

Conditional Random Fields as Recurrent Neural Networks

Shuai Zheng¹ Sadeep Jayasumana¹ Bernardino Romera-Paredes¹ Vibhav Vineet²
 Zizhong Su³ Dalong Du³ Chang Huang³ Philip H. S. Torr¹
¹University of Oxford ²Stanford University ³Baidu Research

Abstract

Pixel-level labelling tasks, such as semantic segmentation and depth estimation, play a central role in image understanding. Recent approaches have attempted to harness the capabilities of deep learning techniques for image recognition to tackle pixel-level labelling tasks. One central issue in this approach is the limited capacity of deep learning techniques to delineate visual objects. To solve this problem, we introduce a new form of convolutional neural network, called CRF-RNN, which expresses a Conditional Random Field (CRF) as a Recurrent Neural Network (RNN). Our short network can be plugged in as a part of a deep Convolutional Neural Network (CNN) to obtain an end-to-end system that has desirable properties of both CNNs and CRFs. Importantly, our system fully integrates CRF modelling with CNNs, making it possible to train the whole system end-to-end with the usual back-propagation algorithm.

We apply this framework to the problem of semantic image segmentation, obtaining competitive results with the state-of-the-art without the need of introducing any post-processing method for object delineation.

1. Introduction

Low-level computer vision problems such as semantic image segmentation or depth estimation often involve assigning labels to each pixel in an image. While the feature representation used to classify individual pixels plays an important role in this task, it is equally important to consider factors such as image edges, appearance consistency and spatial consistency while assigning labels in order to obtain accurate and precise results.

Designing a strong feature representation is a key challenge in pixel-level labelling problems. Works on this topic include: TextonBoost [28], TextonForest [27], and Random Forest-based human body pose estimation for Kinect [26]. Recently, supervised deep learning approaches such as large-scale deep Convolutional Neural Networks (CNNs) have been immensely successful in many high-level com-

puter vision tasks such as image recognition [13] and object detection [8]. This motivates the successive explorations in employing CNNs for pixel-level labelling problems. The key insight is to learn a strong feature representation end-to-end for the pixel-level labelling task instead of hand-crafting features with heuristic parameter tuning. In fact, a number of approaches including Farabet *et al.* [7], FCN [19], and Zoom-out [20] have shown a significant accuracy boost by adapting state-of-the-art CNN based image classifiers to the semantic segmentation problem.

However, there are significant challenges in adapting CNNs designed for high level computer vision tasks such as object recognition to pixel-level labelling tasks. Firstly, traditional CNNs have convolutional filters with large receptive fields and hence produce coarse outputs when restructured to produce pixel-level labels [19, 4]. Presence of max-pooling layers further increases the coarseness of the output. This, for instance, can result in non-sharp boundaries and blob-like shapes in semantic segmentation tasks. Secondly, CNNs lack smoothness constraints that encourage label agreement between similar pixels and spatial and appearance consistency of the labelling output. Lack of such smoothness constraints can result in poor object delineation and small spurious regions in the segmentation output [31, 30, 14, 21].

On a separate track to the progress of deep learning techniques, probabilistic graphical models have been developed as effective methods to enhance the accuracy of pixel-level labelling tasks. Conditional Random Fields (CRFs), in particular, have observed widespread success in this area [17, 14, 12] and have become one of the most successful graphical models used in computer vision. CRF inference is able to refine weak and coarse pixel-level label predictions to produce sharp boundaries and fine-grained segmentations. Therefore, intuitively, CRF refinement can be used to overcome the drawbacks in utilizing CNNs for pixel-level labelling tasks.

In this paper, we propose an end-to-end deep learning solution for pixel-level semantic image segmentation problem, that combines the strengths of both CNNs and CRFs. More specifically, we formulate fully-connected dense CRF

inference as a Recurrent Neural Network (RNN) which can refine coarse outputs from a traditional CNN in the forward pass, while passing error differentials back to the CNN during training. Importantly, with our formulation, the whole network, which comprises a traditional CNN and an RNN for CRF inference, can be trained end-to-end utilizing the usual back propagation algorithm. We show that our method is able to achieve competitive results on challenging datasets.

Our contribution is based on the observation that filter-based approximate mean-field inference approach for fully-connected CRFs [12] relies on applying Gaussian and bilateral filters on the mean-field approximates in each iteration. We unroll an iteration of the algorithm as a stack of layers in a CNN. Unlike the standard convolutional layer in which filters are fixed after the training stage, this new layer involves an edge-preserving filter [29, 23] that depends on the original spatial and appearance information of the image. These filters have the additional advantages of requiring a small set of learning parameters, despite the filter size being potentially as big as the image. The parameters of the final combined CNN-CRF deep network are learned end-to-end, using back-propagation [18] to minimise the structured loss which captures the contextual information of the pixel-level labelling.

2. Related Work

In this section we review approaches that make use of CNN for low-level computer vision tasks, with a special focus on semantic image segmentation. In order to address the pixel-wise semantic image segmentation, a wide variety of approaches have been proposed. These approaches can be categorized into two main strategies.

The first strategy is based on utilizing separated mechanisms for both feature extraction, and image segmentation exploiting the edges of the image [2, 20]. One representative instance of this scheme is the application of an CNN for the extraction of meaningful features, and using superpixels to account for the structural pattern of the image. One such example is given in [20], where the authors first apply superpixels and then use a feature extraction process on each of them. Another example is described in Farabet et al [7], where the authors propose a parallel scheme in which the combination of both contributions, superpixels and features, is made at a latter stage. The main disadvantage of this strategy is that errors in the initial proposals (e.g: super-pixels) may lead to poor predictions, no matter how good the feature extraction process is.

The second strategy is to directly learn a nonlinear model from the images to the label map. This has been shown in Eigen & Fergus *et al.* [6], where the authors replaced the last fully connected layers of a CNN by convolutional layers in order to keep the spatial information. Other works, like

Hariharan *et al.* [9], and Long et al [19], build on the notion that top layers obtain meaningful features for object recognition, whereas low layers keep information about the structure of the image, such as edges. These works include some form of upsampling of the feature representation of the high layers so that they can be combined with the low layers features. The combination of features at different layers is done by summation. Chen *et al.* [4] combined a CNN with a fully connected conditional random fields (CRF). In contrast to Chen et al, which employs CRF inference as a standalone post-processing step disconnected from the CNN, our approach is an end-to-end learning system that could jointly learn the parameters for the pixel-wise CNN regression, deconvolution layers and fully-connected layers together.

In addition to the previous strategies, other works have exploited the commonalities between semantic image segmentation and other low-level computer vision tasks, such as depth and surface normals estimation. Ladicky *et al.* [15] recently showed that it is possible to employ a discriminative supervised approach to simultaneously predict pixel-wise labels with semantic image segmentation, depth estimation, and surface normal [16]. However, the handcrafted features, the use of superpixels, and the shallow classifiers limit the performance of their approach. Eigen *et al.* [6] propose a coarse-to-fine approach to achieve depth prediction from single images. It is based on refining the prediction of a first CNN by using another CNN. The latter takes as input both the original image, and the coarse output of the former CNN. They later extended this work, and developed a joint approach [5] that simultaneously pixel-wisely predict the depth, surface normal and semantic labels.

3. Conditional Random Fields

In this section we provide a brief overview of Conditional Random Fields (CRF) for pixel-wise labelling and introduce the notation used in the paper. A CRF, used in the context of pixel-wise label prediction, models pixels labels as random variables conditioned upon a global observation, namely, the image.

Let X_i be the random variable of the pixel i , which can take any value from a set of labels $\mathcal{L} = \{l_1, l_2, \dots, l_L\}$. Let \mathbf{X} be the vector formed by the variables X_1, X_2, \dots, X_N , where N is the number of pixels in the image. Given a graph $G = (V, E)$, where $V = \{X_1, X_2, \dots, X_N\}$, and a global observation (image) \mathbf{I} , the pair (\mathbf{I}, \mathbf{X}) can be modelled as a CRF characterized by a Gibbs distribution of the form $P(\mathbf{X} = \mathbf{x} | \mathbf{I}) = \frac{1}{Z(\mathbf{I})} \exp(-E(\mathbf{x} | \mathbf{I}))$. Here $E(\mathbf{x})$ is called the energy of the configuration $\mathbf{x} \in \mathcal{L}^N$ and $Z(\mathbf{I})$ is the partition function [17]. From now on, we drop the conditioning on \mathbf{I} in the notation for convenience.

In the fully connected pairwise CRF model of [12], the

energy of a label assignment \mathbf{x} is given by:

$$E(\mathbf{x}) = \sum_i \psi_u(x_i) + \sum_{i < j} \psi_p(x_i, x_j),$$

where the unary potentials $\psi_u(x_i)$ measure the inverse likelihood (and therefore, the cost) of the pixel i taking the label x_i , and pairwise potentials $\psi_p(x_i, x_j)$ measure the cost of assigning labels x_i, x_j to pixels i, j simultaneously. In our model, unary potentials are obtained from a CNN, which, roughly speaking, predicts labels for pixels without considering the smoothness and the consistency of the label assignments. The pairwise potential provides an image data-dependent smoothing term that encourages assigning similar labels to similar pixels. As done in [12], we model pairwise potentials as weighted Gaussians:

$$\psi_p(x_i, x_j) = \mu(x_i, x_j) \sum_{m=1}^M w^{(m)} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j),$$

where each $k^{(m)}$ for $m = 1, \dots, M$, is a Gaussian kernel applied on feature vectors. The feature vector of pixel i , \mathbf{f}_i , is derived from image features such as spatial location and RGB values [12].

Inference of the above CRF yields the most probable label assignment \mathbf{x} for the given image. Since the direct inference is intractable, a mean-field approximation to the CRF distribution is used. It consists in approximating the CRF distribution $P(\mathbf{X})$ by a simpler distribution Q , which is the product of marginals, i.e., $Q(\mathbf{X}) = \prod_i Q_i(X_i)$. The steps of the exact iterative algorithm for mean-field inference and its reformulation as an RNN are discussed next.

4. A Mean-field Iteration as a Stack of CNN Layers

We now present the main contribution of this paper. In particular, we show that the mean-field CRF inference algorithm proposed in [12] can be reformulated as a Recurrent Neural Network (RNN). To this end, we first consider individual steps of the algorithm summarized in Algorithm 3, and describe them as CNN layers. In that and in the remainder of this paper we use $U_i(l)$ to denote the negative of the unary potentials introduced in the previous section, $U_i(l) = -\psi_u(X_i = l)$, for the sake of brevity.

In order to reformulate the steps of the inference algorithm as CNN layers, it is essential to be able to calculate error differentials in each layer w.r.t. its inputs in order to be able to back propagate the error differentials to previous layers during training. We also discuss how to calculate error differentials w.r.t. the parameters in each layer, enabling their optimization through the back-propagation algorithm. Therefore, in our formulation, CRF parameters such as the weights of the Gaussian kernels can also be optimized automatically during the training phase of the full network.

Once the individual steps of the algorithm are broken down as CNN layers, the full algorithm can then be formulated as an RNN by arranging the layers in such a way that each iteration takes inputs from the previous iteration. We come back to this in Section 5 after discussing the steps of Algorithm 3 in detail below.

4.1. Initialization

In the initialization step of the algorithm, the operation $Q_i(l) \leftarrow \frac{1}{Z_i} \exp(U_i(l))$, where $Z_i = \sum_l \exp(U_i(l))$, is performed. Note that this is equivalent to applying a softmax function over the unary log-likelihood values U across all the labels at each pixel. The softmax function has been extensively used in CNN architectures before and is therefore well known in the deep learning community.

This step does not include any parameters and the error differentials received at the output of the step during back-propagation could be passed down to the unary inputs after performing usual backward pass calculations of the softmax transformation.

4.2. Message Passing

In the fully connected CRF formulation, message passing is implemented by applying M Gaussian filters on Q values. Gaussian filter coefficients are derived based on image features such as the pixel locations and RGB values, that reflect how strongly a pixel is related to other pixels. Since the CRF is fully connected, each filter’s receptive field spans the whole image, making it infeasible to use a brute-force implementation of the filters. Fortunately, several approximation techniques exist to make computation of high dimensional Gaussian filtering significantly faster. We use the Permutohedral lattice implementation [1], which can compute the filter response in $O(N)$ time, where N is the number of pixels of the image [1].

During back propagation, error derivatives with respect to the filter inputs are calculated by applying the M Gaussian filters in reverse on the error derivatives w.r.t. the filter outputs. Therefore, back-propagation through this filtering stage can also be performed in $O(N)$ time. In this work, for simplicity, we keep the bandwidth values of the filters fixed.

4.3. Linear Combination of Filter Outputs

The next step of the mean-field iteration is taking a weighted sum of the M filter outputs from the previous step. This can be viewed as usual convolution with a $1 \times 1 \times M$ filter. Since both inputs and the outputs to this step are known during back-propagation, the error derivative w.r.t. the filter weights can be computed, making it possible to automatically learn the filter weights (relative contributions from each Gaussian filter output from the previous stage). Error derivative w.r.t. the inputs can also be computed in the usual

Algorithm 1 Mean-field in fully CRFs, broken down as CNN layers.

$Q_i \leftarrow \frac{1}{Z_i} \exp(U_i(l))$ for all i ▷ Initialization
while not converged do
 $\tilde{Q}_i^{(m)}(l) \leftarrow \sum_{j \neq i} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j) Q_j(l)$ for all m ▷ Message passing
 $\check{Q}_i(l) \leftarrow \sum_m w^{(m)} \tilde{Q}_i^{(m)}(l)$ ▷ Linear combination of message passing outputs
 $\hat{Q}_i(l) \leftarrow \sum_{l' \in \mathcal{L}} \mu(l, l') \check{Q}_i(l')$ ▷ Compatibility transform
 $Q_i(l) \leftarrow U_i(l) - \hat{Q}_i(l)$ ▷ Adding unary values
 Exponentiate and normalize Q_i for all i ▷ Normalizing
end while

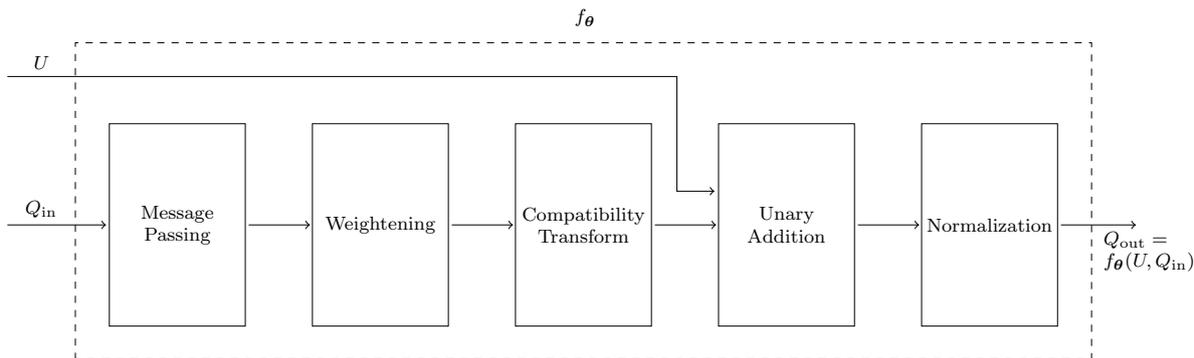


Figure 1. **A mean-field iteration as a CNN.** A single iteration of the mean-field CRF inference algorithm can be modelled as a stack of common CNN layers.

manner to pass the error derivatives down to the previous stage.

4.4. Compatibility Transform

In the compatibility transform step, outputs from the previous step (denoted by \check{Q} in Algorithm 3) are shared between the labels to a varied extent, depending on the compatibility between these labels. Compatibility between the two labels l and l' is parameterized by the function $\mu(l, l')$. The Potts model, given by $\mu(l, l') = [l \neq l']$, assigns a penalty if different labels are assigned to pixels with similar properties. A limitation of this model is that it assigns the same penalty for all different pairs of labels. Intuitively, better results can be obtained by taking the compatibility between different label pairs into account and penalizing the assignments accordingly. For example, assigning labels “person” and “bicycle” to nearby pixels should have a lesser penalty than assigning labels “sky” and “bicycle”. Therefore, learning the function μ from data is preferred to fixing it in advance with Potts model.

Compatibility transform step can be viewed as another convolution layer where the spatial receptive field of the filter is 1×1 , and the number of input and output channels are both L . Learning the weights of this filter is equivalent to learning the label compatibility function μ . Transferring error differentials from the output of this step to the input

could also be done trivially.

4.5. Adding Unary Values

In this step, the output from the compatibility transform stage is subtracted element-wise from the unary inputs U . While no parameters are involved in this step, transferring error differentials can be done trivially by copying the differentials at the output of this step to both inputs with the appropriate sign.

4.6. Normalization

Finally, the normalization step of the iteration can be considered as another softmax layer with no tunable parameters. Differentials at the output of this step can be passed on to the input using the softmax layer’s backward pass.

5. The CRF-RNN Network

In the previous section, it was shown that one iteration of the algorithm can be formulated as a stack of common CNN layers (see Fig. 1). We use the function f_θ to denote the transformation done by one mean-field iteration: given unary log-likelihood values U and an estimation Q_{in} , the next estimation after one mean-field iteration is given by $f_\theta(U, Q_{in})$. The vector θ represents the CRF parameters described in Section 4.

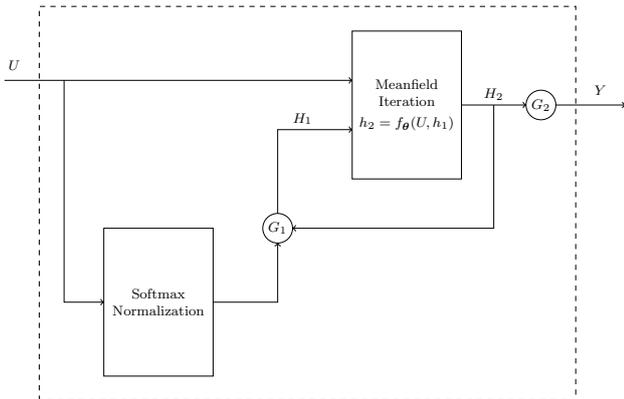


Figure 2. **The CRF-RNN Network.** We formulate the mean-field CRF inference algorithm as a Recurrent Neural Network (RNN). Gating functions G_1 and G_2 are fixed as described in the text.

Multiple mean-field iterations can be implemented by repeating the above stack of layers in such a way that each iteration takes Q value estimations from the previous iteration and the unary values in their original form. This is equivalent to treating the iterative mean-field inference as a Recurrent Neural Network (RNN) as shown in Fig. 2. Using the notation in the figure, the behaviour of the network is given by the following equations where T is the number of iterations:

$$\begin{aligned}
 H_1(t) &= \begin{cases} \text{softmax}(U), & t = 0 \\ H_2(t-1), & 0 < t \leq T, \end{cases} \\
 H_2(t) &= f_\theta(H_1(t), U), \quad 0 \leq t \leq T, \\
 Y(t) &= \begin{cases} 0, & 0 \leq t < T \\ H_2(t), & t = T. \end{cases}
 \end{aligned}$$

We name this RNN structure CRF-RNN. Parameters of the CRF-RNN are same as the mean-field parameters described in Section 4 and denoted by θ here. Since the calculation of error differentials w.r.t. these parameters in a single iteration was described in Section 4, they can be learnt in the RNN setting using the standard back-propagation through time algorithm [25, 22]. It was shown in [12] that the iterative algorithm for fully connected CRF converges very fast in less than 10 iterations. Furthermore, in practice, after 3 to 5 iterations, increasing the number of iterations usually does not significantly improve results [12]. We used 3 iterations in all our experiments. This means that our CRF-RNN does not learn way back through time. Therefore, it does not suffer from the infamous vanishing and exploding gradient problem inherent to temporarily deep RNNs [3, 24].

Let us now consider the case where the CRF-RNN is plugged in as part of a deeper neural network. For instance, this will be the case when the proposed CRF-RNN is used

to refine unary log-likelihood values U provided by a CNN. In the forward pass, once the computation enters the CRF-RNN, it takes T iterations for the data to leave the carousel created by the RNN. Neither the CNN that provides unary values nor the other layers after the CRF-RNN need to perform any computations during this time since the refinement happens only inside the RNN’s carousel. Once the output Y leaves the carousel, next stages of the deep network after the CRF-RNN can continue the forward pass. In our experiments, we use a traditional feed-forward CNN to obtain unary values and use the CRF-RNN as the last stage of the deep network.

During the backward pass, once the error differentials reach the CRF-RNN’s output Y , they similarly spend T iterations within the carousel before reaching the RNN input U in order to propagate to the CNN which provides the unary input. In each iteration inside the carousel, error differentials are computed inside each component of the mean-field iterations as described in Section 4. We note that unnecessarily increasing the number of mean-field iterations T could result in the vanishing and exploding gradient problems in the CRF-RNN.

6. Experiments

In this section, we present some preliminary experimental results with the proposed CRF-RNN framework, where we optimize the parameters of the RNN stage.

Our approach has been implemented using the Caffe [11] deep learning library. Our deep network is formed by using the FCN-8s network of [19] to provide unary potentials to the CRF-RNN. We initialize this first part of the network, which computes unary potentials, using the publicly available models of [19]. Initial values for the CRF parameters are obtained from a validation process using the validation set. These are then optimized within the network using the back-propagation algorithm. In all our experiments we set the number of mean-field iterations T in the CRF-RNN to 3. It was observed that increasing the number of iterations beyond this did not significantly improve the results.

We evaluated our approach CRF-RNN on PASCAL VOC Segmentation datasets. In the VOC-2011 Segmentation dataset, there are 1,112 images in the training set, 1,111 in the validation set, and 1,111 in the test set. The ground truth labels of the latter are not publicly available. We augmented this training set with the additional data provided by Hariharan *et al.* [10], which consists of 9,800 images.

We trained our system using this augmented VOC-2011 train set and use the PASCAL VOC evaluation server to obtain the results on VOC-2011 and 2012 test sets. As shown in Table 1 and Table 2, our CRF-RNN obtains the best performance on the VOC-2011 test set and the second place in the VOC-2012 test set. The CRF-RNN significantly improves the segmentation accuracy over the FCN-8s network,

which is used to obtain unary potentials (by 2.7 percentage points on VOC-2011 and 3.0 on VOC-2012). The best performer on VOC-2012 uses a different CNN, trained weights of which have not been made publicly available yet.

Our approach also produces qualitatively better results in PASCAL VOC 2011. For example, as evidenced by Fig. 3, the segmentation outputs of our approach have better boundaries and fewer spurious regions.

Method	Test Set	Validation Set
R-CNN [8]	47.9	-
SDS [10]	52.6	53.9
FCN-8s [19]	62.7	67.5
Zoom-out [20]	64.1	-
CRF-RNN	65.3	69.0

Table 1. Mean IU accuracy of our approach, CRF-RNN, compared to the other state-of-the-art approaches on the Pascal VOC-2011 dataset.

Method	Test Set
R-CNN [8]	47.9
SDS [10]	52.6
FCN-8s [19]	62.7
Zoom-out [20]	64.1
DeepLab [4]	67.1
CRF-RNN	65.2

Table 2. Mean IU accuracy of our approach, CRF-RNN, compared to the other state-of-the-art approaches on the Pascal VOC-2012 dataset.

7. Conclusion

We presented CRF-RNN, an interpretation of dense CRFs as Recurrent Neural Networks. Our formulation fully integrates CRF inference with emerging deep learning techniques. In particular, the proposed CRF-RNN can be plugged in as a part of a traditional deep neural network: It is capable of passing on error differentials from its outputs to inputs during back-propagation based training of the deep network while learning CRF parameters. We demonstrate the use of the proposed CRF-RNN by utilizing it for the semantic segmentation task, where it is used to perform CRF inference on unary potentials obtained from a traditional deep CNN. Introducing CRF inference to end-to-end deep learning approaches significantly improves semantic segmentation outputs. This improvement can be attributed to the uniting of the strengths of CNNs and CRFs using the proposed CRF-RNN network.

References

- [1] A. Adams, J. Baek, and M. A. Davis. Fast high-dimensional filtering using the permutohedral lattice. *Computer Graphics Forum*, 29(2):753–762, 2010.
- [2] P. Arbeláez, B. Hariharan, C. Gu, S. Gupta, L. Bourdev, and J. Malik. Semantic segmentation using regions and parts. In *IEEE CVPR*, 2012.
- [3] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 1994.
- [4] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *arXiv:1412.7062*, 2014.
- [5] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *arXiv:1411.4734*, 2014.
- [6] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. In *NIPS*, 2014.
- [7] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *IEEE TPAMI*, 2013.
- [8] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE CVPR*, 2014.
- [9] B. Hariharan, P. Arbelaez, R. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In *arXiv:1411.5752*, 2014.
- [10] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Simultaneous detection and segmentation. In *ECCV*, 2014.
- [11] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [12] P. Krähenbühl and V. Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *NIPS*, 2011.
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [14] L. Ladicky, C. Russell, P. Kohli, and P. H. Torr. Associative hierarchical crfs for object class image segmentation. In *IEEE ICCV*, 2009.
- [15] L. Ladicky, J. Shi, and M. Pollefeys. Pulling things out of perspective. In *IEEE CVPR*, 2014.
- [16] L. Ladicky, B. Zeisl, and M. Pollefeys. Discriminatively trained dense surface normal estimation. In *ECCV*, 2014.
- [17] J. D. Lafferty, A. McCallum, and F. C. N. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *ICML*, 2001.
- [18] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. 86(11):2278–2324, 1998.
- [19] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *arXiv:1411.4038*, 2014.

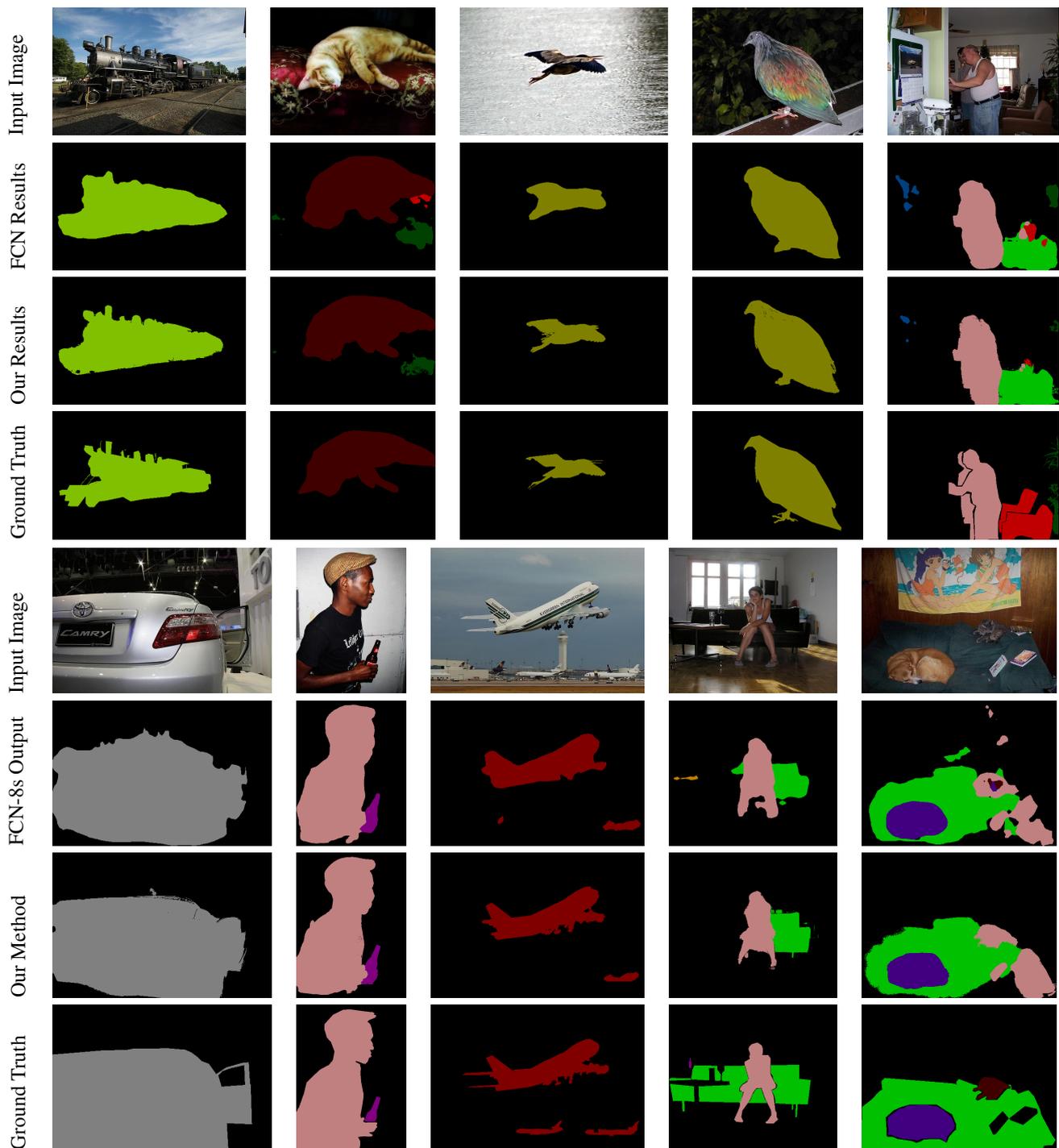


Figure 3. **Segmentation results.** Illustration of sample results on the validation set of the Pascal VOC-2011 dataset.

- [20] M. Mostajabi, P. Yadollahpour, and G. Shakhnarovich. Feed-forward semantic segmentation with zoom-out features. In *arXiv:1412.0774*, 2014.
- [21] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object

detection and semantic segmentation in the wild. In *IEEE CVPR*, 2014.

- [22] M. C. Mozer. Backpropagation. chapter A Focused Back-propagation Algorithm for Temporal Pattern Recognition. L. Erlbaum Associates Inc., 1995.

- [23] S. Paris and F. Durand. A fast approximation of the bilateral filter using a signal processing approach. 81(1):24–52, 2013.
- [24] R. Pascanu, C. Gulcehre, K. Cho, and Y. Bengio. On the difficulty of training recurrent neural networks. In *ICML*, 2013.
- [25] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Neuro-computing: Foundations of research. chapter Learning Internal Representations by Error Propagation. MIT Press, 1988.
- [26] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from single depth images. 2011.
- [27] J. Shotton, M. Johnson, and R. Cipolla. Semantic texton forests for image categorization and segmentation. 2008.
- [28] J. Shotton, J. Winn, C. Rother, and A. Criminisi. Textonboost for image understanding: Multi-class object recognition and segmentation by jointly modeling texture, layout, and context. *IJCV*, 81(1):2–23, 2009.
- [29] C. Tomasi and R. Manduchi. Bilateral filtering for gray and color images. In *IEEE CVPR*, 1998.
- [30] Z. Tu. Auto-context and its application to high-level vision tasks. In *IEEE CVPR*, 2008.
- [31] Z. Tu, X. Chen, A. L. Yuille, and S.-C. Zhu. Image parsing: Unifying segmentation, detection, and recognition. 63(2):113–140, 2005.