuning maximizes segmentation accuracies by optimizing all components together, extremely small objects, which rarely occur in COCO training set, are discarded. As shown in (b), accuracies of object localization are significantly improved in the second and the final stages. This is intuitive because the unary prediction can be refined by long range and high order pixel relations, and joint training further improves results. (c) discloses that the second stage also captures object boundary, since it measures dissimilarities between pixels.

**Per-class Analysis** Table 1(a) reports the per-class accuracies of four evaluation metrics, where the first four rows represent the IoU of four stages, while the last three rows represent TA, LA, and BA, respectively. We have several valuable observations, which motivate future researches. (i) Joint training benefits most of the categories,

	areo	bi e	bird	boat	bott e	bus	car	cat	chair	cow	tab e	dog	horse	bi e person p ant	sheep sofa	train t†	Avg.			
Unary Ter. Io																				
[- Trip e Penalty																				
[- Labe Contexts			82.6					82.6		67.9			70.3			77.8				
[- Joint Tuning	84.8	37.5		66.3	67.5	84.2	76.4		33.8		50.4	76.8		74.9	81.1	48.3	41.8	76.6	60.4	67.8
TA, tagging Acc.																				
LA, labe																				
BA, boundary Acc.																				

a Per c ass resu ts on OG

	areo	bi e	bird	boat	bott e	bus	car	cat	chair	cow	tab e	dog	horse	bi e person p ant	sheep sofa	train t†	Io			
FCN																				
Zoo out																				
Piecewise																				
DeepLab																				
RNN																				
SSL†																				
RNN†														86.4						
BoxSup†	89.8		89.2	68.9				87.7			67.1							81.2	70.7	
DPN																				
DPN†		61.6			74.7	91.2	84.3		36.5	86.3		84.4	87.8		85.4	63.6	87.3	61.3		77.5

b Per c ass resu ts on OG The approaches pre trained on COCO are ar ed with †

Tab e Per c ass resu ts on OG

except ani a s such as bird, cat, and cow. So e instances of these categories are extre e y s a so that joint training discards the for s oother resu ts. Training DPN with pixe wise labe ap s i p i c i t y o d e s i i age e r e tags since it achieves a high averaged TA of \ . . . . Object oca ization a ways he ps. However, for the object with co p ex boundary such as bi e, its Io is ow even it can be oca ized. bi e has high LA but ow BA and Io. Fai ures of different categories ha e different factors. ith these three etrics, they can be easi y identi ed. For exa p e the fai ures of chair, tab e, and p ant are caused by the dif cu ties to accurately capture their bounding boxes and boundaries. A though, bott e and t† are a so dif cu t to oca ize, they achieve o derate Io because of their regu ar shapes. In other words, Io of bott e and t† can be signi cant y i p roved if they can be accurately oca ized.

## 4.2. Overall Performance

As shown in Tab e . b, we co pare DPN with the best perfor ing ethods on OG test set based on two settings - with and without pre training on COCO. The approaches pre trained on COCO are ar ed with †. e e r a uate DPN on se e ra sca es of the i ages and then average the resu ts fo owing . . . .

DPN outperfor s a the existing ethods that were trained on OG, but DPN needs on y one MF iteration to so e MRE other than . iterations of RNN, DeepLab, and Piecewise. By averaging the resu ts of two DPNs, we achieve . accuracy on OG without outside training data. As discussed in Sec. . MF iteration is the ost

† The resu ts of these ethods were presented in either the pub ished papers or arXiv pre prints.

co p ex step even when it is i p e e nted as con o tions. Therefore DPN at east reduces ( . ) x runti e co pared to previous wor s.

Fo owing . . . . we pre train DPN with COCO where . object categories that are a so presented in OG are selected for training. A sing e DPN† has achieved . . . .

Io on OG test set. As shown in Tab e . b, we observe that DPN† achieves best perfor ances on ore than ha f of the object c asses.

## 5. Conclusion

we proposed Deep Parsing Network (DPN) to address se antic i age seg entation, which has se e ra appealing properties. First, DPN uni es the inference and earning of unary ter and pairwise ter s in a sing e con o tiona network. No iterative inference are required during bac propagation. Second, high order re ations and ixtures of labe contexts are incorporated to its pairwise ter s o de ing a ing existing wor s se e ra as special cases. Third, DPN is bui t upon con e tiona operations of CNN, thus easy to be para e ized and speeded up.

DPN achieves state of the art perfor ance on OG, and u tip e r a uab e facts about se antic i age seg entation are re e a ed through extensive experi ents. Future directions include investigating the genera izabi ty of DPN to ore cha enging scenarios, arge nu ber of object c asses and substantia appearance sca e r variations.

## References

[1] A. Ada s, J. Bae, and M. A. Da is. Fast high di ensiona l ter ing using the per t ohedra attice. In *Co r*, or . o u e . pages . . . .

P. Arbelaez, M. Maire, C. Fowlkes, and J. Sivic. Contour detection and hierarchical image segmentation. In *CVPR*, pages 1689–1696, 2004.

L. C. Chen, G. Papandreou, I. Kostas, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICCV*, pages 3431–3438, 2015.

S. Chetverov, C. Sminchisescu, P. Andrusch, J. Cohen, J. Tran, B. Catanzaro, and E. Shechtman. Efficient priors for deep learning. In *CVPR*, pages 1037–1044, 2015.

J. Dai, K. He, and J. Sun. Boxesup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *CVPR*, pages 1718–1725, 2015.

J. Deng, J. Dong, R. Socher, L. J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009.

M. Everingham, L. Van Gool, C. Kiciak, S. John, and A. Zisserman. The pascal3d+ visual object classes challenge. In *CVPR*, pages 1–6, 2015.

C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. In *CVPR*, pages 3197–3204, 2014.

P. F. Felzenszfeld and D. P. Huttenlocher. Efficient belief propagation for early vision. In *CVPR*, pages 157–164, 2004.

T. Freese, E. C. Pasztor, and O. T. Carmichael. Learning low-level vision. In *CVPR*, pages 1–6, 2015.

B. Fu, Y. S. Ong, A. L. Yuille, and S. Soatto. Class segmentation and object localization with superpixel neighborhoods. In *CVPR*, pages 1–6, 2015.

B. Hariharan, P. Arbelaez, L. Bourdev, S. Maji, and J. Sivic. Semantic contours from inverse detectors. In *CVPR*, pages 1–6, 2015.

G. E. Hinton, O. S. Yoon, and J. Dean. Distilling the knowledge in a neural network. In *ICML*, pages 1–6, 2015.

M. Jaderberg, A. L. Yuille, and A. Zisserman. Speeding up convolutional neural networks with low-rank expansions. In *CVPR*, pages 1–6, 2015.

M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. In *ICML*, pages 1–6, 2001.

P. Krähenbühl and P. Kohler. Efficient inference in fully connected crfs with gaussian edge potentials. In *CVPR*, pages 1–6, 2015.

P. Krähenbühl and P. Kohler. Parallel learning and convergent inference for dense random fields. In *CVPR*, pages 1–6, 2015.

A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *CVPR*, pages 1–6, 2012.

G. Lin, C. Shen, I. Reid, and A. Hengele. Efficient piecewise training of deep structured models for semantic segmentation. In *CVPR*, pages 1–6, 2014.

T. Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *CVPR*, pages 1–6, 2014.

C. Liu, J. Yuen, and A. Torralba. Nonparametric scene parsing via label transfer. In *CVPR*, pages 1–6, 2015.

J. Long, E. Shelton, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3106–3114, 2015.

P. Luo, X. Tang, and X. Tang. Hierarchical face parsing via deep learning. In *CVPR*, pages 1–6, 2015.

P. Luo, X. Tang, and X. Tang. Pedestrian parsing via deep convolutional neural networks. In *CVPR*, pages 1–6, 2015.

M. Mostajabi, P. Yadoorani, and G. Shakhnarovich. Feed forward semantic segmentation with zoom-out features. In *CVPR*, pages 1–6, 2015.

R. Mottaghi, X. Chen, X. Liu, N. G. Cho, S. Lee, S. Fidler, R. Szeliski, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, pages 1–6, 2015.

M. Opper, O. S. Yoon, et al. From naive learned theory to the tap equations. In *CVPR*, pages 1–6, 2015.

G. Papandreou, L. C. Chen, K. Murphy, and A. L. Yuille. Deeply supervised learning of a dcnn for semantic image segmentation. In *CVPR*, pages 1–6, 2015.

X. Ren and J. Sivic. Learning a classification code for segmentation. In *CVPR*, pages 1–6, 2015.

A. G. Schwing and R. Szeliski. Fully connected deep structured networks. In *CVPR*, pages 1–6, 2015.

J. Shi and J. Sivic. Normalized cuts and image segmentation. In *CVPR*, pages 1–6, 2015.

K. Siyonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *CVPR*, pages 1–6, 2015.

Y. Sun, X. Tang, and X. Tang. Deep learning face representation by joint identification and classification. In *CVPR*, pages 1–6, 2015.

M. Sze, P. Kohler, and D. Hoies. Learning crfs using graph cuts. In *CVPR*, pages 1–6, 2015.

Y. Taig, M. Yang, M. Ranzato, and L. O. F. DeepFace: Closing the gap to human performance in face verification. In *CVPR*, pages 1–6, 2015.

T. Tenenbaum, G. Sheehy, J. J. Freese, and P. H. Torr. Pose-free: An efficient learned-based method for joint estimation of human pose, segmentation and depth. In *CVPR*, pages 1–6, 2015.

T. Tenenbaum, J. J. Freese, and P. H. Torr. Filter-based learned inference for random fields with higher order terms and productable spaces. In *CVPR*, pages 1–6, 2015.

J. Yang, B. Price, S. Cohen, and M. H. Yang. Context-driven scene parsing with attention to rare classes. In *CVPR*, pages 1–6, 2015.

S. Zheng, S. Jayasuman, B. Rother, P. Paredes, T. Tenenbaum, Z. Su, D. Du, C. Huang, and P. Torr. Conditional random fields as recurrent neural networks. In *CVPR*, pages 1–6, 2015.

Z. Su, D. Du, C. Huang, and P. Torr. Conditional random fields as recurrent neural networks. In *CVPR*, pages 1–6, 2015.