

Revisiting Deep Structured Models for Pixel-Level Labelling with Gradient-Based Inference ^{*†}

Måns Larsson [‡], Anurag Arnab [§], Shuai Zheng [§], Philip Torr [§], and Fredrik Kahl [‡]

Abstract. Pixel-level labelling tasks such as semantic segmentation have witnessed 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, 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 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 modelled 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

1. Introduction. Markov Random Fields (MRFs), Conditional Random Fields (CRFs) and more generally, probabilistic graphical models are a ubiquitous tool used in a variety of domains spanning Computer Vision, Computer Graphics and Image Processing [32, 9, 4]. In this paper, we focus on the application of MRFs for Computer Vision problems involving per-pixel labelling such as image segmentation. There are many successful approaches in this line of research, such as the interactive segmentation of [42] using graph cuts and the semantic segmentation works of [34, 45] 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 [27, 44, 41]. However, we observe that several leading segmentation approaches still include CRFs, either as a post-processing step [14, 15, 24, 13], or as part of the deep neural network itself [48, 37, 3, 39, 31, 47].

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 [48, 3, 6, 47], we use a recurrent neural network to unroll the iterative inference steps of a CRF. This was first used in [48] and [43] to imitate mean-field inference and to train a fully

^{*}Submitted to the editors on 2018-01-30.

[†]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" [35].

[‡]Department of Electrical Engineering, Chalmers University of Technology, Gothenburg (mans.larsson@chalmers.se, fredrik.kahl@chalmers.se).

[§]Department of Engineering Science at University of Oxford (aarnab@robots.ox.ac.uk, szheng@robots.ox.ac.uk, phst@robots.ox.ac.uk).

convolutional network [40, 14] along with a CRF end-to-end via back propagation. In contrast to mean field, we do not optimize the KL-divergence between the true probability distribution and a fully-factorised 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 followup work [7] for multi-label classification problems in machine learning with impressive results. Moreover, [20, 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 [33, 48, 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 [16], a general framework for learning arbitrary potentials in deep structured model was proposed based on approximate ML learning. One of the advantages with 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 [31] 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 [37] and [13] also learn arbitrary pairwise potentials to model contextual relations between parts of the image. However, their approaches still perform post-processing with a CRF model with parametric Gaussian potentials. In [29], 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 [29], 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 is that 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 [21] 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-labelling 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 by [47].

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 approximate 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 a 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

80 filter far away from semantic edges.

81 Previous approaches trying to find alternatives to the computation heavy bilateral CRF
 82 include [17] where they use discriminatively trained domain transform as an edge-preserving
 83 filtering method. The authors show that the domain transform can be applied as a Recurrent
 84 Neural Network (RNN) applied across the image across all directions. Another example is [8]
 85 where they add a final layer that performs random graph walk across the image refining the
 86 segmentation.

87 In summary, our contributions are as follows.

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

97 Our framework has been implemented in CAFFE [30] and all source code is publicly avail-
 98 able to facilitate further research. ¹

99 **2. CRF Formulation.** Consider a Conditional Random Field over N discrete random
 100 variables $\mathcal{X} = \{X_1, \dots, X_N\}$ conditioned on an observation \mathbf{I} and let $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ be an
 101 undirected graph whose vertices are the random variables $\{X_1, \dots, X_N\}$. Each random vari-
 102 able corresponds to a pixel in the image and takes values from a predefined set of L labels
 103 $\mathcal{L} = \{0, \dots, L - 1\}$. The pair $(\mathcal{X}, \mathbf{I})$ is modelled as a CRF characterized by the Gibbs distri-
 104 bution

$$105 \quad (2.1) \quad P(\mathcal{X} = \mathbf{x} | \mathbf{I}) = \frac{1}{Z(\mathbf{I})} \exp(-E(\mathbf{x} | \mathbf{I})),$$

106 where $E(\mathbf{x} | \mathbf{I})$ denotes the Gibbs energy function with respect to the labeling $\mathbf{x} \in \mathcal{L}^N$ and
 107 $Z(\mathbf{I})$ is the partition function. To simplify notation the conditioning on \mathbf{I} will from now on
 108 be dropped. The MAP inference problem for the CRF model is equivalent to the problem of
 109 minimizing the energy $E(\mathbf{x})$. In this paper, we only consider energies containing unary and
 110 pairwise terms. The energy function can hence be written as

$$111 \quad (2.2) \quad E(\mathbf{x}) = \sum_{i \in \mathcal{V}} \psi_i(x_i) + \sum_{(i,j) \in \mathcal{E}} \psi_{ij}(x_i, x_j)$$

112 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.
 113 We now describe these potentials before discussing inference in Sec. 3.

114 **2.1. Potentials.** The unary potential $\psi_i(x_i)$ specifies the energy cost of assigning label x_i
 115 to pixel i . In this work we obtain our unary potentials from a CNN. Roughly speaking, the

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

116 CNN outputs a probability estimate of each pixel containing each class. Denoting the output
 117 of the CNN for pixel i and class x_i as $z_{i:x_i}$, the unary potential is

$$118 \quad (2.3) \quad \psi_i(x_i) = -w_u \log(z_{i:x_i} + \epsilon)$$

119 where w_u is a parameter controlling the impact of the unary potentials, and ϵ is introduced
 120 to avoid numerical problems.

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

125 **2.1.1. Combined.** The *combined* version has pairwise potentials that consist of a sum of
 126 one spatial term and one bilateral term. It has the following form

$$127 \quad (2.4) \quad \psi_{ij}(x_i, x_j) = k_{x_i, x_j}^{spatial}(\mathbf{p}_i - \mathbf{p}_j) + k_{x_i, x_j}^{bilateral}(\mathbf{f}_i - \mathbf{f}_j)$$

128 Here $k_{x_i, x_j}^{spatial}$ denote a spatial kernel with compact support. Its value depends on the relative
 129 position coordinates $\mathbf{p}_i - \mathbf{p}_j$ between pixels i and j . We do not restrict these spatial terms to
 130 any specific shape. However we restrict the support of the potential meaning that if pixels i
 131 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
 132 kernels with compact support in contrast to the commonly used dense Gaussian potential since
 133 this allows the inference calculations to be performed using standard 2D convolutions. The
 134 CRFs with Gaussian potentials do not in theory have compact support, and therefore, they
 135 are often referred to as dense. However, in practice, the exponential function in the kernel
 136 drops off quickly and effectively, the interactions between pixels far apart are negligible.

137 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
 138 for pixels i and j , respectively. Following several previous works on random fields, we let
 139 the vector depend on pixel coordinates \mathbf{p}_i and RGB values associated to the pixel, hence \mathbf{f}_i
 140 is a 5-dimensional vector. Note that for both the spatial and the bilateral kernels, there is
 141 one kernel for each label-to-label (x_i and x_j) interaction to enable the model learn differently
 142 shaped kernels for each of these interactions.

143 **2.1.2. Filterbank.** The pairwise potentials of the *filterbank* version has the following form

$$144 \quad (2.5) \quad \psi_{ij}(x_i, x_j) = \sum_{f=1}^F g_f(\mathbf{p}_i, I) k_{x_i, x_j, f}^{spatial}(\mathbf{p}_i - \mathbf{p}_j),$$

145 where $k_{x_i, x_j, f}^{spatial}$ denote a spatial kernel with compact support similar to the case of the *combined*
 146 version. The weights g_f depends both on the position of the pixel as well as the image I .
 147 These weights are taken as the output of a CNN. Hence, this gives rise to an image-dependent
 148 potential, but one only needs convolve with a bank of F 2D filters to evaluate it during
 149 inference. For example, the CNN outputting the weights can learn to detect semantic edges
 150 meaning that we would apply a different spatial filter close to a semantic edge than at the
 151 center of an semantic object. Setting the last layer of the CNN as a softmax the features \mathbf{g}_f
 152 act as "filter selectors" deciding how the several 2d-filters describing the pairwise term should
 153 be weighted for each pixel individually.

154 **2.2. Multi-label Graph Expansion and Relaxation.** To be able to explain our inference
 155 method we reformulate the original minimization of $E(\mathbf{x})$ as a real-valued optimization prob-
 156 lem. To facilitate a continuous relaxation of the energy minimisation problem we start off by
 157 expanding our original graph in the following manner. Each vertex in the original graph \mathcal{G} will
 158 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
 159 variable X_i is equivalent to an assignment of boolean labels 0 or 1 to each node $X_{i:\lambda}$, whereby
 160 an assignment of label 1 to $X_{i:\lambda}$ means that in the multi-label assignment, X_i receives label
 161 λ . To ensure that only one label is assigned to each node, an additional constraint is needed
 162 saying that, for each i , only one of $X_{i:\lambda}$ are allowed to be labeled 1. This enables to rewrite
 163 the energy minimization problem $\min E(\mathbf{x})$ as the following equivalent integer program

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

165 As a next step, we relax the integer program by allowing real values on the unit interval
 166 $[0, 1]$ instead of booleans only. We denote the relaxed variables $q_{i:\lambda} \in [0, 1]$. We can now write
 167 our problem as a quadratic program

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

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

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

$$174 \quad E(\mathbf{x}^*) = E(\mathbf{q}^*).$$

175 In the supplementary material, we show that for *any* real \mathbf{q} , one can obtain a binary \mathbf{x} such
 176 that $E(\mathbf{x}) \leq E(\mathbf{q})$. In particular, it will be true for \mathbf{x}^* and \mathbf{q}^* , which implies $E(\mathbf{x}^*) = E(\mathbf{q}^*)$.
 177 Note that the proof is constructive.

178 **3. MAP Inference via Gradient Descent Minimization.** To solve the program stated
 179 in (2.7) we propose an optimization scheme based on projected gradient descent, see Algo-
 180 rithm 3.1. It was designed with an extra condition in mind, that all operations should be
 181 differentiable to enable back propagation during training.

Algorithm 3.1 Algorithm 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  $\tilde{\mathbf{q}}_{i:\lambda}^{t+1}$  to the probability simplex  $\Delta^L$ .  $\mathbf{q}^{t+1} = \text{Proj}_{\Delta^L}(\tilde{\mathbf{q}})$ .
end for
return  $\mathbf{q}^{T-1}$ 

```

182 **3.1. Gradient Computations.** The gradient $\nabla_{\mathbf{q}}E$ of the objective function $E(\mathbf{q})$ in (2.7)
183 has the following elements

$$184 \quad (3.1) \quad \frac{\partial E}{\partial q_{i:\lambda}} = \psi_i(\lambda) + \sum_{\substack{j:(i,j) \in \mathcal{E} \\ \mu \in \mathcal{L}}} \psi_{ij}(\lambda, \mu) q_{j:\mu}.$$

185 The contribution from the spatial kernel in ψ_{ij} , cf. (2.4), can be written as

$$186 \quad (3.2) \quad v_{i:\lambda}^{spatial} = \sum_{\substack{j:(i,j) \in \mathcal{E} \\ \mu \in \mathcal{L}}} k_{\lambda,\mu}^{spatial}(\mathbf{p}_i - \mathbf{p}_j) q_{j:\mu}.$$

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

191 The contribution from the bilateral term is

$$192 \quad (3.3) \quad v_{i:\lambda}^{bilateral} = \sum_{\substack{j:(i,j) \in \mathcal{E} \\ \mu \in \mathcal{L}}} k_{\lambda,\mu}^{bilateral}(\mathbf{f}_i - \mathbf{f}_j) q_{j:\mu}.$$

193 For this computation we utilize the method presented by Jampani *et al.* [29] which is based on
194 the permutohedral lattice introduced by Adams *et al.* [1]. Efficient computations are obtained
195 by using the fact that the feature space is generally sparsely populated. Similar to the spatial
196 filter we get $L \times L$ filters, each having size of $(s + 1)^{d+1} - s^{d+1}$ where s is the number of
197 neighbours each pixel interacts with in each dimension in the sparse feature space.

198 For the *filter bank* version the contribution of the pairwise term can be calculated as

$$199 \quad (3.4) \quad v_{i:\lambda}^{bank} = \sum_{f=1}^F g_f(\mathbf{p}_i, I) \sum_{\substack{j:(i,j) \in \mathcal{E} \\ \mu \in \mathcal{L}}} k_{x_i, x_j, f}^{spatial}(\mathbf{p}_i - \mathbf{p}_j) q_{j:\mu},$$

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

202 **3.2. Update Step and Projection to Feasible Set.** Given the energy gradient and a
 203 previous estimate of the solution we want to improve our solution by taking a step which
 204 decreases the energy while still keeping the solution feasible. A straightforward approach of
 205 doing this would be to start by taking a step in the negative direction of the gradient according
 206 to

$$207 \quad (3.5) \quad \tilde{\mathbf{q}}^{t+1} = \mathbf{q}^t - \gamma \nabla_{\mathbf{q}} \mathbf{E},$$

208 where γ is the the step size. After taking the step the values are projected onto the simplex
 209 Δ^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 *et al.* [19].
 210 This method is used by by Larsson *et al.* in [35]. A drawback with this approach is that, if
 211 $\tilde{\mathbf{q}}^{t+1}$ is outside of the simplex, backpropagation through the projection method will give zero
 212 gradients.

213 An alternative method is to use the entropic descent algorithm proposed by Beck *et al.*
 214 [5]. In this method, the distance measure for the projection is the Kullback-Leibler divergence
 215 in contrast to the Euclidean distance. Beck *et al.* showed that the update step can be written
 216 on the following closed form

$$217 \quad (3.6) \quad q_{ij}^{k+1} = \frac{q_{ij}^k \exp(-t_k \nabla_{q^k} E)}{\sum q_{ij}^k \exp(-t_k \nabla_{q^k} E)}, \quad t_k = \frac{\sqrt{2 \ln n}}{L_f} \frac{1}{\sqrt{k}}$$

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

221 **3.3. Comparison to Mean-Field.** In recent years, a popular choice for CRF inference is
 222 to apply the mean-field algorithm. One reason is that the kernel evaluations can be computed
 223 with fast bilateral filtering [34]. As we have seen in this section, it can be accomplished with
 224 our framework as well, with formulas that are less involved. The main difference is that our
 225 framework directly optimizes the Gibbs energy which corresponds to MAP while mean-field
 226 optimizes KL-divergence which does not.

227 **4. Integration in a Deep Neural Network.** In this section we will describe how the steps
 228 of Algorithm 3.1 can be formulated as layers in a neural network. For this, we need to be able
 229 to calculate error derivatives with respect to the input given error derivatives with respect to
 230 the output. In addition we need to be able to calculate the error derivatives with respect to the
 231 network parameters, i.e. the filter weights for the pairwise kernels as well as the unary weight.
 232 This will enable us to unroll the entire gradient descent process as a Recurrent Neural Network
 233 (RNN) making it possible to train both the parameters of the CRF as well as the parameters
 234 of the CNN that gives the unary potentials as well as \mathbf{g} , the filter weighting function. A
 235 schematic of the data flow for one step is shown in Fig. 1. In the supplementary material, all
 236 derivative formulas are written out in detail.

237 **4.1. Initialization.** The variables \mathbf{q}^0 are set as the output of the CNN, which has been
 238 pretrained to estimate the probability of each pixel containing each class and has a softmax
 239 layers as the last layer to ensure that the variables lies within zero and one.

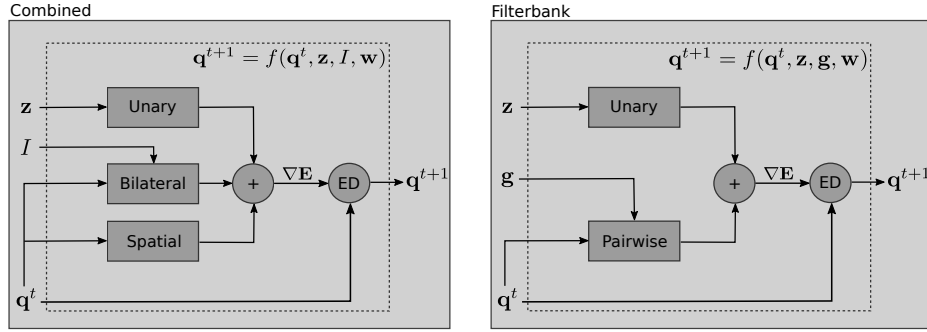


Figure 1. The data flow of one iteration of the projected gradient descent algorithm. Each rectangle or circle represent an operation that can be performed within a deep learning framework, the ED component performs an entropic descent update step according to equation (3.6). Left: Combined version of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

240 **4.2. Gradient Computations.** We have previously explained the gradient computations
 241 in Section 3 for the forward pass. To describe the calculation of the error derivatives we first
 242 notice that the gradient is calculated by summing the unary term and the pairwise term. We
 243 can hence treat these separately and combine them using an element-wise summing operation.

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

249 **Pairwise Term - Combined version.** The spatial pairwise term of the gradient can be cal-
 250 culated efficiently using standard 2D convolution. In addition to giving us an efficient way
 251 of performing the forward pass we can also utilize the 2D convolution layer to perform the
 252 backward pass, calculating the error derivatives with respect to the input and parameters.
 253 Similar to the spatial term, the bilateral term is also calculated utilizing a bilateral filtering
 254 technique. Jampani *et al.* [29] also presented a way to calculate the error derivatives with
 255 respect to the parameters for an arbitrary shaped bilateral filter.

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

263 **4.3. Entropic Descent Update.** The entropic descent step is done individually for each
 264 pixel. Since we have the update step on closed form we can easily implement it as a layer in
 265 a deep learning framework. Regarding the error derivative we are required to calculate both
 266 the error derivatives with respect to the values of the previous iteration, q^t and with respect
 267 to gradient, $\nabla_{q^t} E$. The error derivatives with respect to the values of the previous iteration

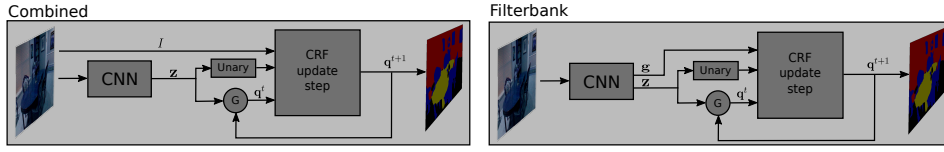


Figure 2. The data flow of the deep structure model. Each rectangle or circle represent an operation that can be performed within a deep learning framework. Note that the CNN outputs both class probabilities \mathbf{z} and filterbank features \mathbf{g} for the filterbank version.

268 are given according to

$$269 \quad (4.1) \quad \frac{\partial L}{\partial q_{ij}^k} = \frac{q_{ij}^{k+1}}{q_{ij}^k} \left(\frac{\partial L}{\partial q_{ij}^{k+1}} - \sum_{l=1}^n \frac{\partial L}{\partial q_{il}^{k+1}} q_{il}^{k+1} \right),$$

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

$$273 \quad (4.2) \quad \frac{\partial L}{\partial y_{ij}} = -t_k q_{ij}^{k+1} \left(\frac{\partial L}{\partial q_{ij}^{k+1}} - \sum_{l=1}^n \frac{\partial L}{\partial q_{il}^{k+1}} q_{il}^{k+1} \right).$$

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

275 **5. Recurrent Formulation of Deep Structured Model.** Our iterative solution to the
 276 CRF energy minimisation problem by projected gradient descent, as described in the previous
 277 sections, is formulated as a Recurrent Neural Network (RNN). The input to the RNN is the
 278 image, and the outputs of the CNN, as shown in Fig. 2. The Unary CNN’s output, \mathbf{z} , are
 279 the unary potentials and obtained after the final softmax layer (since the CNN is initially
 280 trained for classification). For the filterbank version the CNN also outputs image-dependent
 281 features, \mathbf{g} , which are ”selecting” which filters to use to compose the pairwise term at each
 282 pixel location.

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

$$285 \quad (5.1) \quad \mathbf{q}^{t+1} = f(\mathbf{q}^t, \mathbf{z}, I, \mathbf{w}).$$

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

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

294 The parameters of the RNN are the filter weights for the pairwise kernels, and also the
 295 weight for the unary terms. Since we are able to compute error derivatives with respect to

296 the parameters, and input of the RNN, we can backpropagate error derivatives through our
 297 RNN to the preceding CNN and train our entire network end-to-end. Furthermore, since the
 298 operations of the RNN are formulated as filtering, training and inference can be performed in
 299 a fully-convolutional manner.

300 The CNN part of our network allows us to leverage the ability of CNNs to learn rich
 301 feature representations from data, whilst the RNN part of the network utilises the CRF’s
 302 ability to model output structure. As we learn the parameters of our pairwise terms, we are
 303 not restricted to Gaussian potentials as in [33, 48], and we show the benefits of this in our
 304 experiments (Section 7).

305 **6. Implementation Details.** Our proposed CRF model has been implemented in the
 306 CAFFE [30] library. The Unary CNN part of our model is initialized from a pre-trained
 307 segmentation network. For all experiments we use the Deeplab-LargeFOV proposed by Chen
 308 *et al.* [18]. For the *combined* version of our model the unary CNN is pre-trained for pixel-wise
 309 classification.

310 For the *filterbank* version we use a modified version of the Deeplab-LargeFOV where a
 311 second head is added to the the network as in [17]. This head is formed by upsampling
 312 and concatenating several intermediate layers of the original network, a final convolution are
 313 applied to the concatenated features and lastly a softmax layer is added. The second head
 314 outputs the filter choosing features \mathbf{g}_f and is pre-trained to classify each pixel as horizontal
 315 semantic edge, vertical semantic edge or no edge. This part of the base network can also be
 316 trained during the final end-to-end training.

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

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

326 The spatial weights of the CRF model are all initialized as zero with the motivation that
 327 we did not want to impose a shape for these filters, but instead see what was learned during
 328 training. The bilateral filters were initialized as Gaussians with the common Potts class
 329 interaction (the filters corresponding to interactions between the same class were set to zero)
 330 [34, 14, 48].

331 **7. Experiments.** We evaluate the proposed approach on three datasets: WEIZMANN
 332 HORSE dataset [10], NYU V2 geometric dataset [46] and PASCAL VOC 2012 [23]. In these
 333 experiments, we show that the proposed approach, has advantages over similar approaches
 334 such as CRF-RNN [48]. In addition we show that adding a CRF-model as proposed in this
 335 paper improves the results on strong unary CNN networks, even for cases where the CNN has
 336 been trained on large amounts of extra data.

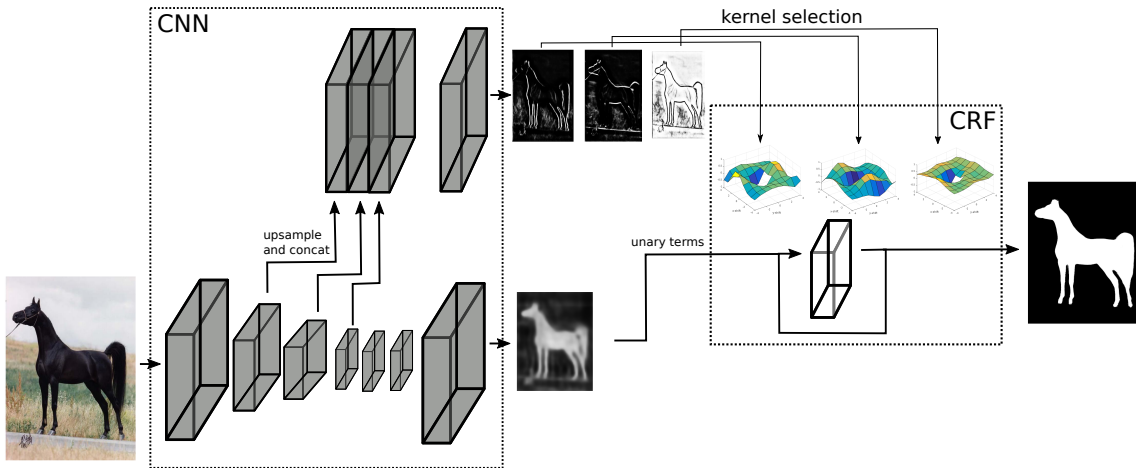


Figure 3. Schematic of the filterbank version of our model. The CNN part outputs initial class probability maps as well as filter selection maps. The structure used for the CNN is a modified version of the Deeplab-LargeFOV [18] with an extra added head.

Method	mIoU (%)
Unary CNN - Deeplab [14]	90.89
CRF-RNN [48]	91.47
Gaussian-ED	92.64
Combined-ED	92.99
Combined-MF	92.73
Combined-PGD	92.79
Filterbank-ED	93.22

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.

337 **7.1. Weizmann Horse.** The WEIZMANN HORSE dataset is widely used for benchmarking
 338 object segmentation algorithms. It contains 328 images of horses in different environments.
 339 We divide these images into a training set of 150 images, a validation set of 50 images and
 340 a test set of 128 images. Our purpose is to verify our ability to learn reasonable kernels and
 341 study the effects of different settings on a relatively small dataset. In addition we use this
 342 dataset to evaluate our proposed inference method as well as the different types of CRF-
 343 models. To compare the different types of inference methods train our *combined* model with
 344 three types of inference methods: Entropic Descent (ED), Projected Gradient Descent (PGD)
 345 and Mean Field (MF). We also train a version with only Gaussian potentials (using the same
 346 potentials as for CRF-RNN [48]). We also trained and evaluated the *filter-bank* version. The
 347 results are summarized in Table 1 and some example segmentations are shown in Fig. 4. As
 348 can be seen from the results, our proposed inference method using entropic descent achieves
 349 slightly better results on the test set for the *combined* CRF model. However, the increase

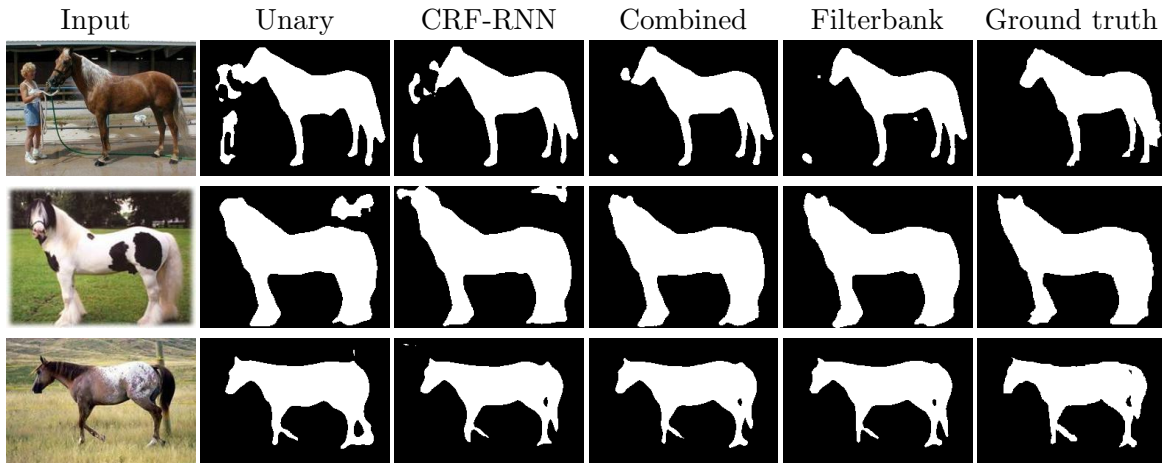


Figure 4. Qualitative results on the WEIZMANN HORSE dataset. Note that the proposed methods capture the shape of the horses better than the baselines, especially compared to the unary network.

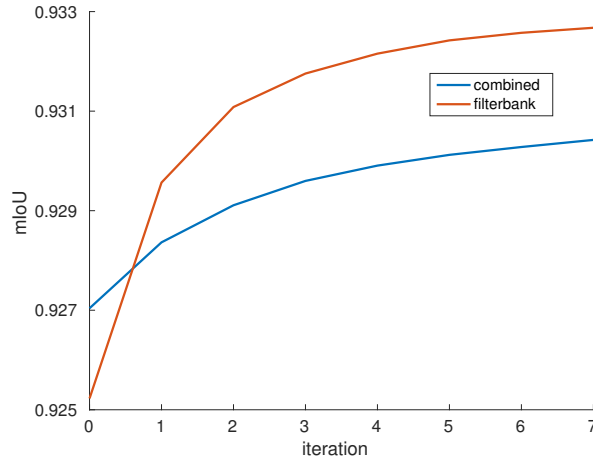


Figure 5. WEIZMANN HORSE test set results in terms of mean Intersection over Union plotted as a function of the number of iterations for the CRF inference method. During training the number of iterations were set to five.

350 over mean field and projected gradient descent inference is minor. For the case where we
 351 used Gaussian CRF potentials we get better results with entropic descent inference compared
 352 to mean field. Comparing entropic descent inference and projected gradient descent the two
 353 methods achieve similar results. Training a model with projected gradient descent is however
 354 problematic due to the zeroing of gradients, to solve this we train with a "leaky" version of
 355 projected gradient descent. This means that the intermediate states and final results might
 356 not lie on the probability simplex, something that is guaranteed for entropic descent inference.

357 In Fig. 5 the mean intersection over union on the test set is plotted as a function of the
 358 number of CRF inference iterations. During training the number of iterations were set to five.
 359 As can be seen in the figure, increasing the number of inference step will only slightly increase

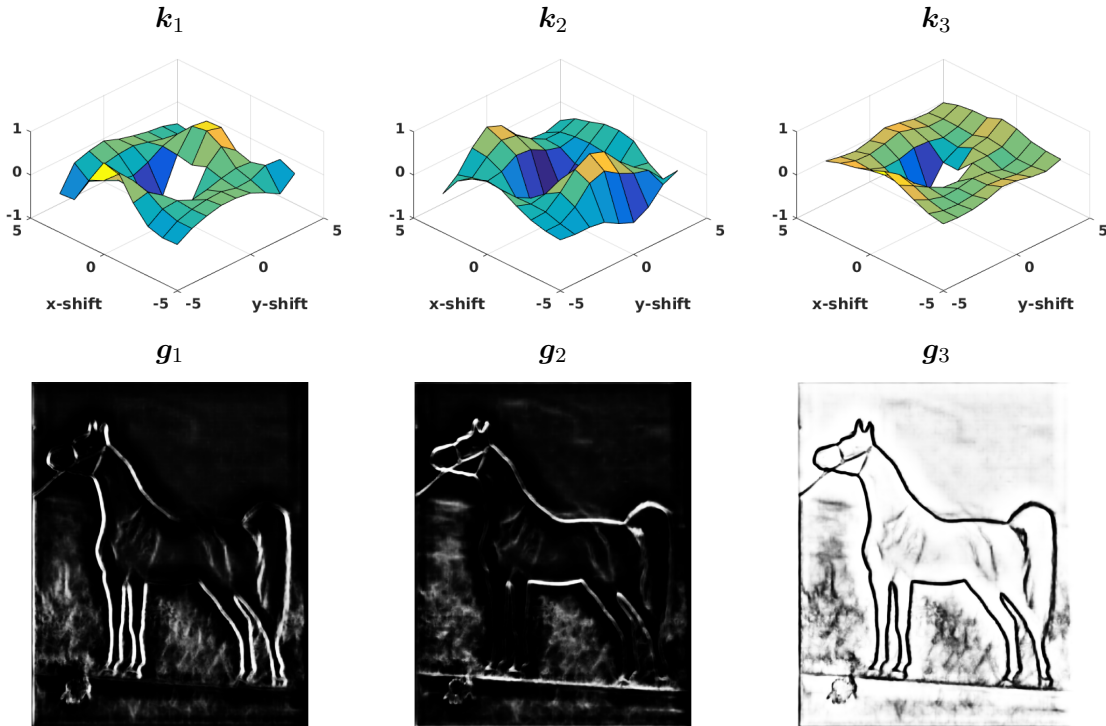


Figure 6. Visualization of the pairwise kernel weights for the filterbank version trained on the WEIZMANN HORSE data set. These weights are for the classes "background" and "horse", the plots can be understood as the energy added when assigning the pixels with the relative positions (x,y) and $(x+x\text{-shift},y+y\text{-shift})$ 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 behaviour 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 can hence not be interpreted in the same way.

360 the segmentation result. In Fig. 6 the pairwise weights of the *filterbank* version is visualized.

361 **7.2. NYU V2.** The NYU V2 dataset contains images taken by Microsoft Kinect V-1
 362 camera in 464 indoor scenes. We use the official training and validation splits consisting of
 363 795 and 654 images, respectively. Following the setting described in Wang *et al.* [46], we also
 364 include additional images for training. These are the images from the NYU V1 dataset that
 365 do not overlap with the images in the official validation set. This gives a total of 894 images
 366 with semantic label annotations for training. As in [46] we consider 5 classes conveying strong
 367 geometric properties: ground, vertical, ceiling, furniture and objects.

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

370 **7.3. PASCAL VOC.** The PASCAL VOC 2012 segmentation benchmark [22] consists of 20
 371 foreground and one background class. The unary network used for these experiments is again

Method	mIoU (%)
R-CNN [25]	40.3
Semantic HCRF [46]	42.7
Joint HCRF [46]	44.2
Modular CNN [28]	54.3
Unary CNN - Deeplab [14]	62.8
CRF-RNN [48]	64.4
Combined	65.4
Filterbank	65.4

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.

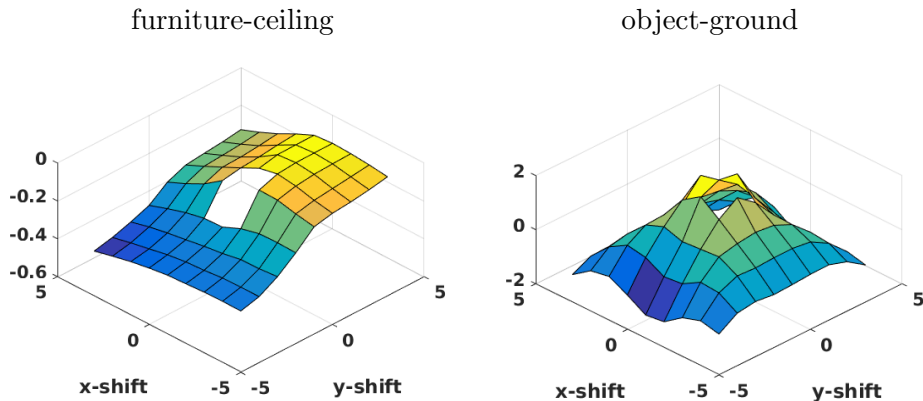


Figure 7. Visualization of the pairwise kernel weights for the filterbank version trained on the NYU V2 data set. These weights are for the classes shown above the plots and for the third filter selection map which usually has a high value for pixels with no semantic edge. The furniture-ceiling kernel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

372 the Deeplab-LargeFOV network [18], this network has been pretrained on the MS-COCO 2014
373 dataset [38] and then trained on the PASCAL VOC training data as well as a training set
374 created from annotations of the semantic boundaries dataset [26]. We add our CRF-models to
375 this baseline network and train only on the PASCAL VOC training data. This to show that
376 we can improve upon really strong baselines, even though we finetune the complete models
377 on only a fraction of the training data used for the baseline. The results for the PASCAL
378 VOC 2012 validation set is shown in Table 3. In addition we evaluate our model on the test
379 set, for this the results are shown in Table 4. As can be seen, our models perform similar
380 to models trained with the same base network. Note that our models are only trained on
381 the training data during end-to-end training. Recently there have been several CNNs with
382 different base architectures presented that perform well, even without a CRF. The top entry
383 at the moment is PSPNet [36] with a mIoU of 85.4, we leave it to future work to explore

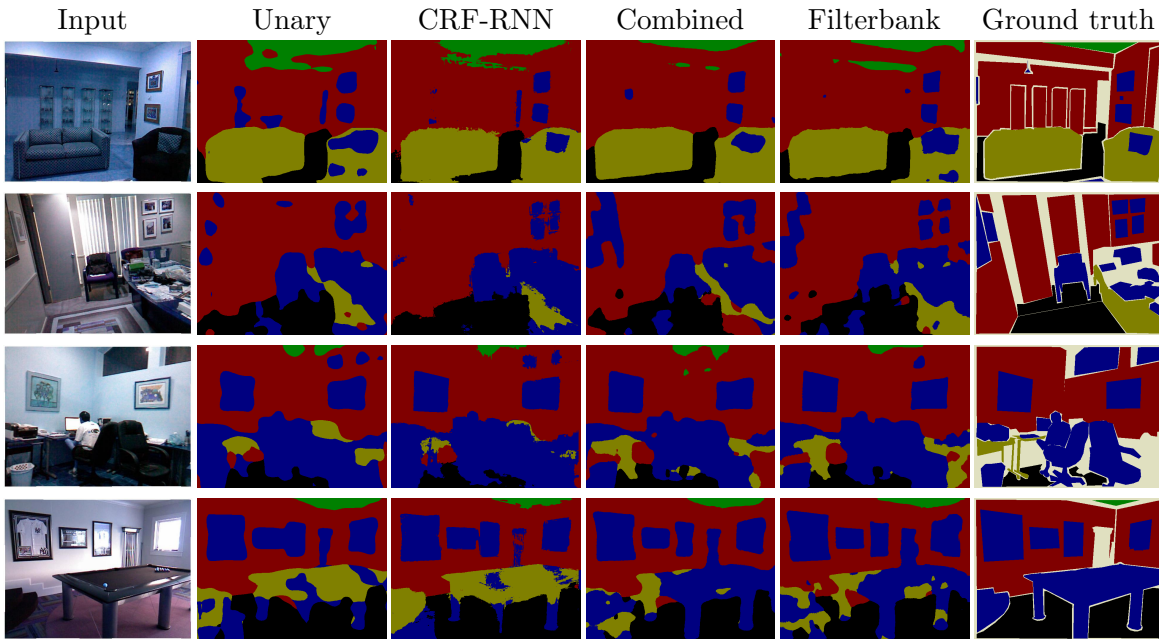


Figure 8. Qualitative results on the NYU V2 dataset. Note that the proposed methods captures the shape of the object instances better than the baselines. This effect is perhaps most pronounced for the paintings hanging on the walls. The pixels colored off-white are "ignore"-pixels, these are not counted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

Method	mIoU (%)
Unary CNN - Deeplab [14]	68.5
CRF-RNN [48]	71.7
Combined	72.0
Filterbank	70.1

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 [38].

384 whether these architectures can be improved using our proposed methodology.

385 **7.4. Execution time.** We also investigated the difference in running time between the
 386 two proposed models. This was done on a computer with a Nvidia Titan X GPU with Pascal
 387 architecture and an Intel i7-5930K processor. The implementation used for the bilateral
 388 filtering used was the one from Jampani *et al.* [29] where most of the computations are done
 389 on the GPU. The initialization of the permutohedral lattice is however done on the CPU.
 390 The runtimes were tested by performing the forward step for a randomized RGB image of
 391 size 640×640 with 21 classes. The numbers presented are the average of 100 runs. For
 392 the *combined* model the forward runtime was 12 seconds while for the *filterbank* it was 0.37

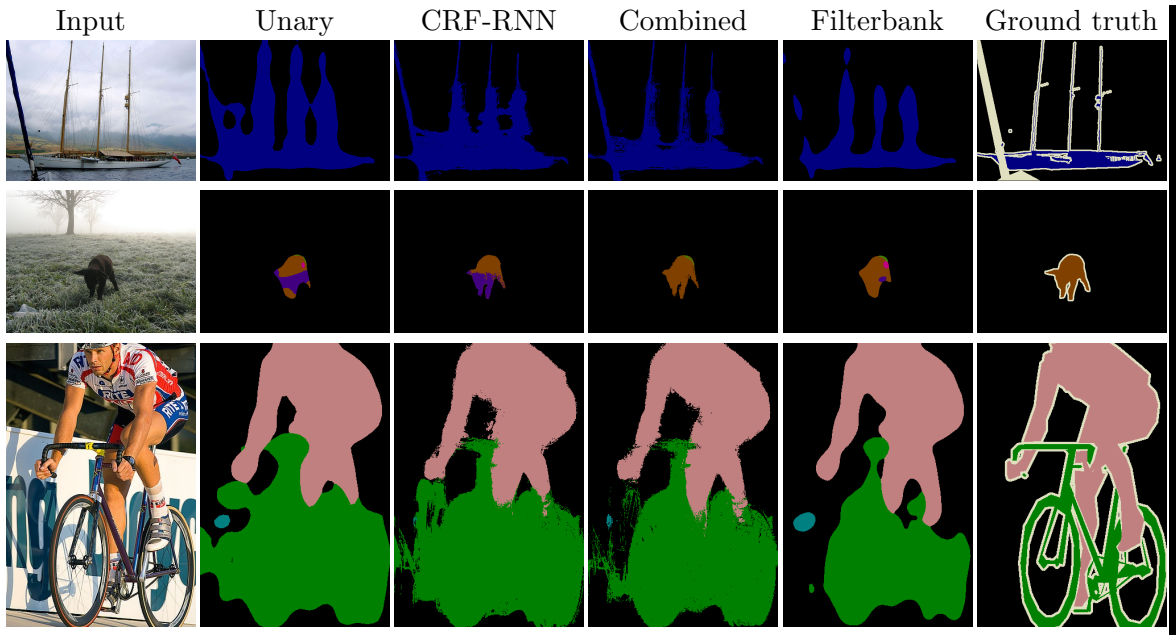


Figure 9. Qualitative results on the PASCAL VOC dataset. [22]. The pixels colored off-white are "ignore"-pixels, these are not counted in the evaluation. The training images have similar "ignore"-pixels.

Method	mIoU (%)
Unary CNN - Deeplab [14]	68.9
DT-EdgeNet [17]	71.7
CRF-RNN [48]	72.2
Combined	72.5
Filterbank	69.5

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 [38], but note that our models were not trained using MS-COCO.

393 second.

394 **8. Conclusion.** In this paper we have presented a gradient descent based method for inference
 395 in Conditional Random Fields. This method allows for backpropagation of error deriva-
 396 tive hence enabling end-to-end training with an Convolutional Neural Network of choice. We
 397 show that this inference method has beneficial properties and performs better on some tasks
 398 compared to other methods such as mean field. In addition, we present two types of Condi-
 399 tional Random Field models tailored for semantic segmentation. The *combined* model that
 400 uses spatial pairwise terms as well as image-dependent bilateral pairwise terms. This model
 401 performs well but is somewhat computational expensive due to the high dimensionality of the
 402 bilateral filtering. We also present the *filterbank* model which is also image dependent but

403 only requires 2D convolutions during inference. The image dependence of the model comes
404 from an output map of the base CNN that acts as a "filter choosing" map. This enables the
405 model to, for example, use one filter representing pairwise interaction between pixel labels at
406 semantic edges and another filter far away from semantic edges. This model gives a speedup
407 by a factor of 32 compared to the *combined* model without losing performance in terms
408 of segmentation quality. For the smaller dataset it achieves similar segmentation quality as
409 the *combined* model. Since the *filterbank* version of the model learns how the pairwise term
410 should depend on the image it is in this aspect more expressive than the *combined* version.
411 The pairwise terms of the *combined* model is however, due to its dependence on the color
412 gradient, hand-crafted to preserve and refine edges. This is beneficial for the PASCAL VOC
413 dataset where the unary network generally capture the context well but outputs "blobby"
414 segmentations. For all the models presented the pairwise kernels can have arbitrary shape,
415 instead of commonly used Gaussian kernels. This enables the models to learn more compli-
416 cated pairwise label interactions.

417

418

419 *Acknowledgements.* This work has been funded by the Swedish Research Council (grant
420 no. 2016-04445), the Swedish Foundation for Strategic Research (Semantic Mapping and
421 Visual Navigation for Smart Robots), Vinnova / FFI (Perceptron, grant no. 2017-01942),
422 ERC grant ERC-2012-AdG 321162-HELIOS, EPSRC grant Seebibyte EP/M013774/1 and
423 EPSRC/MURI grant EP/N019474/1.

424

REFERENCES

- 425 [1] A. ADAMS, J. BAEK, AND M. A. DAVIS, *Fast high-dimensional filtering using the permutohedral lattice*,
426 Computer Graphics Forum, (2010).
- 427 [2] T. AJANTHAN, A. DESMAISON, R. BUNEL, M. SALZMANN, P. H. TORR, AND M. P. KUMAR, *Efficient*
428 *linear programming for dense crfs*.
- 429 [3] A. ARNAB, S. JAYASUMANA, S. ZHENG, AND P. H. S. TORR, *Higher order conditional random fields in*
430 *deep neural networks*, in European Conference on Computer Vision, 2016.
- 431 [4] A. ARNAB, S. ZHENG, S. JAYASUMANA, B. ROMERA-PAREDES, M. LARSSON, A. KIRILLOV, B. SAVCHYN-
432 SKYY, C. ROTHER, F. KAHL, AND P. H. S. TORR, *Conditional random fields meet deep neural*
433 *networks for semantic segmentation: Combining probabilistic graphical models with deep learning for*
434 *structured prediction*, IEEE Signal Processing Magazine, 35 (2018), pp. 37–52, <https://doi.org/10.1109/MSP.2017.2762355>.
- 435 [5] A. BECK AND M. TEBoulLE, *Mirror descent and nonlinear projected subgradient methods for convex*
436 *optimization*, Operations Research Letters, 31 (2003), pp. 167–175.
- 437 [6] D. BELANGER AND A. MCCALLUM, *Structured prediction energy networks*, in International Conference
438 on Machine Learning, 2016.
- 439 [7] D. BELANGER, B. YANG, AND A. MCCALLUM, *End-to-end learning for structured prediction energy*
440 *networks*, arXiv preprint arXiv:1703.05667, (2017).
- 441 [8] G. BERTASIUS, L. TORRESANI, S. X. YU, AND J. SHI, *Convolutional random walk networks for semantic*
442 *image segmentation*, in The IEEE Conference on Computer Vision and Pattern Recognition (CVPR),
443 July 2017.
- 444 [9] A. BLAKE, P. KOHLI, AND C. ROTHER, *Markov Random Fields for Vision and Image Processing*, MIT
445 Press, 2011.
- 446 [10] E. BORENSTEIN AND S. ULLMAN, *Class-specific, top-down segmentation*, in European Conference on
447 Computer Vision, 2002.
- 448 [11] E. BOROS AND P. HAMMER, *Pseudo-boolean optimization*, Discrete Appl. Math., 123 (2002), pp. 155–225.
- 449 [12] L. BOTTOU, Y. BENGIO, AND Y. LE CUN, *Global training of document processing systems using graph*
450 *transformer networks*, in IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 1997,
451 pp. 489–494.
- 452 [13] S. CHANDRA AND I. KOKKINOS, *Fast, exact and multi-scale inference for semantic image segmentation*
453 *with deep gaussian crfs*, in European Conference on Computer Vision, 2016.
- 454 [14] L. CHEN, G. PAPANDREOU, I. KOKKINOS, K. MURPHY, AND A. L. YUILLE, *Semantic image segmenta-*
455 *tion with deep convolutional nets and fully connected crfs*, in International Conference on Learning
456 Representations, 2015.
- 457 [15] L. CHEN, G. PAPANDREOU, I. KOKKINOS, K. MURPHY, AND A. L. YUILLE, *Deeplab: Semantic image*
458 *segmentation with deep convolutional nets, atrous convolution, and fully connected crfs*, arXiv preprint
459 arXiv:1606.00915, (2016).
- 460 [16] L. CHEN, A. SCHWING, A. YUILLE, AND R. URTASUN, *Learning deep structured models*, in Int. Conf.
461 Machine Learning, Lille, France, 2015.
- 462 [17] L.-C. CHEN, J. T. BARRON, G. PAPANDREOU, K. MURPHY, AND A. L. YUILLE, *Semantic image segmen-*
463 *tation with task-specific edge detection using cnns and a discriminatively trained domain transform*, in
464 Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 4545–
465 4554.
- 466 [18] L.-C. CHEN, G. PAPANDREOU, I. KOKKINOS, K. MURPHY, AND A. L. YUILLE, *Deeplab: Semantic*
467 *image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs*, IEEE
468 Transactions on Pattern Analysis and Machine Intelligence, (2017).
- 469 [19] Y. CHEN AND X. YE, *Projection onto a simplex*, arXiv preprint arXiv:1101.6081, (2011).
- 470 [20] A. DESMAISON, R. BUNEL, P. KOHLI, P. H. S. TORR, AND M. P. KUMAR, *Efficient continuous relax-*
471 *ations for dense CRF*, in European Conference on Computer Vision, 2016.
- 472 [21] J. DOMKE, *Learning graphical model parameters with approximate marginal inference*, IEEE Transactions
473 on Pattern Analysis and Machine Intelligence, 35 (2013), pp. 2454–2467.
- 474 [22] M. EVERINGHAM, S. M. A. ESLAMI, L. VAN GOOL, C. K. I. WILLIAMS, J. WINN, AND A. ZISSER-
475 MAN, *The pascal visual object classes challenge: A retrospective*, International Journal of Computer
476

- 477 Vision, 111 (2015), pp. 98–136, <https://doi.org/10.1007/s11263-014-0733-5>, <https://doi.org/10.1007/s11263-014-0733-5>.
- 478
- 479 [23] M. EVERINGHAM, L. V. GOOL, C. K. I. WILLIAMS, J. WINN, AND A. ZISSERMAN, *The pascal visual*
480 *object classes (voc) challenge*, Int. Journal Computer Vision, 88 (2010), pp. 303–338.
- 481 [24] G. GHIASI AND C. FOWLKES, *Laplacian reconstruction and refinement for semantic segmentation*, in
482 European Conference on Computer Vision, 2016.
- 483 [25] R. GIRSHICK, J. DONAHUE, T. DARRELL, AND J. MALIK, *Rich feature hierarchies for accurate object*
484 *detection and semantic segmentation*, in IEEE Conference on Computer Vision and Pattern Recog-
485 nition, 2014.
- 486 [26] B. HARIHARAN, P. ARBELEZ, L. BOURDEV, S. MAJI, AND J. MALIK, *Semantic contours from inverse*
487 *detectors*, in 2011 International Conference on Computer Vision, Nov 2011, pp. 991–998, <https://doi.org/10.1109/ICCV.2011.6126343>.
- 488
- 489 [27] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in IEEE Conference
490 on Computer Vision and Pattern Recognition, 2016.
- 491 [28] O. H. JAFARI, O. GROTH, A. KIRILLOV, M. Y. YANG, AND C. ROTHER, *Analyzing modular cnn architec-*
492 *tures for joint depth prediction and semantic segmentation*, in International Conference on Robotics
493 and Automation, 2017.
- 494 [29] V. JAMPANI, M. KIEFEL, AND P. V. GEHLER, *Learning sparse high dimensional filters: Image filter-*
495 *ing, dense crfs and bilateral neural networks*, in IEEE Conference on Computer Vision and Pattern
496 Recognition, June 2016.
- 497 [30] Y. JIA, E. SHELHAMER, J. DONAHUE, S. KARAYEV, J. LONG, R. GIRSHICK, S. GUADARRAMA,
498 AND T. DARRELL, *Caffe: Convolutional architecture for fast feature embedding*, arXiv preprint
499 arXiv:1408.5093, (2014).
- 500 [31] A. KIRILLOV, D. SCHLESINGER, S. ZHENG, B. SAVCHYNSKY, P. TORR, AND C. ROTHER, *Joint training*
501 *of generic cnn-crf models with stochastic optimization*, in Asian Conference on Computer Vision,
502 2016.
- 503 [32] D. KOLLER AND N. FRIEDMAN, *Probabilistic Graphical Models*, MIT Press, 2009.
- 504 [33] P. KRAEHNBUHL AND V. KOLTUN, *Parameter learning and convergent inference for dense random*
505 *fields*, in Proceedings of The 30th International Conference on Machine Learning, 2013, pp. 513–521.
- 506 [34] P. KRÄHENBÜHL AND V. KOLTUN, *Efficient inference in fully connected CRFs with gaussian edge poten-*
507 *tials*, in Neural Information Processing Systems, 2011.
- 508 [35] M. LARSSON, A. ARNAB, F. KAHL, S. ZHENG, AND P. H. S. TORR, *A projected gradient descent method*
509 *for crf inference allowing end-to-end training of arbitrary pairwise potentials*, in 11th International
510 Conference on Energy Minimization Methods in Computer Vision and Pattern Recognition, (EMM-
511 CVPR), Springer, 2017.
- 512 [36] G. LIN, A. MILAN, C. SHEN, AND I. REID, *Refinenet: Multi-path refinement networks with identity*
513 *mappings for high-resolution semantic segmentation*, arXiv preprint arXiv:1611.06612, (2016).
- 514 [37] G. LIN, C. SHEN, A. HENGEL, AND I. REID, *Efficient piecewise training of deep structured models for*
515 *semantic segmentation*, in IEEE Conference on Computer Vision and Pattern Recognition, June 2016.
- 516 [38] T.-Y. LIN, M. MAIRE, S. BELONGIE, J. HAYS, P. PERONA, D. RAMANAN, P. DOLLÁR, AND
517 C. L. ZITNICK, *Microsoft COCO: Common Objects in Context*, Springer International Publishing,
518 Cham, 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.
- 519
- 520 [39] Z. LIU, X. LI, P. LUO, C. C. LOY, AND X. TANG, *Semantic image segmentation via deep parsing network*,
521 in International Conference on Computer Vision, 2015.
- 522 [40] J. LONG, E. SHELHAMER, AND T. DARRELL, *Fully convolutional networks for semantic segmentation*, in
523 IEEE Conference on Computer Vision and Pattern Recognition, 2015.
- 524 [41] S. REN, K. HE, R. GIRSHICK, AND J. SUN, *Faster R-CNN: Towards real-time object detection with region*
525 *proposal networks*, in Neural Information Processing Systems, 2015.
- 526 [42] C. ROTHER, V. KOLMOGOROV, AND A. BLAKE, *“GrabCut”: Interactive foreground extraction using*
527 *iterated graph cuts*, in ACM Transactions on Graphics, 2004, pp. 309–314.
- 528 [43] A. SCHWING AND R. URTASUN, *Fully connected deep structured networks*, in arXiv preprint
529 arXiv:1503.02351, 2015.
- 530 [44] K. SIMONYAN AND A. ZISSERMAN, *Very deep convolutional networks for large-scale image recognition*, in

- 531 International Conference on Learning Representations, 2015.
- 532 [45] V. VINEET, J. WARRELL, AND P. TORR, *Filter-based mean-field inference for random fields with higher*
533 *order terms and product label-spaces*, in European Conference on Computer Vision, 2012.
- 534 [46] P. WANG, X. SHEN, Z. LIN, S. COHEN, B. PRICE, AND A. YUILLE, *Towards unified depth and semantic*
535 *prediction from a single image*, in IEEE Conference on Computer Vision and Pattern Recognition,
536 2014.
- 537 [47] W. WANG, S. FIDLER, AND R. URTASUN, *Proximal deep structured models*, in Neural Information Pro-
538 cessing Systems, 2016.
- 539 [48] S. ZHENG, S. JAYASUMANA, B. ROMERA-PAREDES, V. VINEET, Z. SU, D. DU, C. HUANG, AND P. TORR,
540 *Conditional random fields as recurrent neural networks*, in International Conference on Computer
541 Vision, 2015.