

## Revisiting Deep Structured Models for Pixel-Level Labeling with Gradient-Based Inference\*

Måns Larsson<sup>†</sup>, Anurag Arnab<sup>‡</sup>, Shuai Zheng<sup>‡</sup>, Philip Torr<sup>‡</sup>, and Fredrik Kahl<sup>†</sup>

**Abstract.** Semantic segmentation and other pixel-level labeling tasks have made significant progress recently due to the deep learning paradigm. Many state-of-the-art structured prediction methods also include a random field model with a hand-crafted Gaussian potential to model spatial priors and label consistencies and feature-based image conditioning. These random field models with image conditioning typically require computationally demanding filtering techniques during inference. In this paper, we present a new inference and learning framework which can learn arbitrary pairwise conditional random field (CRF) potentials. Both standard spatial and high-dimensional bilateral kernels are considered. In addition, we introduce a new type of potential function which is image-dependent like the bilateral kernel, but an order of magnitude faster to compute since only spatial convolutions are employed. It is empirically demonstrated that such learned potentials can improve segmentation accuracy and that certain label-class interactions are indeed better modeled by a non-Gaussian potential. Our framework is evaluated on several public benchmarks for semantic segmentation with improved performance compared to previous state-of-the-art CNN+CRF models.

**Key words.** deep structured models, conditional random fields, deep learning, semantic segmentation

**AMS subject classifications.** 68T45, 68R10

**DOI.** 10.1137/18M1167267

**1. Introduction.** Markov random fields (MRFs), conditional random fields (CRFs), and, more generally, probabilistic graphical models are ubiquitous tools used in a variety of domains spanning computer vision, computer graphics, and image processing [31, 9, 4]. In this paper, we focus on the application of MRFs for computer vision problems involving per-pixel labeling such as image segmentation. There are many successful approaches in this line of research, such as the interactive segmentation of [41] using graph cuts and the semantic segmentation works of [33, 44] where the parallel mean-field algorithm was applied for fast inference. Recently, convolutional neural networks (CNNs) have dominated the field in a variety of recognition tasks [26, 43, 40]. However, we observe that several leading segmentation approaches still include CRFs either as a postprocessing step [14, 17, 23, 13] or as part of the deep neural

\*Received by the editors January 30, 2018; accepted for publication (in revised form) August 27, 2018; published electronically November 6, 2018. This article is an extended version of the conference paper “A projected gradient descent method for CRF inference allowing end-to-end training of arbitrary pairwise potentials” [34].

<http://www.siam.org/journals/siims/11-4/M116726.html>

**Funding:** This work has been funded by the Swedish Research Council (grant 2016-04445), the Swedish Foundation for Strategic Research (Semantic Mapping and Visual Navigation for Smart Robots), Vinnova / FFI (Perceptron, grant 2017-01942), ERC grant ERC-2012-AdG 321162-HELIOS, EPSRC grant Seebibyte EP/M013774/1, and EPSRC/MURI grant EP/N019474/1.

<sup>†</sup>Department of Electrical Engineering, Chalmers University of Technology, SE-412 96 Gothenburg, Sweden (mans.larsson@chalmers.se, fredrik.kahl@chalmers.se).

<sup>‡</sup>Department of Engineering Science, University of Oxford, Oxford OX1 3PJ, United Kingdom (aarnab@robots.ox.ac.uk, szheng@robots.ox.ac.uk, phst@robots.ox.ac.uk).

network itself [47, 36, 3, 38, 30, 46].

We also leverage this idea of embedding inference of graphical models into a neural network. An early example of this idea was presented in [12], where the authors back-propagated through the Viterbi algorithm when designing a document recognition system. Similar to [47, 3, 6, 46], we use a recurrent neural network to unroll the iterative inference steps of a CRF. This was first used in [47] and [42] to imitate mean-field inference and to train a fully convolutional network [39, 14] along with a CRF end-to-end via back propagation. In contrast to mean-field, we do not optimize the Kullback–Leibler (KL) divergence between the true probability distribution and a fully factorized approximation. Instead, we use a gradient descent approach for the inference that directly minimizes the Gibbs energy of the random field and hence avoids the approximations of mean-field. A similar framework was recently suggested in [6] and the follow-up work [7] for multilabel classification problems in machine learning with impressive results. Moreover, [19, 2] have recently shown that one can obtain lower energies compared to mean-field inference using gradient descent–based optimization schemes.

In many works, the pairwise potentials consist of parameterized Gaussians [32, 47, 3], and it is only the parameters of this Gaussian which are learned. Our framework can learn arbitrary pairwise potentials which need not be Gaussian. In [15], a general framework for learning arbitrary potentials in deep structured models was proposed based on approximate maximum likelihood learning. One of the advantages of that framework is that data likelihood is maximized in the learning process. However, this involves approximating the partition function which is otherwise intractable. This hinders the handling of large structured output spaces like in our case.

Another approach to learning arbitrary pairwise potentials was presented in [30], which uses Gibbs sampling. Again they struggle with the difficulty of computing the partition function. In the end, only experiments on synthetic data restricted to learned 2D potentials are presented.

The authors of [36] and [13] also learn arbitrary pairwise potentials to model contextual relations between parts of the image. However, their approaches still perform postprocessing with a CRF model with parametric Gaussian potentials. In [28], a pairwise potential is learned based on sparse bilateral filtering. Applying such a filter can be regarded as one iteration in the CRF inference step. In [28], the bilateral filter is applied twice, mimicking the first two iterations of inference. Our method is not restricted to a limited number of iterations. Perhaps more importantly, we not only learn sparse high-dimensional bilateral filters but also learn arbitrary spatial filters. Such spatial 2D potentials are computationally much more efficient and easier to analyze and interpret compared to their high-dimensional counterparts. We also note that [20] proposed back-propagating through mean-field inference to learn parameters. However, this was not in the context of neural networks as in the aforementioned approaches and our work. For pixel-labeling tasks, we focus on discrete random fields. We note that learning arbitrary pairwise potentials for deep structured models with continuous valued output variables has recently been explored in [46].

A major drawback with using image-dependent dense CRFs is the relatively high computation cost. Calculating the contribution of a bilateral kernel requires a filtering operation in 5D space—something that is very computationally expensive, even utilizing sophisticated ap-

proximate filtering techniques such as the permutohedral lattice filtering technique [1]. Since the image-dependent CRF usually performs very well, especially when it comes to aligning object boundaries in segmentation tasks, it is still used for these tasks. In this paper we also propose an alternative CRF model which is also image-dependent but only requires 2D convolutions during inference. The image dependence of the model comes from an output map of the base CNN that acts as a “filter selection” map. This enables the model to, for example, use one filter representing pairwise interaction between pixel labels at semantic edges and another filter far away from semantic edges.

Previous approaches trying to find alternatives to the computation-heavy bilateral CRF include [16], where they use discriminatively trained domain transform as an edge-preserving filtering method. The authors show that the domain transform can be applied as a recurrent neural network (RNN) applied across the image across all directions. Another example is [8], where they add a final layer that performs a random graph walk across the image refining the segmentation.

In summary, our contributions are as follows.

- We present a new model for a pairwise CRF potential which is image-dependent like the bilateral kernel, but does not require high-dimensional filtering. It is based on a learned 2D filter bank which makes both inference and learning an order of magnitude faster than high-dimensional filtering approaches.
- We introduce a new optimization method for CRF inference based on gradient descent that enables end-to-end training.
- We show that our inference method supports learning pairwise kernels of arbitrary shape. The learned kernels are empirically analyzed, and it is demonstrated that in many cases non-Gaussian potentials are preferred.

Our framework has been implemented in CAFFE [29], and all source code is publicly available to facilitate further research.<sup>1</sup>

**2. CRF formulation.** Consider a conditional random field (CRF) over  $N$  discrete random variables  $\mathcal{X} = \{X_1, \dots, X_N\}$  conditioned on an observation  $\mathbf{I}$ , and let  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  be an undirected graph whose vertices are the random variables  $\{X_1, \dots, X_N\}$ . Each random variable corresponds to a pixel in the image and takes values from a predefined set of  $L$  labels  $\mathcal{L} = \{0, \dots, L-1\}$ . The pair  $(\mathcal{X}, \mathbf{I})$  is modeled as a CRF characterized by the Gibbs distribution

$$(2.1) \quad P(\mathcal{X} = \mathbf{x} | \mathbf{I}) = \frac{1}{Z(\mathbf{I})} \exp(-E(\mathbf{x} | \mathbf{I})),$$

where  $E(\mathbf{x} | \mathbf{I})$  denotes the Gibbs energy function with respect to the labeling  $\mathbf{x} \in \mathcal{L}^N$  and  $Z(\mathbf{I})$  is the partition function. To simplify notation the conditioning on  $\mathbf{I}$  will from now on be dropped. The MAP inference problem for the CRF model is equivalent to the problem of minimizing the energy  $E(\mathbf{x})$ . In this paper, we only consider energies containing unary and pairwise terms. The energy function can hence be written as

$$(2.2) \quad E(\mathbf{x}) = \sum_{i \in \mathcal{V}} \psi_i(x_i) + \sum_{(i,j) \in \mathcal{E}} \psi_{ij}(x_i, x_j),$$

---

<sup>1</sup><https://github.com/maunzzz/caffe-crfgd>

where  $\psi_i : \mathcal{L} \rightarrow \mathbb{R}$  and  $\psi_{ij} : \mathcal{L} \times \mathcal{L} \rightarrow \mathbb{R}$  are the unary and pairwise potentials, respectively. We now describe these potentials before discussing inference in section 3.

**2.1. Potentials.** The unary potential  $\psi_i(x_i)$  specifies the energy cost of assigning label  $x_i$  to pixel  $i$ . In this work we obtain our unary potentials from a CNN. Roughly speaking, the CNN outputs a probability estimate of each pixel containing each class. Denoting the output of the CNN for pixel  $i$  and class  $x_i$  as  $z_{i:x_i}$ , the unary potential is

$$(2.3) \quad \psi_i(x_i) = -w_u \log(z_{i:x_i} + \epsilon),$$

where  $w_u$  is a parameter controlling the impact of the unary potentials, and  $\epsilon$  is introduced to avoid numerical problems.

The pairwise potential  $\psi_{ij}(x_i, x_j)$  specifies the energy cost of assigning label  $x_i$  to pixel  $i$  while pixel  $j$  is assigned label  $x_j$ . Introducing pairwise terms in our model enables us to take dependencies between output variables into account. We consider two alternative types, the *combined* and the *filterbank* versions.

**2.1.1. Combined.** The *combined* version has pairwise potentials that consist of a sum of one spatial term and one bilateral term. It has the following form:

$$(2.4) \quad \psi_{ij}(x_i, x_j) = k_{x_i, x_j}^{spatial}(\mathbf{p}_i - \mathbf{p}_j) + k_{x_i, x_j}^{bilateral}(\mathbf{f}_i - \mathbf{f}_j).$$

Here  $k_{x_i, x_j}^{spatial}$  denotes a spatial kernel with compact support. Its value depends on the relative position coordinates  $\mathbf{p}_i - \mathbf{p}_j$  between pixels  $i$  and  $j$ . We do not restrict these spatial terms to any specific shape. However, we restrict the support of the potential, meaning that if pixels  $i$  and  $j$  are far apart, then the value of  $k_{x_i, x_j}^{spatial}(\mathbf{p}_i - \mathbf{p}_j)$  will be zero. We choose to use spatial kernels with compact support instead of the commonly used dense Gaussian potential since this allows the inference calculations to be performed using standard 2D convolutions. The CRFs with Gaussian potentials do not in theory have compact support, and therefore, they are often referred to as dense. However, in practice, the exponential function in the kernel drops off quickly, making the interactions between pixels that are far apart negligible.

The term  $k_{x_i, x_j}^{bilateral}$  is a bilateral kernel which depends on the feature vectors  $\mathbf{f}_i$  and  $\mathbf{f}_j$  for pixels  $i$  and  $j$ , respectively. Following several previous works on random fields, we let the vector depend on pixel coordinates  $\mathbf{p}_i$  and RGB values associated to the pixel; hence  $\mathbf{f}_i$  is a 5D vector. Note that for both the spatial and the bilateral kernels, there is one kernel for each label-to-label ( $x_i$  and  $x_j$ ) interaction to enable the model to learn differently shaped kernels for each of these interactions.

**2.1.2. Filterbank.** The pairwise potentials of the *filterbank* version have the following form:

$$(2.5) \quad \psi_{ij}(x_i, x_j) = \sum_{f=1}^F g_f(\mathbf{p}_i, I) k_{x_i, x_j, f}^{spatial}(\mathbf{p}_i - \mathbf{p}_j),$$

where  $k_{x_i, x_j, f}^{spatial}$  denotes a spatial kernel with compact support similar to the case of the *combined* version. The weights  $g_f$  depend both on the position of the pixel as well as the image  $I$ . These

weights are taken as the output of a CNN. Hence, this gives rise to an image-dependent potential, but one only needs to convolve with a bank of  $F$  2D filters to evaluate it during inference. For example, the CNN outputting the weights can learn to detect semantic edges, meaning that we would apply a different spatial filter close to a semantic edge than at the center of a semantic object. Setting the last layer of the CNN as a softmax, the features  $\mathbf{g}_f$  act as “filter selectors” deciding how the several 2D filters describing the pairwise term should be weighted for each pixel individually.

**2.2. Multilabel graph expansion and relaxation.** To be able to explain our inference method, we reformulate the original minimization of  $E(\mathbf{x})$  as a real-valued optimization problem. To facilitate a continuous relaxation of the energy minimization problem, we start off by expanding our original graph in the following manner. Each vertex in the original graph  $\mathcal{G}$  will now be represented by  $L$  vertices  $X_{i:\lambda}$ ,  $\lambda \in \mathcal{L}$ . In this way, an assignment of labels in  $\mathcal{L}$  to each variable  $X_i$  is equivalent to an assignment of boolean labels 0 or 1 to each node  $X_{i:\lambda}$ , whereby an assignment of label 1 to  $X_{i:\lambda}$  means that in the multilabel assignment,  $X_i$  receives label  $\lambda$ . To ensure that only one label is assigned to each node, an additional constraint is needed saying that, for each  $i$ , only one of  $X_{i:\lambda}$  is allowed to be labeled 1. This enables us to rewrite the energy minimization problem  $\min E(\mathbf{x})$  as the following equivalent integer program:

$$(2.6) \quad \begin{aligned} \min \quad & \sum_{i \in \mathcal{V}, \lambda \in \mathcal{L}} \psi_i(\lambda) x_{i:\lambda} + \sum_{\substack{(i,j) \in \mathcal{E} \\ \lambda, \mu \in \mathcal{L}}} \psi_{ij}(\lambda, \mu) x_{i:\lambda} x_{j:\mu} \\ \text{s.t.} \quad & x_{i:\lambda} \in \{0, 1\} \quad \forall i \in \mathcal{V}, \lambda \in \mathcal{L}, \\ & \sum_{\lambda \in \mathcal{L}} x_{i:\lambda} = 1 \quad \forall i \in \mathcal{V}. \end{aligned}$$

As a next step, we relax the integer program by allowing real values on the unit interval  $[0, 1]$  instead of booleans only. We denote the relaxed variables  $q_{i:\lambda} \in [0, 1]$ . We can now write our problem as a quadratic program:

$$(2.7) \quad \begin{aligned} \min \quad & \sum_{i \in \mathcal{V}, \lambda \in \mathcal{L}} \psi_i(\lambda) q_{i:\lambda} + \sum_{\substack{(i,j) \in \mathcal{E} \\ \lambda, \mu \in \mathcal{L}}} \psi_{ij}(\lambda, \mu) q_{i:\lambda} q_{j:\mu} \\ \text{s.t.} \quad & q_{i:\lambda} \geq 0 \quad \forall i \in \mathcal{V}, \lambda \in \mathcal{L}, \\ & \sum_{\lambda \in \mathcal{L}} q_{i:\lambda} = 1 \quad \forall i \in \mathcal{V}. \end{aligned}$$

The two constraints can be summarized as  $\mathbf{q}_i \in \Delta^L \forall i \in \mathcal{V}$ , where  $\Delta^L$  is the probability simplex and  $L$  is the number of classes. A natural question is what happens when the domain is enlarged. Somewhat surprisingly, the relaxation is tight [11].

**Proposition 2.1.** *Let  $E(\mathbf{x}^*)$  and  $E(\mathbf{q}^*)$  denote the optimal values of (2.6) and (2.7), respectively. Then,*

$$E(\mathbf{x}^*) = E(\mathbf{q}^*).$$

In the supplementary material (M116726\_01.pdf [local/web 264KB]), we show that for *any* real  $\mathbf{q}$ , one can obtain a binary  $\mathbf{x}$  such that  $E(\mathbf{x}) \leq E(\mathbf{q})$ . In particular, it will be true for  $\mathbf{x}^*$  and  $\mathbf{q}^*$ , which implies  $E(\mathbf{x}^*) = E(\mathbf{q}^*)$ . Note that the proof is constructive.

**3. MAP inference via gradient descent minimization.** To solve the program stated in (2.7), we propose an optimization scheme based on projected gradient descent; see Algorithm 3.1. It was designed with an extra condition in mind—that all operations should be differentiable to enable back-propagation during training.

---

**Algorithm 3.1.** Projected gradient descent algorithm.

---

```

Initialize  $\mathbf{q}^0$ 
for  $t$  from 0 to  $T - 1$  do
  Compute the gradient  $\nabla_{\mathbf{q}} E(\mathbf{q}^t)$ .
  Take a step in the negative direction,  $\tilde{\mathbf{q}}^{t+1} = \mathbf{q}^t - \gamma \nabla_{\mathbf{q}} E$ .
  Project  $\tilde{\mathbf{q}}^{t+1}$  to the probability simplex  $\Delta^L$ .  $\mathbf{q}^{t+1} = \text{Proj}_{\Delta^L}(\tilde{\mathbf{q}})$ .
end for
return  $\mathbf{q}^{T-1}$ 

```

---

**3.1. Gradient computations.** The gradient  $\nabla_{\mathbf{q}} E$  of the objective function  $E(\mathbf{q})$  in (2.7) has the following elements:

$$(3.1) \quad \frac{\partial E}{\partial q_{i:\lambda}} = \psi_i(\lambda) + \sum_{\substack{j:(i,j) \in \mathcal{E} \\ \mu \in \mathcal{L}}} \psi_{ij}(\lambda, \mu) q_{j:\mu}.$$

The contribution from the spatial kernel in  $\psi_{ij}$  (cf. (2.4)) can be written as

$$(3.2) \quad v_{i:\lambda}^{spatial} = \sum_{\substack{j:(i,j) \in \mathcal{E} \\ \mu \in \mathcal{L}}} k_{\lambda,\mu}^{spatial}(\mathbf{p}_i - \mathbf{p}_j) q_{j:\mu}.$$

Since the value of the kernel  $v_{i:\lambda}^{spatial}$  only depends on the relative position of pixels  $i$  and  $j$ , the contribution for all pixels and classes can be calculated by passing  $q_{j:\mu}$  through a standard convolution layer consisting of  $L \times L$  filters of size  $(2s + 1) \times (2s + 1)$ , where  $L$  is the number of labels and  $s$  is the number of neighbors each pixel interacts with in each dimension.

The contribution from the bilateral term is

$$(3.3) \quad v_{i:\lambda}^{bilateral} = \sum_{\substack{j:(i,j) \in \mathcal{E} \\ \mu \in \mathcal{L}}} k_{\lambda,\mu}^{bilateral}(\mathbf{f}_i - \mathbf{f}_j) q_{j:\mu}.$$

For this computation we utilize the method presented by Jampani, Kiefel, and Gehler [28] which is based on the permutohedral lattice introduced by Adams, Baek, and Davis [1]. Efficient computations are obtained by using the fact that the feature space is generally sparsely populated. Similar to the spatial filter, we get  $L \times L$  filters, each having size of  $(s + 1)^{d+1} - s^{d+1}$ , where  $s$  is the number of neighbors each pixel interacts with in each dimension in the sparse feature space.

For the *filterbank* version the contribution of the pairwise term can be calculated as

$$(3.4) \quad v_{i:\lambda}^{bank} = \sum_{f=1}^F g_f(\mathbf{p}_i, I) \sum_{\substack{j:(i,j) \in \mathcal{E} \\ \mu \in \mathcal{L}}} k_{x_i, x_j, f}^{spatial}(\mathbf{p}_i - \mathbf{p}_j) q_{j:\mu},$$

which, similar to the other spatial kernel, can be efficiently calculated using a standard convolution layer. The number of filters needed is  $L \times FL$ .

**3.2. Update step and projection to feasible set.** Given the energy gradient and a previous estimate of the solution, we want to improve our solution by taking a step which decreases the energy while still keeping the solution feasible. A straightforward approach to doing this would be to start by taking a step in the negative direction of the gradient according to

$$(3.5) \quad \tilde{\mathbf{q}}^{t+1} = \mathbf{q}^t - \gamma \nabla_{\mathbf{q}} E,$$

where  $\gamma$  is the step size. After taking the step, the values are projected onto the simplex  $\Delta^L$  satisfying  $\sum_{\lambda \in \mathcal{L}} q_{i:\lambda} = 1$  and  $0 \leq q_{i:\lambda} \leq 1$  by following the method by Chen and Ye [18]. This method is used by Larsson et al. in [34]. A drawback of this approach is that, if  $\tilde{\mathbf{q}}^{t+1}$  is outside of the simplex, back-propagation through the projection method will give zero gradients.

An alternative method is to use the entropic descent algorithm proposed by Beck and Teboulle [5]. In this method, the distance measure for the projection is the KL divergence instead of the Euclidean distance. Beck and Teboulle showed that the update step can be written on the following closed form:

$$(3.6) \quad q_{ij}^{k+1} = \frac{q_{ij}^k \exp(-t_k \nabla_{q^k} E)}{\sum q_{ij}^k \exp(-t_k \nabla_{q^k} E)}, \quad t_k = \frac{\sqrt{2 \ln n}}{L_f} \frac{1}{\sqrt{k}},$$

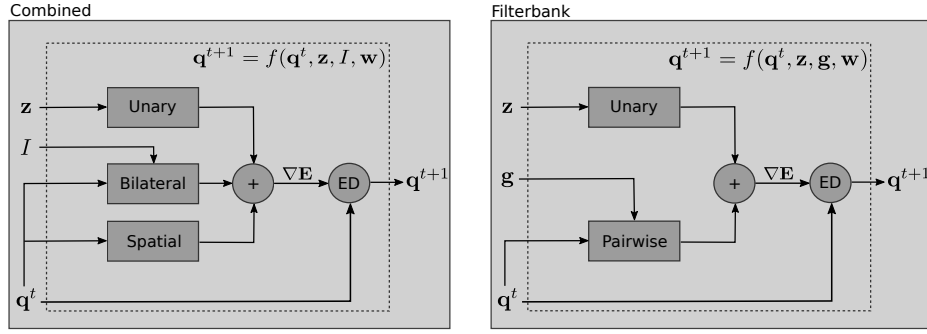
where  $n$  is the number of dimensions (the number of classes in our case),  $k$  is the iteration number, and  $L_f$  is a tunable parameter. Note that this projection is done individually for each pixel  $i$ .

**3.3. Comparison to mean-field.** In recent years, a popular choice for CRF inference has been to apply the mean-field algorithm. One reason is that the kernel evaluations can be computed with fast bilateral filtering [33]. As we have seen in this section, it can be accomplished with our framework as well, with formulas that are less involved. The main difference is that our framework directly optimizes the Gibbs energy, which corresponds to MAP, while mean-field optimizes KL divergence, which does not.

**4. Integration in a deep neural network.** In this section we will describe how the steps of Algorithm 3.1 can be formulated as layers in a neural network. For this, we need to be able to calculate error derivatives with respect to the input given error derivatives with respect to the output. In addition we need to be able to calculate the error derivatives with respect to the network parameters, i.e., the filter weights for the pairwise kernels as well as the unary weight. This will enable us to unroll the entire gradient descent process as a Recurrent Neural Network (RNN), making it possible to train both the parameters of the CRF as well as the parameters of the CNN that give the unary potentials as well as  $\mathbf{g}$ , the filter weighting function. A schematic of the data flow for one step is shown in Figure 1. In the supplementary material (M116726.01.pdf [local/web 264KB]), all derivative formulas are written out in detail.

**4.1. Initialization.** The variables  $\mathbf{q}^0$  are set as the output of the CNN, which has been pretrained to estimate the probability of each pixel containing each class and has a softmax layer as the last layer to ensure that the variables lie within zero and one.





**Figure 1.** The data flow of one iteration of the projected gradient descent algorithm. Each rectangle or circle represents an operation that can be performed within a deep learning framework; the ED component performs an entropic descent update step according to (3.6). Left: Combined version of CRF. Right: Filterbank version of CRF.

**4.2. Gradient computations.** We explained the gradient computations in section 3 for the forward pass. To describe the calculation of the error derivatives, we first notice that the gradient is calculated by summing the unary term and the pairwise term. We can hence treat these separately and combine them using an elementwise summing operation.

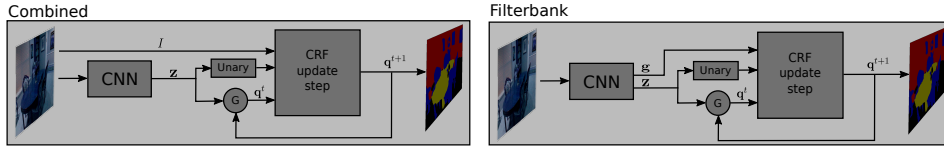
**Unary term.** The unary term in (2.3) is an elementwise operation with the CNN output as input and the unary weight  $w_u$  as parameter. The operation is obviously differentiable with respect to both the layer input as well as its parameter. Note that for  $w_u$  we get a summation over all class and pixel indexes for the error derivatives, while for the input the error derivatives are calculated elementwise.

**Pairwise term—combined version.** The spatial pairwise term of the gradient can be calculated efficiently using standard 2D convolution. In addition to giving us an efficient way of performing the forward pass, the 2D convolution layer can also be utilized to perform the backward pass, calculating the error derivatives with respect to the input and parameters. Similar to the spatial term, the bilateral term is also calculated utilizing a bilateral filtering technique. Jampani, Kiefel, and Gehler [28] also presented a way to calculate the error derivatives with respect to the parameters for an arbitrarily shaped bilateral filter.

**Pairwise term—filterbank version.** For the filterbank version we also use standard 2D convolution operations to calculate the pairwise part of the gradient. This makes the process of propagating the error derivatives similar to that for the spatial term of the combined version. Interpreting the calculations as two separate steps, one convolution with  $L \times FL$  filters and one weighted summation over the feature weights, the error derivative can be calculated with standard network layers. Note that the error derivatives with respect to the feature weights  $g_f$  are also calculated and propagated further back through the pairwise CNN.

**4.3. Entropic descent update.** The entropic descent step is done individually for each pixel. Since we have the update step on closed form, we can easily implement it as a layer in a deep learning framework. Regarding the error derivative, we are required to calculate the error derivatives both with respect to the values of the previous iteration,  $q^t$ , and with respect to gradient,  $\nabla_{q^t} E$ . The error derivatives with respect to the values of the previous iteration





**Figure 2.** The data flow of the deep structure model. Each rectangle or circle represents an operation that can be performed within a deep learning framework. Note that the CNN outputs both class probabilities  $\mathbf{z}$  and filterbank features  $\mathbf{g}$  for the filterbank version.

are given according to

$$(4.1) \quad \frac{\partial L}{\partial q_{ij}^k} = \frac{q_{ij}^{k+1}}{q_{ij}^k} \left( \frac{\partial L}{\partial q_{ij}^{k+1}} - \sum_{l=1}^n \frac{\partial L}{\partial q_{il}^{k+1}} q_{il}^{k+1} \right),$$

where the index  $i$  is over all pixels and  $j$  is over all the  $n$  number of classes. Note that the error derivatives with respect to  $q_{ij}^{k+1}$  are given by the previous iteration. The error derivatives with respect to the gradient are given according to

$$(4.2) \quad \frac{\partial L}{\partial y_{ij}} = -t_k q_{ij}^{k+1} \left( \frac{\partial L}{\partial q_{ij}^{k+1}} - \sum_{l=1}^n \frac{\partial L}{\partial q_{il}^{k+1}} q_{il}^{k+1} \right).$$

Note that, for ease of notation, we have used  $y_{ij}$  as the energy derivative of pixel  $i$  and class  $j$ .

**5. Recurrent formulation of deep structured model.** Our iterative solution to the CRF energy minimization problem by projected gradient descent, as described in the previous sections, is formulated as an RNN. The input to the RNN is the image and the outputs of the CNN, as shown in Figure 2. The unary CNN’s output,  $\mathbf{z}$ , is the unary potentials and is obtained after the final softmax layer (since the CNN is initially trained for classification). For the filterbank version the CNN also outputs image-dependent features,  $\mathbf{g}$ , which are “selecting” which filters to use to compose the pairwise term at each pixel location; see Figure 3 for a detailed schematic of the filterbank version.

Each iteration of the RNN performs one projected gradient descent step to approximately solve (2.7). Thus, one update step can be represented by

$$(5.1) \quad \mathbf{q}^{t+1} = f(\mathbf{q}^t, \mathbf{z}, I, \mathbf{w}).$$

As illustrated in Figure 2, the gating function  $G$  sets  $\mathbf{q}^t$  to  $\mathbf{z}$  at the first time step, and to  $\mathbf{q}^{t-1}$  at all other time steps. In our iterative energy minimization, the output of one step is the input to the next step. We initialize at  $t = 0$  with the output of the unary CNN.

The output of the RNN can be read off  $\mathbf{q}^T$ , where  $T$  is the total number of steps taken. In practice, we perform a set number of  $T$  steps where  $T$  is a hyperparameter. It is possible to run the RNN until convergence for each image (thus a variable number of iterations per image), but we observed minimal benefit in the final intersection over union (IoU) from doing so, as opposed to fixing the number of iterations to  $T = 5$ .

The parameters of the RNN are the filter weights for the pairwise kernels, and also the weight for the unary terms. Since we are able to compute error derivatives with respect to the parameters, and input of the RNN, we can back-propagate error derivatives through our RNN to the preceding CNN and train our entire network end-to-end. Furthermore, since the operations of the RNN are formulated as filtering, training and inference can be performed in a fully convolutional manner.

The CNN part of our network allows us to leverage the ability of CNNs to learn rich feature representations from data, while the RNN part of the network utilizes the CRF’s ability to model output structure. As we learn the parameters of our pairwise terms, we are not restricted to Gaussian potentials as in [32, 47], and we show the benefits of this in our experiments (section 7).

**6. Implementation details.** Our proposed CRF model has been implemented in the Caffe [29] library. The unary CNN part of our model is initialized from a pretrained segmentation network. For all experiments we use the Deeplab-LargeFOV proposed by Chen et al. [17]. For the *combined* version of our model the unary CNN is pretrained for pixelwise classification.

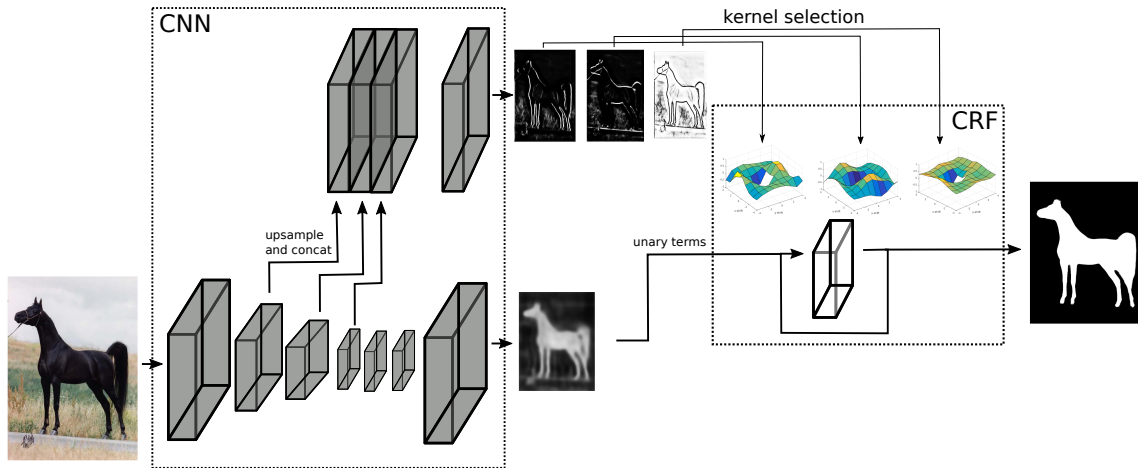
For the *filterbank* version we use a modified version of the Deeplab-LargeFOV where a second head is added to the network as in [16]. This head is formed by upsampling and concatenating several intermediate layers of the original network, then a final convolution is applied to the concatenated features, and finally a softmax layer is added. The second head outputs the filter choosing features  $\mathbf{g}_f$  and is pretrained to classify each pixel as a horizontal semantic edge, a vertical semantic edge, or no edge. This part of the base network can also be trained during the final end-to-end training.

Both the *combined* and the *filterbank* models can be trained from scratch; however, the training converges faster and more reliably when the unary part is pretrained.

The CRF model has several tunable hyperparameters. The parameter  $L_f$  and the number of iterations  $T$  specify the properties of the gradient descent algorithm.  $L_f$  influences the step size (larger  $L_f$  gives a smaller step size); too high a step size might make the algorithm not end up in a minimum, while setting a low step size and a low number of iterations might not give the algorithm a chance to converge. The kernel sizes for the pairwise kernels also need to be set. Choosing the value of these parameters gives a trade-off between model expression ability and number of parameters, which may cause (or hinder) overfitting.

The spatial weights of the CRF model are all initialized as zero with the motivation that we did not want to impose a shape for these filters, but instead see what was learned during training. The bilateral filters were initialized as Gaussians with the common Potts class interaction (the filters corresponding to interactions within the same class were set to zero) [33, 14, 47].

**7. Experiments.** We evaluate the proposed approach on three datasets: the WEIZMANN HORSE dataset [10], the NYU V2 geometric dataset [45], and the PASCAL VOC 2012 dataset [22]. In these experiments, we show that the proposed approach has advantages over similar approaches such as CRF-RNN [47]. In addition we show that adding a CRF model as proposed in this paper improves the results on strong unary CNN networks, even for cases where the CNN has been trained on large amounts of extra data.



**Figure 3.** Schematic of the filterbank version of our model. The CNN part outputs initial class probability maps as well as filter selection maps. The structure used for the CNN is a modified version of the Deeplab-LargeFOV [17] with an extra added head.

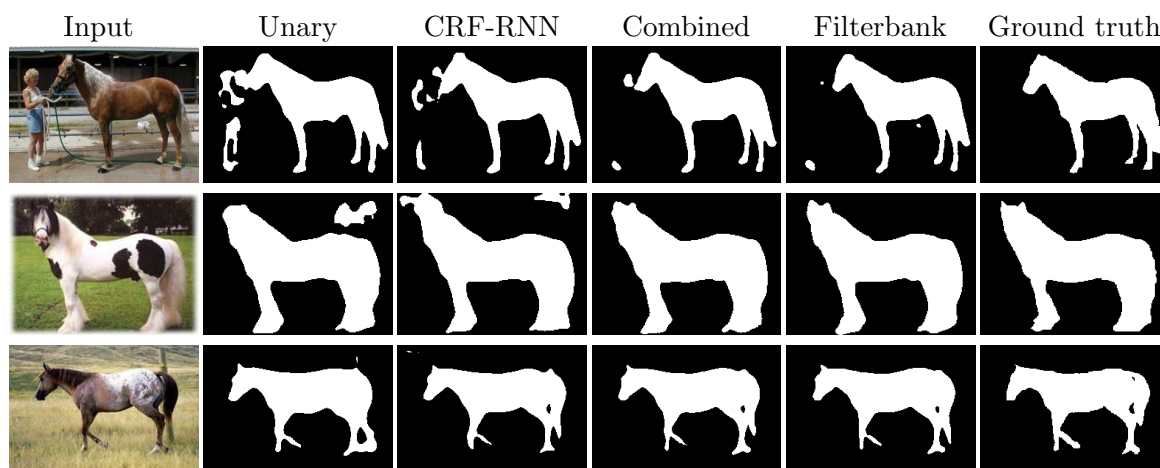
**7.1. Weizmann Horse.** The WEIZMANN HORSE dataset is widely used for benchmarking object segmentation algorithms. It contains 328 images of horses in different environments. We divide these images into a training set of 150 images, a validation set of 50 images, and a test set of 128 images. Our purpose is to verify our ability to learn reasonable kernels and study the effects of different settings on a relatively small dataset. In addition we use this dataset to evaluate our proposed inference method as well as the different types of CRF models. To compare the different types of inference methods we train our *combined* model with three types of inference methods: entropic descent (ED), projected gradient descent (PGD), and mean-field (MF). We also train a version with only Gaussian potentials (using the same potentials as for CRF-RNN [47]). We also trained and evaluated the *filterbank* version. The results are summarized in Table 1, and some example segmentations are shown in Figure 4. As can be seen from the results, our proposed inference method using entropic descent achieves slightly better results on the test set for the *combined* CRF model. However, the increase over mean-field and projected gradient descent inference is minor. For the case where we used Gaussian CRF potentials we get better results with entropic descent inference compared to mean-field. When comparing entropic descent inference and projected gradient descent, we find that the two methods achieve similar results. Training a model with projected gradient descent is, however, problematic due to the zeroing of gradients; to solve this we train with a “leaky” version of projected gradient descent. This means that the intermediate states and final results might not lie on the probability simplex—something that is guaranteed for entropic descent inference.

In Figure 5 the mean intersection over union on the test set is plotted as a function of the number of CRF inference iterations. During training the number of iterations was set to five. As can be seen in the figure, increasing the number of inference steps will only slightly increase the segmentation result. In Figure 6 the pairwise weights of the *filterbank* version are visualized.

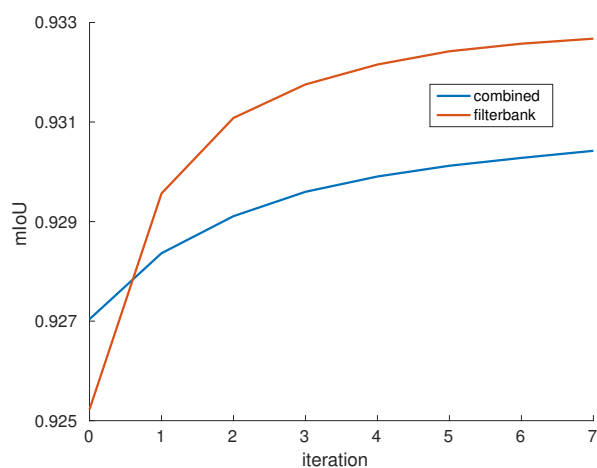
**Table 1**

Quantitative results on the WEIZMANN HORSE dataset comparing our method to baselines as well as comparison of different inference methods. Mean intersection over union for the test set is shown.

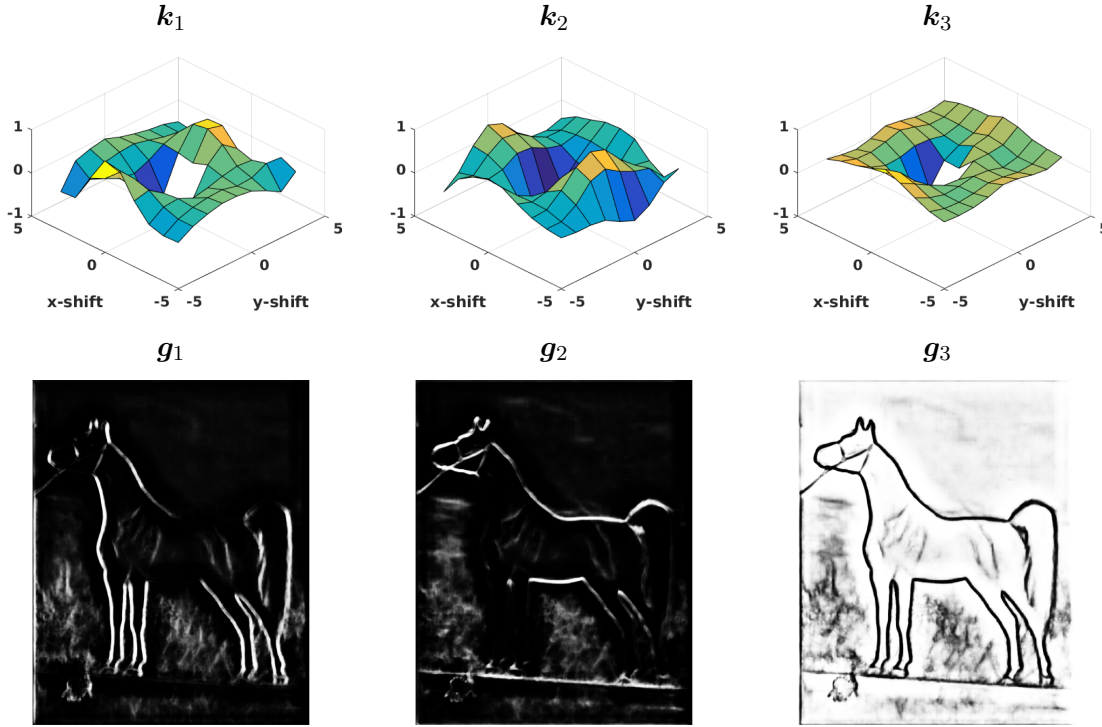
Method	mIoU (%)
Unary CNN - Deeplab [14]	90.89
CRF-RNN [47]	91.47
Gaussian-ED	92.64
Combined-ED	92.99
Combined-MF	92.73
Combined-PGD	92.79
Filterbank-ED	93.22



**Figure 4.** Qualitative results on the WEIZMANN HORSE dataset. Note that the proposed methods capture the shape of the horses better than the baselines, especially compared to the unary network.



**Figure 5.** WEIZMANN HORSE test set results in terms of mean intersection over union plotted as a function of the number of iterations for the CRF inference method. During training the number of iterations was set to five.



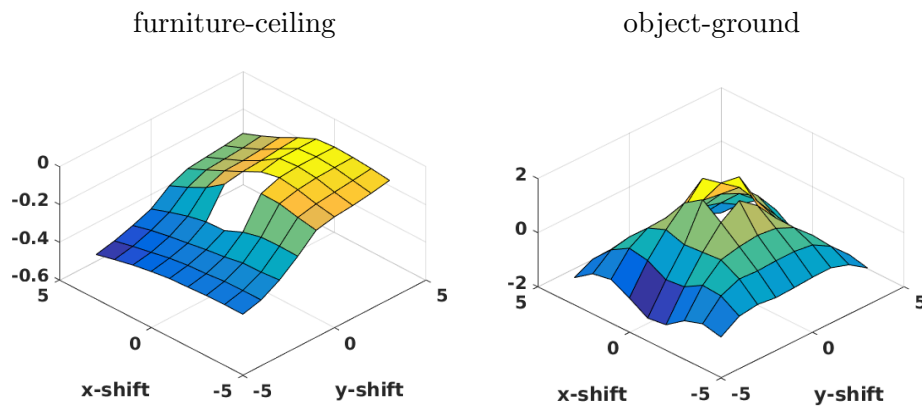
**Figure 6.** Visualization of the pairwise kernel weights for the filterbank version trained on the WEIZMANN HORSE data set. These weights are for the classes “background” and “horse.” The plots can be understood as the energy added when assigning the pixels with the relative positions  $(x,y)$  and  $(x+x\text{-shift},y+y\text{-shift})$  as background and horse. This energy is then multiplied by the “filter selection” map  $\mathbf{g}$  for each pixel and then summed. The first map of  $\mathbf{g}$  has high values at edges in the horizontal direction; looking at  $\mathbf{k}_1$  we see that changing classes in this direction does not add as much energy as changing classes in the vertical direction. Similar behavior can be seen for the second map. Note that the middle position has been removed from the kernel plots since it does not provide the same structural information as the other weights and hence cannot be interpreted in the same way.

**7.2. NYU V2.** The NYU V2 dataset contains images taken by a Microsoft Kinect V-1 camera in 464 indoor scenes. We use the official training and validation splits consisting of 795 and 654 images, respectively. Following the setting described in Wang et al. [45], we also include additional images for training. These are the images from the NYU V1 dataset that do not overlap with the images in the official validation set. This gives a total of 894 images with semantic label annotations for training. As in [45], we consider five classes conveying strong geometric properties: ground, vertical, ceiling, furniture, and objects.

As shown in Table 2, we achieved superior results for semantic image segmentation on the NYU V2 dataset. Some example segmentations are shown in Figure 8.

In Figure 7 the pairwise weights of the filterbank version are visualized.

**7.3. PASCAL VOC.** The PASCAL VOC 2012 segmentation benchmark [21] consists of 20 foreground classes and one background class. The unary network used for these experiments



**Figure 7.** Visualization of the pairwise kernel weights for the filterbank version trained on the NYU V2 data set. These weights are for the classes shown above the plots and for the third filter selection map which usually has a high value for pixels with no semantic edge. The furniture-ceiling kernel favors putting furniture labels below ceiling labels, while the object-ground kernel has a more Gaussian-like shape.

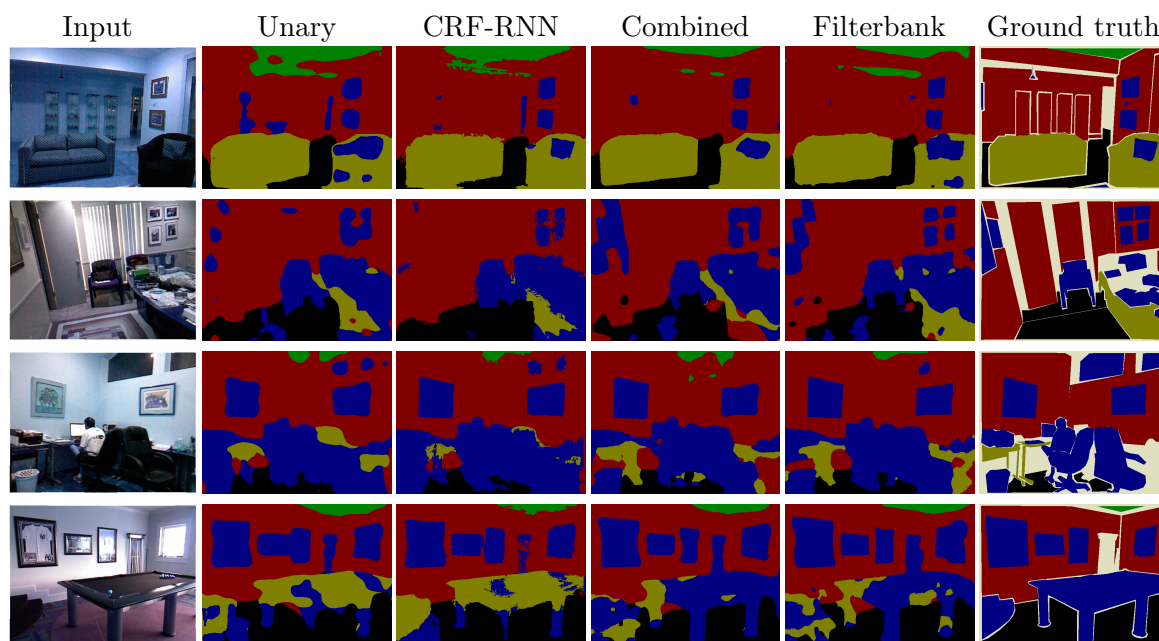
**Table 2**

Quantitative results comparing our method to baselines as well as state-of-the-art methods. Mean intersection over union for the validation set is shown for the NYU V2 dataset. The CRF-RNN baseline was initialized with the same unary network as the proposed models.

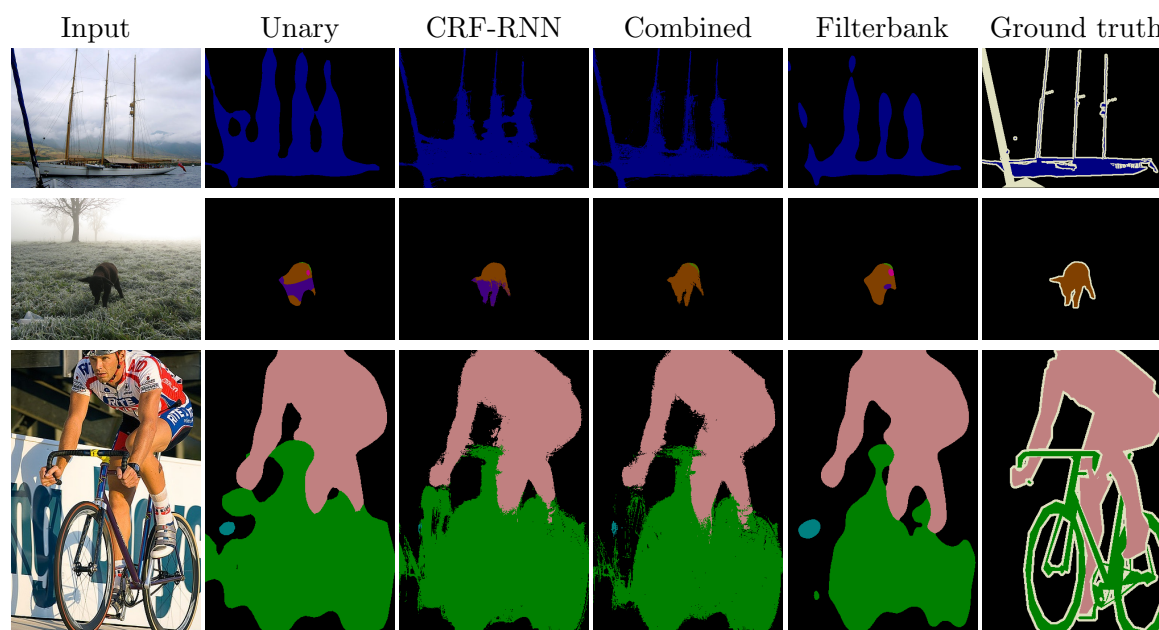
Method	mIoU (%)
R-CNN [24]	40.3
Semantic HCRF [45]	42.7
Joint HCRF [45]	44.2
Modular CNN [27]	54.3
Unary CNN - Deeplab [14]	62.8
CRF-RNN [47]	64.4
Combined	65.4
Filterbank	65.4

is again the Deeplab-LargeFOV network [17]; this network has been pretrained on the MSCOCO 2014 dataset [37] and then trained on the PASCAL VOC training data as well as a training set created from annotations of the semantic boundaries dataset [25]. We add our CRF models to this baseline network and train only on the PASCAL VOC training data. This is to show that we can improve upon really strong baselines, even though we fine-tune the complete models on only a fraction of the training data used for the baseline. The results for the PASCAL VOC 2012 validation set are shown in Table 3; in Figure 9 some example segmentations are shown. In addition we evaluate our model on the test set; for this the results are shown in Table 4. As can be seen, our models perform similarly to models trained with the same base network. Note that our models are only trained on the training data during end-to-end training. Recently there have been several CNNs with different base architectures presented that perform well, even without a CRF. The top entry at the moment is PSPNet [35] with an mIoU of 85.4; we leave it to future work to explore whether these architectures can be improved using our proposed methodology.





**Figure 8.** Qualitative results on the NYU V2 dataset. Note that the proposed methods capture the shape of the object instances better than the baselines. This effect is perhaps most pronounced for the paintings hanging on the walls. The pixels colored off-white are “ignore”-pixels; these are not counted in the evaluation. The training images have similar “ignore”-pixels.



**Figure 9.** Qualitative results on the PASCAL VOC dataset [21]. The pixels colored off-white are “ignore”-pixels; these are not counted in the evaluation. The training images have similar “ignore”-pixels.



**Table 3**

Quantitative results on the PASCAL VOC 2012 validation set. The CRF-RNN baseline was initialized with the same unary network as the proposed models. The unary model was pretrained on the MS-COCO 2014 dataset [37].

Method	mIoU (%)
Unary CNN - Deeplab [14]	68.5
CRF-RNN [47]	71.7
Combined	72.0
Filterbank	70.1

**Table 4**

Quantitative results on the PASCAL VOC 2012 test set. The three top entries use the same base network as our models. The unary model was pretrained on the MS-COCO 2014 dataset [37], but note that our models were not trained using MS-COCO.

Method	mIoU (%)
Unary CNN - Deeplab [14]	68.9
DT-EdgeNet [16]	71.7
CRF-RNN [47]	72.2
Combined	72.5
Filterbank	69.5

**7.4. Execution time.** We also investigated the difference in running time between the two proposed models. This was done on a computer with a Nvidia Titan X GPU with Pascal architecture and an Intel i7-5930K processor. The implementation used for the bilateral filtering was the one from Jampani, Kiefel, and Gehler [28], where most of the computations are done on the GPU. The initialization of the permutohedral lattice, however, is done on the CPU. The runtimes were tested by performing the forward step for a randomized RGB image of size  $640 \times 640$  with 21 classes. The numbers presented are the average of 100 runs. For the *combined* model the forward runtime was 12 seconds, while for the *filterbank* it was 0.37 second.

**8. Conclusion.** In this paper we have presented a gradient descent-based method for inference in conditional random fields. This method allows for back-propagation of error derivative, hence enabling end-to-end training with a convolutional neural network of choice. We show that this inference method has beneficial properties and performs better on some tasks compared to other methods such as mean-field. In addition, we present two types of conditional random field models tailored for semantic segmentation. The *combined* model uses spatial pairwise terms as well as image-dependent bilateral pairwise terms. This model performs well but is somewhat computationally expensive due to the high dimensionality of the bilateral filtering. We also present the *filterbank* model, which is also image-dependent but only requires 2D convolutions during inference. The image dependence of the model comes from an output map of the base CNN that acts as a “filter choosing” map. This enables the model to, for example, use one filter representing pairwise interaction between pixel labels at semantic edges and another filter far away from semantic edges. This model gives a speedup by a factor of 32 compared to the *combined* model without losing performance

in terms of segmentation quality. For the smaller dataset it achieves similar segmentation quality as the *combined* model. Since the *filterbank* version of the model learns how the pairwise term should depend on the image, it is in this aspect more expressive than the *combined* version. The pairwise terms of the *combined* model, however, due to its dependence on the color gradient, are hand-crafted to preserve and refine edges. This is beneficial for the PASCAL VOC dataset, where the unary network generally captures the context well but outputs “blobby” segmentations. For all the models presented, the pairwise kernels can have arbitrary shape instead of commonly used Gaussian kernels. This enables the models to learn more complicated pairwise label interactions.

## REFERENCES

- [1] A. ADAMS, J. BAEK, AND M. A. DAVIS, *Fast high-dimensional filtering using the permutohedral lattice*, Computer Graphics Forum, 29 (2010), pp. 753–762.
- [2] T. AJANTHAN, A. DESMAISON, R. BUNEL, M. SALZMANN, P. H. S. TORR, AND M. P. KUMAR, *Efficient linear programming for dense CRFs*, in 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Washington, DC, 2017, pp. 2934–2942.
- [3] A. ARNAB, S. JAYASUMANA, S. ZHENG, AND P. H. S. TORR, *Higher order conditional random fields in deep neural networks*, in European Conference on Computer Vision 2016, Lecture Notes in Comput. Sci. 9906, Springer, Berlin, 2016, pp. 524–540.
- [4] A. ARNAB, S. ZHENG, S. JAYASUMANA, B. ROMERA-PAREDES, M. LARSSON, A. KIRILLOV, B. SAVCHYNSKY, C. ROTHER, F. KAHL, AND P. H. S. TORR, *Conditional random fields meet deep neural networks for semantic segmentation: Combining probabilistic graphical models with deep learning for structured prediction*, IEEE Signal Process. Mag., 35 (2018), pp. 37–52, <https://doi.org/10.1109/MSP.2017.2762355>.
- [5] A. BECK AND M. TEBoulLE, *Mirror descent and nonlinear projected subgradient methods for convex optimization*, Oper. Res. Lett., 31 (2003), pp. 167–175.
- [6] D. BELANGER AND A. MCCALLUM, *Structured prediction energy networks*, in Proceedings of the International Conference on Machine Learning (ICML’16), ACM, New York, 2016, pp. 983–992.
- [7] D. BELANGER, B. YANG, AND A. MCCALLUM, *End-to-End Learning for Structured Prediction Energy Networks*, preprint, <https://arxiv.org/abs/1703.05667>, 2017.
- [8] G. BERTASIUS, L. TORRESANI, S. X. YU, AND J. SHI, *Convolutional random walk networks for semantic image segmentation*, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (Honolulu, HI), IEEE, Washington, DC, 2017.
- [9] A. BLAKE, P. KOHLI, AND C. ROTHER, *Markov Random Fields for Vision and Image Processing*, MIT Press, Cambridge, MA, 2011.
- [10] E. BORENSTEIN AND S. ULLMAN, *Class-specific, top-down segmentation*, in European Conference on Computer Vision 2002, Lecture Notes in Comput. Sci. 2351, Springer, Berlin, 2002, pp. 109–122.
- [11] E. BOROS AND P. HAMMER, *Pseudo-Boolean optimization*, Discrete Appl. Math., 123 (2002), pp. 155–225.
- [12] L. BOTTOU, Y. BENGIO, AND Y. LE CUN, *Global training of document processing systems using graph transformer networks*, in IEEE Conference on Computer Vision and Pattern Recognition, IEEE, Washington, DC, 1997, pp. 489–494.
- [13] S. CHANDRA AND I. KOKKINOS, *Fast, exact and multi-scale inference for semantic image segmentation with deep Gaussian CRFs*, in European Conference on Computer Vision 2016, Lecture Notes in Comput. Sci. 9911, Springer, Berlin, pp. 402–418.
- [14] L. CHEN, G. PAPANDREOU, I. KOKKINOS, K. MURPHY, AND A. L. YUILLE, *Semantic image segmentation with deep convolutional nets and fully connected CRFs*, in International Conference on Learning Representations, San Diego, CA, 2015.
- [15] L. CHEN, A. SCHWING, A. YUILLE, AND R. URTASUN, *Learning deep structured models*, in Proceedings of the International Conference on Machine Learning (Lille, France), ACM, New York, 2015, pp. 1785–1794.

- [16] L.-C. CHEN, J. T. BARRON, G. PAPANDREOU, K. MURPHY, AND A. L. YUILLE, *Semantic image segmentation with task-specific edge detection using CNNs and a discriminatively trained domain transform*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, Washington, DC, 2016, pp. 4545–4554.
- [17] L.-C. CHEN, G. PAPANDREOU, I. KOKKINOS, K. MURPHY, AND A. L. YUILLE, *DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs*, IEEE Trans. Pattern Anal. Mach. Intell., 40 (2018), pp. 834–848.
- [18] Y. CHEN AND X. YE, *Projection onto a Simplex*, preprint, <https://arxiv.org/abs/1101.6081>, 2011.
- [19] A. DESMAISON, R. BUNEL, P. KOHLI, P. H. S. TORR, AND M. P. KUMAR, *Efficient continuous relaxations for dense CRF*, in European Conference on Computer Vision 2016, Lecture Notes in Comput. Sci. 9906, Springer, Berlin, 2016, pp. 818–833.
- [20] J. DOMKE, *Learning graphical model parameters with approximate marginal inference*, IEEE Trans. Pattern Anal. Mach. Intell., 35 (2013), pp. 2454–2467.
- [21] M. EVERINGHAM, S. M. A. ESLAMI, L. VAN GOOL, C. K. I. WILLIAMS, J. WINN, AND A. ZISSERMAN, *The PASCAL visual object classes challenge: A retrospective*, Internat. J. Comput. Vis., 111 (2015), pp. 98–136, <https://doi.org/10.1007/s11263-014-0733-5>.
- [22] M. EVERINGHAM, L. V. GOOL, C. K. I. WILLIAMS, J. WINN, AND A. ZISSERMAN, *The PASCAL visual object classes (VOC) challenge*, Internat. J. Comput. Vis., 88 (2010), pp. 303–338.
- [23] G. GHIASI AND C. FOWLKES, *Laplacian reconstruction and refinement for semantic segmentation*, in European Conference on Computer Vision 2016, Lecture Notes in Comput. Sci. 9906, Springer, Berlin, 2016, pp. 519–534.
- [24] R. GIRSHICK, J. DONAHUE, T. DARRELL, AND J. MALIK, *Rich feature hierarchies for accurate object detection and semantic segmentation*, in IEEE Conference on Computer Vision and Pattern Recognition (Columbus, OH), IEEE, Washington, DC, 2014.
- [25] B. HARIHARAN, P. ARBELEZ, L. BOURDEV, S. MAJI, AND J. MALIK, *Semantic contours from inverse detectors*, in 2011 International Conference on Computer Vision (Barcelona, Spain), IEEE, Washington, DC, 2011, pp. 991–998, <https://doi.org/10.1109/ICCV.2011.6126343>.
- [26] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in IEEE Conference on Computer Vision and Pattern Recognition (Las Vegas, NV), IEEE, Washington, DC, 2016.
- [27] O. H. JAFARI, O. GROTH, A. KIRILLOV, M. Y. YANG, AND C. ROTHER, *Analyzing modular CNN architectures for joint depth prediction and semantic segmentation*, in Proceedings of the IEEE International Conference on Robotics and Automation, IEEE, Washington, DC, 2017, pp. 4620–4627.
- [28] V. JAMPANI, M. KIEFEL, AND P. V. GEHLER, *Learning sparse high dimensional filters: Image filtering, dense CRFs and bilateral neural networks*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, Washington, DC, 2016, pp. 4452–4461.
- [29] Y. JIA, E. SHELHAMER, J. DONAHUE, S. KARAYEV, J. LONG, R. GIRSHICK, S. GUADARRAMA, AND T. DARRELL, *Caffe: Convolutional Architecture for Fast Feature Embedding*, preprint, <https://arxiv.org/abs/1408.5093>, 2014.
- [30] A. KIRILLOV, D. SCHLESINGER, S. ZHENG, B. SAVCHYNSKY, P. TORR, AND C. ROTHER, *Joint training of generic CNN-CRF models with stochastic optimization*, in Asian Conference on Computer Vision 2016, Lecture Notes in Comput. Sci. 10112, Springer, Berlin, pp. 221–236.
- [31] D. KOLLER AND N. FRIEDMAN, *Probabilistic Graphical Models*, MIT Press, Cambridge, MA, 2009.
- [32] P. KRAEHNBUHL AND V. KOLTUN, *Parameter learning and convergent inference for dense random fields*, in Proceedings of the 30th International Conference on Machine Learning, ACM, New York, 2013, pp. 513–521.
- [33] P. KRÄHNBUHL AND V. KOLTUN, *Efficient inference in fully connected CRFs with Gaussian edge potentials*, in Proceedings of Neural Information Processing Systems 2011, Curran Associates, Red Hook, NY, 2011, pp. 109–117.
- [34] M. LARSSON, A. ARNAB, F. KAHL, S. ZHENG, AND P. H. S. TORR, *A projected gradient descent method for CRF inference allowing end-to-end training of arbitrary pairwise potentials*, in 11th International Conference on Energy Minimization Methods in Computer Vision and Pattern Recognition (EMM-CVPR), Springer, Berlin, 2017.
- [35] G. LIN, A. MILAN, C. SHEN, AND I. REID, *RefineNet: Multi-Path Refinement Networks for High-Resolution Semantic Segmentation*, preprint, <https://arxiv.org/abs/1611.06612>, 2016.

- [36] G. LIN, C. SHEN, A. HENGEL, AND I. REID, *Efficient piecewise training of deep structured models for semantic segmentation*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, Washington, DC, 2016, pp. 3194–3203.
- [37] 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 Proceedings of the European Conference on Computer Vision, Lecture Notes in Comput. Sci. 8693, Springer, Berlin, 2014, pp. 740–755, [https://doi.org/10.1007/978-3-319-10602-1\\_48](https://doi.org/10.1007/978-3-319-10602-1_48).
- [38] Z. LIU, X. LI, P. LUO, C. C. LOY, AND X. TANG, *Semantic image segmentation via deep parsing network*, in Proceedings of the International Conference on Computer Vision, IEEE, Washington, DC, 2015, pp. 1377–1385.
- [39] J. LONG, E. SHELHAMER, AND T. DARRELL, *Fully convolutional networks for semantic segmentation*, IEEE Trans. Pattern Recognition, 39 (2017), pp. 640–651.
- [40] S. REN, K. HE, R. GIRSHICK, AND J. SUN, *Faster R-CNN: Towards real-time object detection with region proposal networks*, in Proceedings of Neural Information Processing Systems 2015, MIT Press, Cambridge, MA, 2015, pp. 91–99.
- [41] C. ROTHER, V. KOLMOGOROV, AND A. BLAKE, *“GrabCut”: Interactive foreground extraction using iterated graph cuts*, ACM Trans. Graphics, 23 (2004), pp. 309–314.
- [42] A. SCHWING AND R. URTASUN, *Fully Connected Deep Structured Networks*, preprint, <https://arxiv.org/abs/1503.02351>, 2015.
- [43] K. SIMONYAN AND A. ZISSERMAN, *Very deep convolutional networks for large-scale image recognition*, in Proceedings of the International Conference on Learning Representations, 2015, pp. 1–14.
- [44] V. VINEET, J. WARRELL, AND P. TORR, *Filter-based mean-field inference for random fields with higher order terms and product label-spaces*, in European Conference on Computer Vision 2012, Lecture Notes in Comput. Sci. 7576, Springer, Berlin, pp. 31–44.
- [45] P. WANG, X. SHEN, Z. LIN, S. COHEN, B. PRICE, AND A. YUILLE, *Towards unified depth and semantic prediction from a single image*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition 2014, IEEE, Washington, DC, 2014.
- [46] W. WANG, S. FIDLER, AND R. URTASUN, *Proximal deep structured models*, in Proceedings of Neural Information Processing Systems 2016, Curran Associates, Red Hook, NY, 2016, pp. 865–873.
- [47] S. ZHENG, S. JAYASUMANA, B. ROMERA-PAREDES, V. VINEET, Z. SU, D. DU, C. HUANG, AND P. TORR, *Conditional random fields as recurrent neural networks*, in Proceedings of the IEEE International Conference on Computer Vision 2015, IEEE, Washington, DC, 2015, pp. 1529–1537.