

Fast, Exact and Multi-Scale Inference for Semantic Image Segmentation with Deep Gaussian CRFs

Anonymous ECCV submission

Paper ID 1438

Abstract. In this work we propose a combination of the Gaussian Conditional Random Field (G-CRF) with Deep Learning for the task of structured prediction. Our method inherits several virtues of G-CRF and Deep Learning: (a) the structured prediction task has a unique global optimum that is obtained exactly from the solution of a linear system (b) structured prediction can be jointly trained in an end-to-end setting in general architectures and with arbitrary loss functions, (c) the pairwise terms do not have to be simple hand-crafted expressions, as in the line of works building on the DenseCRF, but can rather be ‘discovered’ from data through deep architectures - in particular we use fully convolutional networks to obtain both the unary and pairwise terms of our G-CRF.

Building on standard tools from numerical analysis we develop very efficient algorithms for inference and learning. This efficiency allows us to explore more sophisticated architectures for structured prediction in deep learning: we introduce multi-resolution architectures to couple information across scales in a joint optimization framework, yielding systematic improvements.

We demonstrate the utility of our approach on the challenging VOC PASCAL 2012 image segmentation benchmark, where an extensive ablation study indicates substantial improvements over strong baselines.

1 Introduction

Over the last few years deep learning has resulted in dramatic progress in the task of semantic image segmentation. Early works on using CNNs as feature extractors [1,2,3] and combining them with standard superpixel-based front-ends gave substantial improvements over well-engineered approaches that used hand-crafted features. The currently mainstream approach is relying on Fully Convolutional Networks (FCNs), introduced into the task in [4], where networks originally trained for classification were retrained to provide fields of outputs used for pixelwise labeling.

A dominant research direction for improving semantic segmentation with deep learning is the combination of the powerful classification capabilities of FCNs with structured prediction, which aims at improving classification by capturing interactions between predicted labels. One of the first works in the direction of combining deep networks with structured prediction was [5] which advocated the use of densely-connected conditional random fields (DenseCRF) [6] to post-process an FCNN output so as to obtain a sharper segmentation that preserves image boundaries. This was then used by Zheng *et al.* [7] who combined DenseCRF with a CNN into a single Recurrent Neural Network (RNN), accommodating the DenseCRF post processing in an end-to-end training procedure.

All of the currently leading techniques for semantic segmentation perform structured prediction using approximate inference and learning. For instance the techniques of [5,7] perform mean-field inference for a fixed number of iterations (10), and convergence (even to a local minimum) is not guaranteed. The current state-of-the-art in semantic segmentation, [8] uses piecewise training and three iterations of mean-field inference, while the runner-up, [9] uses highly-sophisticated modules, effectively learning to approximate mean-field inference. In all of these works a more pragmatic approach is taken to inference, considering it as a sequence of operations that need to be learned, as proposed in [7]. This ‘infering’ approach of combining learning and inference may be liberating, in the sense that one acknowledges and accommodates the approximations in the inference through end-to-end training. But it would certainly be preferable to avoid approximations if this can be done efficiently, without sacrificing expressive power.

Motivated by [10,11], our starting point in this work is the observation that a particular type of graphical model, the Gaussian Conditional Random Field (G-CRF), allows us to perform exact and efficient Maximum-A-Posteriori inference. The log-likelihood of the G-CRF posterior has the form of a quadratic energy function which captures unary and pairwise interactions between random variables. There are two advantages to using a quadratic function: (a) unlike the energy of general graphical models, a quadratic function has a unique global minimum if the system matrix is positive definite, and (b) this unique minimum can be efficiently found by solving a system of linear equations. We can actually discard the probabilistic underpinning of the G-CRF and understand G-CRF inference as an energy-based model, casting structured prediction as quadratic optimization.

Our work aims at bringing this insight into deep learning; as in [8,9] we augment FCNs with discriminatively trained convolutional layers that model and enforce pairwise consistencies between neighbouring regions. But unlike [8,9,5,7], we do not use approximate inference; we show that a global optimum can be efficiently recovered in a small number (typically 10) iterations of conjugate gradients, with a negligible cost (0.02 seconds) when implemented on the GPU. Back-propagating also amounts to (exactly) solving a linear system.

The connection with linear system solvers allows us to leverage on the wealth of tools from numerical analysis. Firstly, as we prove in the supplemental material, for Gaussian CRFs sequential/parallel mean-field inference amounts to solving a linear system using the classic Gauss-Seidel/Jacobi algorithms respectively. So the currently default algorithm for inference would amount to a naive method for solving our task, when compared to conjugate gradients.

Secondly, having a toolkit for exact inference allows us to perform a systematic comparison of different graph topologies, without confounding the effects of approximate inference and representational power. Rather than using the traditional 4-connected neighbourhoods, our framework allows us the freedom to work with arbitrary neighbourhoods. Most importantly, we explore the merit of using multi-scale networks, where variables computed from different image scales interact with each other. This gives rise to a flow of information across different-sized neighborhoods. We show experimentally that this yields substantially improved results over single-scale baselines.

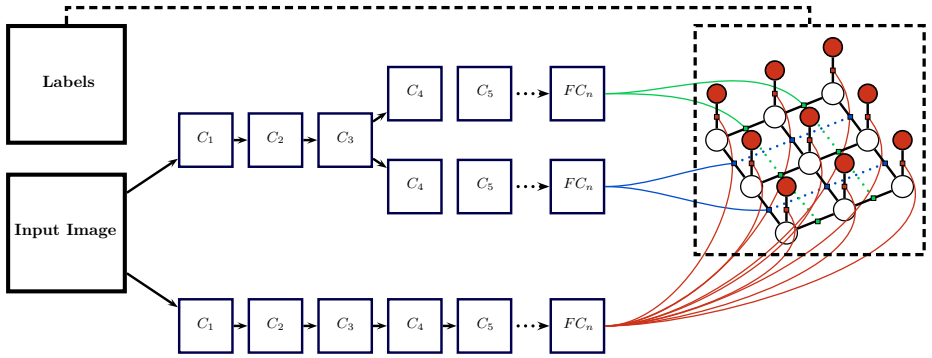


Fig. 1: A detailed schematic representation of our fully convolutional neural network with a Quadratic Optimization module. The network has a unary stream which populates unary terms, and a pairwise stream which populates horizontal and vertical pairwise terms. The unit outlined by dotted lines represents the quadratic optimization module. The quadratic optimization module collects unary and pairwise terms from the network and proposes an image hypothesis after inference. Additionally, it computes the error gradients and back propagates them through the network during the training phase. The factor graph inside the quadratic optimization module shows a 3×3 image, with a 4-connected neighbourhood. The white blobs represent the pixels, or random variables. The red blobs represent unary factors, and these are fed by a convolutional layer from the unary stream of the network. The green and blue squares represent vertical and horizontal connectivity factors, and these are fed by convolutional layers from the pairwise stream of the network. Please see in colour.

We demonstrate the merit of our approach alongside strong, publicly available baselines. In particular, we use a modified variant of the Deeplab-LargeFOV network structure proposed in [5], extending it with a Gaussian Conditional Random Field with discriminatively trained unary and pairwise terms. We compare with the improvements attained by using DenseCRF, as well as with the improvements obtained by combining our inference with DenseCRF in a cascade, demonstrating substantial improvements. We cannot directly compete with [8] and [9], since we do not use their additional (and orthogonal) extensions for higher-resolution networks and global context. However, our ablation study on the PASCAL validation set indicates the merits of each of the contributions outlined above, while we obtain competitive preliminary results on the PASCAL test set; the final results using converged networks will be in the supplemental material.

Finally, our implementation is efficient, fully GPU based, and built on top of the popular *Caffe* library [12] - we intend to make the distribution available if the paper gets accepted.

In Sec. 2 we describe our approach in detail, and derive the expressions for weight update rules for parameter learning. In Sec. 3 we analyze the efficiency of the linear system solvers and present our multi-resolution structured prediction algorithm. In Sec. 4 we report consistent improvements over well-known baselines and state-of-the-art results on the PASCAL VOC test set.

2 Quadratic Optimization Formulation

We now describe our formulation. Consider an image \mathcal{I} containing P pixels. Each pixel $p \in \{p_1, \dots, p_P\}$ can take a label $l \in \{1, \dots, L\}$. Although our objective is to assign discrete labels to the pixels, we phrase our problem as a continuous inference task. Rather than performing a discrete inference task that delivers one label per variable, we use a continuous function of the form $\mathbf{x}(p, l)$ which gives an affinity score for each (pixel, label) pair. This score can be understood as being proportional to the log-likelihood of the pixel p taking the label l , if a softmax unit is used to process \mathbf{x} .

We denote the pixel-level ground-truth labeling by $\mathbf{y} \in \mathbb{R}^P$, and the inferred hypothesis by $\mathbf{x} \in \mathbb{R}^N$, where $N = P \times L$. Our formulation is posed as an energy minimization problem. In the following subsections, we describe the form of the energy function, the inference procedure, and the parameter learning approach, followed by some technical details pertinent to using our framework in a fully convolutional neural network.

2.1 Energy of a hypothesis

We define the energy of a hypothesis by a function of the following form:

$$E(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T (A + \lambda \mathbf{I}) \mathbf{x} - B \mathbf{x} \quad (1)$$

where A denotes the symmetric $N \times N$ matrix of pairwise terms, and B denotes the $N \times 1$ vector of unary terms. We can understand $E(\mathbf{x})$ as being the posterior log-likelihood of the labelling \mathbf{x} according to a gaussian CRF with pairwise terms $U_{i,j}(x_i, x_j) = \exp(-\frac{1}{2} x_i (A_{i,j} + \lambda d_{i,j}) x_j)$ and unary terms $V_i(x_i) = \exp(-B_i x_i)$.

To make the matrix A strictly positive definite, we add to it λ times the Identity Matrix \mathbf{I} . Equation 1 is a standard way of expressing the energy of a system with unary and pair-wise interactions among the random variables [11] in a vector labeling task. We chose this function primarily because it has a unique global minimum and allows for exact inference, alleviating the need for approximate inference.

Figure 2 illustrates a toy example, showing the construction procedure of matrices A and B . In this example we consider a 3×3 input image containing white and black pixels. The goal here is to model discontinuities in pixel colour among neighbouring pixels. The matrix B , in Fig. 2(a) containing the unary terms is red for black, and blue for white pixels. The matrix A , shown in Fig. 2(c), contains the pairwise terms between every pair of pixels. We assume a 4-connected neighbourhood. The pairwise term is shown to be blue for any two pixels which have the same colour, and red if they have different colours. A white pairwise term indicates the two pixels are not neighbours of each other. Kindly note that the width and height of matrices A and B respectively are N -dimensional, where N is the number of pixels times the number of labels. This is a one-label toy example, hence N is the number of pixels. As shown in figure 1 in our case, these matrices, A and B are outputs of convolutional layers of a deep network.

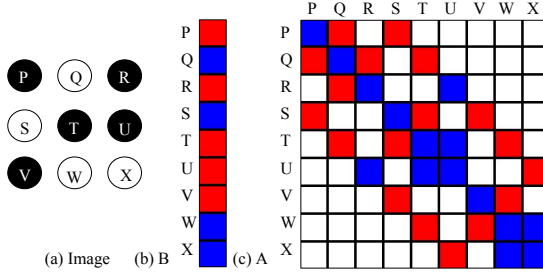


Fig. 2: Toy example illustrating the matrices A and B showing up in our cost function. The input image (a) contains black and white pixels. (b) shows matrix B , containing the unary terms: red for black pixels, blue otherwise. (c) illustrates A , containing the pairwise terms for a 4-connected neighbourhood: blue for two pixels having the same colour, red for different colours, white if pixels are not neighbours.

2.2 Inference

Given A and B , inference involves solving for the value of \mathbf{x} that minimizes the energy function in 1. If $(A + \lambda \mathbf{I})$ is symmetric positive definite, then $E(\mathbf{x})$ has a unique global minimum [13] at:

$$(A + \lambda \mathbf{I})\mathbf{x} = B. \quad (2)$$

As such, inference is exact and efficient, only involving a system of linear equations.

2.3 Learning A and B

In this work, we learn model parameters A and B by stochastic gradient descent. The loss between the prediction \mathbf{x} and the ground truth \mathbf{y} is denoted by $\mathcal{L}(\mathbf{x}, \mathbf{y})$, and we denote the derivative of the loss with respect to the prediction by $\nabla_{\mathcal{L}}$:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} \equiv \nabla_{\mathcal{L}}. \quad (3)$$

While we use the softmax loss function in our experiments, theoretically, any differentiable loss function can be used in this formulation. To learn the parameters A and B via gradient descent, we need the derivatives of the loss $\mathcal{L}(\mathbf{x}, \mathbf{y})$ with respect to A and B , i.e. $\frac{\partial \mathcal{L}}{\partial A}$ and $\frac{\partial \mathcal{L}}{\partial B}$ respectively.

Derivative of Loss with respect to B The derivative of the loss with respect to B , $\frac{\partial \mathcal{L}}{\partial B}$ can be obtained by solving the system of linear equations in Eq. 4. The derivation of this expression is presented in the supplementary material.

$$(A + \lambda \mathbf{I}) \frac{\partial \mathcal{L}}{\partial B} = \nabla_{\mathcal{L}}. \quad (4)$$

Derivative of Loss with respect to A The derivative of the loss with respect to A , $\frac{\partial \mathcal{L}}{\partial A}$ is given by negative of the kronecker product of the solution \mathbf{x} and the derivative of the loss with respect to B . This is expressed in the following equation,

$$\frac{\partial \mathcal{L}}{\partial A} = -\frac{\partial \mathcal{L}}{\partial B} \otimes \mathbf{x}. \quad (5)$$

The derivation of this expression is presented in the supplementary material.

2.4 Integration in a Deep Architectures

Having described the quadratic optimization problem in detail in the previous section, we now describe its use in a fully convolutional network. Our aim is to learn the model parameters A and B from the data with a deep architecture, and use these to then infer the image hypothesis for test images. To this end, we implement a neural network layer *QO* (Quadratic Optimization) using the *Caffe* library [12] which takes as input the unary and pairwise terms, and returns as output the image hypothesis \mathbf{x} . Additionally, to allow learning of the model parameters, it computes the derivatives of the loss with respect to the model parameters, and back-propagates them through the network. This is shown in a schematic diagram in Fig. 1.

2.5 Continuous-to-Discrete Labeling

Our image hypothesis prediction, defined in Sec. 2 is a scoring function of the form $\mathbf{x}(p, l)$ which gives an indication of how likely the pixel p is to take the label l . We use the softmax function to convert these scores into probabilities, and assign the most probable label l^* to the pixel p .

3 Linear Systems for Efficient and Effective Structured Prediction

Having identified that both the inference problem in Eq. 2 and computation of pairwise gradients in Eq. 4 require the solution of a linear system of equations, we now discuss methods for accelerated inference that rely on standard numerical analysis techniques for linear systems [14,15]. Our main contributions consist in (a) using fast linear system solvers exhibiting fast convergence (Sec. 3.1) and (b) performing inference on multi-scale graphs by constructing block-structured linear systems (Sec. 3.2).

Our contributions in (a) indicate that our extremely simple linear system solvers can be up to 2.5 faster than the solutions one could get by a naive application of parallel mean-field when implemented on the GPU. Our contribution in (b) aims at accuracy rather than efficiency, and is experimentally validated in Sec. 4

3.1 Fast Linear System Solvers

The computational cost of solving the linear system of equations in Eq. 2 and Eq. 4 depends on the size of the matrix A , i.e. $N \times N$, and its sparsity pattern. In our experiments, while $N \sim 10^5$, the matrix A is quite sparse, since we deal with small

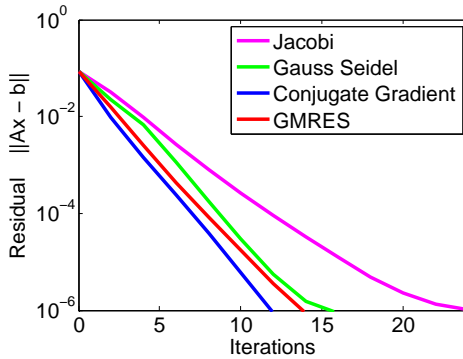


Fig. 3: A plot showing the convergence of various iterative solvers, namely Jacobi, Gauss Seidel, Conjugate Gradient, and Generalized Minimal Residual (GMRES) iterative methods on a sample image for the inference problem. The conjugate gradient method outperforms the other competitors in terms of number of iterations taken to converge to a residual of tolerance 10^{-6} .

4-connected, 8-connected and 12-connected neighbourhoods. While a number of direct linear system solver methods exist, the sheer size of the system matrix A renders them prohibitive, because of large memory requirements. For large problems, a number of iterative methods exist, which require less memory, come with convergence (to a certain tolerance) guarantees under certain conditions, and can be faster than direct methods.

In this work, we considered the *Jacobi*, *Gauss-Seidel*, *Conjugate Gradient*, and *Generalized Minimal Residual* (GMRES) methods [14], as candidates for iterative solvers. Table 1 shows the average number of iterations required by the aforementioned methods for solving the inference problem in equation 2. We used 25 images in this analysis, and a tolerance of 10^{-6} . We initialized the solution \mathbf{x} by the matrix B for each of these methods. Figure 3 shows the convergence of these methods for one of these images. Conjugate gradients clearly stand out as being the fastest of these methods, so our following results use the conjugate gradient method.

We note that inference for the Gaussian CRF can be performed exactly using standard naive mean-field, i.e. we can understand mean-field as solving the linear system of Eq. 2. As we demonstrate in the supplemental material, parallel mean-field amounts to using the Jacobi algorithm while sequential mean-field amounts to using the Gauss-Seidel algorithm, which are the two weakest baselines in our comparisons. When implementing mean-field on the GPU we are left with single weakest baseline, the Jacobi method. Similarly, back-propagating the error signal by unrolled mean-field inference amounts to solving the system in Eq. 4 with the Jacobi/Gauss-Seidel method. So we can understand all of the other algorithms, as being competitive alternatives to ‘vanilla’ mean-field inference and unrolled back-propagation for Eq. 4 respectively.

Method	Iterations
Jacobi	24.8
Gauss Siedel	16.4
GMRES	14.8
Conjugate Gradient	13.2

Table 1: Linear solver statistics

Method	Inference time (s)	Iterations
Zero Initialization	0.0228	15.4
B Initialization	0.0203	13.2
Coarse Inference	0.0041	12.4
Coarse Initialization	0.0158	8.8

Table 2: Multi-grid statistics

Initialization and multi-grid experiments Iterative methods typically set the initial solution \mathbf{x} to a zero-vector in case an initialization is not available. In this work, the matrix B represents the unary terms, and A represents the pairwise interactions, and we expect the pairwise interactions to be small, compared to the unary terms, to ensure the matrix $A + \lambda \mathbf{I}$ remains positive definite. Thus, our matrix is diagonally- dominant, and B serves as a good initial solution. Table 2 shows that the convergence of the inference problem is faster with this initialization compared to zero initialization.

Inspired by the improvement in convergence by initializing the solution with the unary terms, we experimented with the standard V-cycle geometric multigrid approach [14], commonly used in several segmentation applications e.g. [16,17]. Here, we first solve a coarse inference problem by downsampling our unary and pairwise terms by a factor of two. The coarse problem, thus, has $\frac{1}{4}^{th}$ the number of variables as the original problem, and can be solved in $\frac{1}{4}^{th}$ of the time, as can be seen in table 2. We then up-sample the solution to the coarse problem and use it as the initial solution to the original problem. This initialization improves the convergence of the original inference problem. Table 2 shows the number of iterations, and inference time in seconds for (a) inference with zero initialization, (b) inference with B initialization, (c) the coarse inference, and (d) inference with initialization by the coarse solution. The reported statistics are averaged over 25 images, use the conjugate gradient method as the linear system solver, and exploit the *cuSparse* library for sparse operations on the GPU.

We realize that the total acceleration obtained over the simple ‘ B ’ initialization ends up being negligible if we add up the coarse- and fine- level results. This suggests exploring more sophisticated multigrid solves, such as the algebraic multigrid [17], which we leave however for future work.

3.2 Multiresolution graph architecture

We now turn to improving the classification accuracy of our structured prediction results by incorporating computation from multiple scales in a single system. Even though CNNs are designed to be largely scale-invariant, it has been repeatedly reported [18,19] that fusing information from a CNN operating at multiple scales can improve image labelling performance. These results have been obtained for feedforward CNNs - we consider how these could be extended to CNNs with lateral connections, as in our case.

A simple way of achieving this would be to use multiple image resolutions, construct one structured prediction module per resolution, train these as disjoint networks, and average the final results. This amounts to solving three decoupled systems - by itself yields a certain improvement as reported in Sec. 4

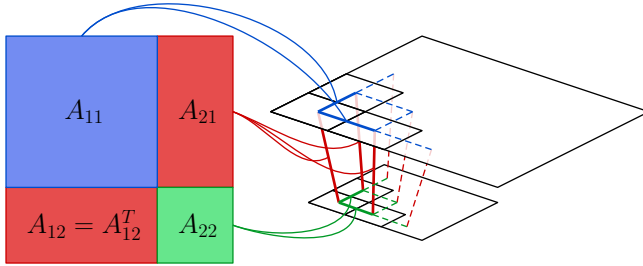


Fig. 4: Schematic diagram of matrix A for the multi-resolution formulation in Sec. 3.2. In this example, we have the input image at 2 resolutions. The pairwise matrix A contains two kinds of pairwise interactions: (a) neighbourhood interactions between pixels at the same resolution (these interactions are shown as the blue and green squares), and (b) interactions between the same image region at two resolutions (these interactions are shown as red rectangles). While interactions of type (a) encourage the pixels in a neighbourhood to take the same or different label, the interactions of type (b) encourage the same image region to take the same labels at different resolutions.

We advocate however a richer connectivity that couples the two systems, allowing information to flow across scales. As illustrated in Fig. 4 the resulting linear system captures the following multi-resolution interactions simultaneously: (a) pairwise constraints between pixels at each resolution, and (b) pairwise constraints between the same image region at two different resolutions. These inter-resolution pairwise terms connect a pixel in the image at one resolution, to the pixel it would spatially correspond to at another resolution. The inter-resolution connections help enforce a different kind of pairwise consistency: rather than encouraging pixels in a neighbourhood to have the same/different label, these encourage image regions to have the same/different labels across resolutions. This is experimentally validated in Sec. 4 to outperform the simpler multi-resolution architecture outlined above.

It is not straightforward to see how this multi-resolution computation scheme could be implemented using existing works on structured prediction, that rely on highly-customized inference algorithms, such as e.g. [6]. However, as soon as this architecture is mapped to a linear system, it becomes straightforward to account for the information flow across and within scales during inference and learning.

3.3 Implementation details

We now discuss the computational efficiency of our quadratic optimization formulation. As mentioned in section 3.1, our approach requires the solution of two linear systems during training and a single system during testing. Our fastest approach (Sec. 3.1) allows the linear system to be solved in ~ 0.02 seconds. Our current implementation is GPU based, and uses the highly optimized *cuBlas* and *cuSparse* libraries for linear algebra on large sparse matrices. The *cuSparse* library requires the matrices to be in the compressed-storage-row (CSR) format in order to fully optimize linear algebra for sparse matrices. Our implementation caches the indices of the CSR matrices, and as

such their computation time is not taken into account in the calculations above, since their computation time is zero for streaming applications, or if the images get warped to a canonical size. In applications where images may be coming at different dimensions, considering that the indexes have been precomputed for the changing dimensions, an additional overhead of ~ 0.1 seconds per image is incurred to read the binary files containing the cached indexes from the hard disk (using an SSD drive could further reduce this). We intend to make our implementation publicly available, we therefore omit further implementation details.

4 Experiments

In this section, we describe our experimental setup, network architecture and results.

Dataset. We evaluate our methods on the *VOC PASCAL 2012 image segmentation benchmark*. This benchmark uses the VOC PASCAL 2012 dataset, which consists of 1464 training and 1449 validation images with manually annotated pixel-level labels for 20 foreground object classes, and 1 background class. In addition, we exploit the additional pixel-level annotations provided by [3], obtaining 10582 training images in total. The test set has 1456 unannotated images. The evaluation criterion is the pixel intersection-over-union (IOU) metric, averaged across the 21 classes. Our baseline network described below was pretrained on the MS-COCO 2014 trainval dataset [20].

Baseline network. Our basenet is based on the Deeplab-LargeFOV network from [5]. As in [19], we extend it to get a multi-resolution network, which operates at three resolutions with tied weights. More precisely, our network downsamples the input image by factors of 2 and 3 and later *fuses* the downsampled activations with the original resolution via concatenation followed by convolution. The layers at three resolutions share weights. This acts like a strong baseline for a purely feedforward network.

QO network. We extend our basenet to accommodate the binary stream of our network. Figure 1 shows a rough schematic diagram of our network. The basenet forms the unary stream of our QO network, while the pairwise terms are drawn from the 3rd pooling layers of the three resolutions via convolutional layers. The hyper-parameter λ (equation 1) was set to 10 for all experiments. A value of 10 ensures that the matrix $(A + \lambda I)$ is positive definite. In our experiments, we observed that using higher values of lambda typically requires more time to converge to the same level of performance.

4.1 Experiments on train+aug - val data

In this set of experiments, we train our methods on the *train+aug* images, and evaluate them on the *val* images for testing various hypotheses.

Image resolutions. Since we have a multi-resolution basenet with two downsampled branches, upsampling the image is necessary [19] to preserve image detail in the lower resolution branches. Here, we experiment with three input image resolutions, 577×385 , 673×481 , and 865×673 . Table 3 shows a consistent improvement in performance as we upsample the input image. Resolutions higher than 865×673 could not be tested due to GPU memory constraints. Kindly note that while a resolution of 865×673 gives better performance, the subsequent experiments in this section are

conducted at a resolution of 673×481 due to limited time resources. However, our preliminary results on the test set (section 4.1) use high resolution (865×673) images.

Multi-loss training. As described, the final activations of our basenet come from the *fusion* of activations of three streams operating at different resolutions. The deeplab network uses one softmax-loss function applied to the final output of the network during training. Here we test the effect of having additional softmax-losses, one per resolution, applied to the activations of the three *resolution* streams. Table 4 shows that while having multi-losses improves the performance of the basenet, it worsens the performance of the 4-connected QO network (QO_4). We attribute this to the fact that having multiple losses on the activations of the *resolution* streams imposes additional constraints on the unary activations of the network. In the absence of these constraints, the quadratic optimizer has access to a larger search space to find the set of parameters that satisfy the global objective. Thus, in the subsequent experiments, we do not use multi-loss for the QO network. However, we always use multi-loss for the basenet, to have a stronger baseline. Table 4 shows that the QO network outperforms the basenet by 1.15%, indicating that including pairwise terms, alongside the unary terms helps improve the performance.

Neighbourhoods. Here we study the effect of having larger neighbourhoods among image regions, thus allowing richer connectivity. More precisely, we study three kinds of connectivities: (a) 4-connected (QO_4), where each pixel is connected to its left, right, top, and bottom neighbours, (b) 8-connected (QO_8), where each pixel is additionally connected to the 4 diagonally adjacent neighbours, and (c) 12-connected (QO_{12}), where each pixel is connected to 2 left, right, top, bottom neighbours besides the diagonally adjacent ones. Table 5 demonstrates that while there are improvements in performance upon increasing connectivities, these are not substantial.

Resolution	IoU
Basenet 577×385	74.89
Basenet 673×481	75.21
Basenet 865×673	75.69

Table 3: Image resolutions

Method	IoU
Basenet	73.44
Basenet multi-loss	75.21
QO_4	76.36
QO_4 multi-loss	76.25

Table 4: Multi-loss

Method	IoU
QO_4	76.36
QO_8	76.40
QO_{12}	76.42

Table 5: Neighbourhoods

Method	$QO_{4 \times 1}$	$QO_{4 \times 2}$	$QO_{4 \times 3}$	$QO_{4 \times 1}^{res}$	$QO_{4 \times 2}^{res}$	$QO_{4 \times 3}^{res}$	$QO_{4 \times 1}^{mres}$	$QO_{4 \times 2}^{mres}$	$QO_{4 \times 3}^{mres}$
IoU	76.36	76.65	76.81	76.69	76.75	76.76	76.93	76.94	76.93

Table 6: Multiple QO: The subscripts $c \times n$ implies c connected neighbourhood, and n units stacked end to end. The superscripts *res* and *mres* indicate one QO per resolution, and a multi-resolution QO respectively.

Given that we obtain diminishing returns, rather than trying even larger neighbourhoods to improve performance, we focus on the effect of having multiple QO units in the network, i.e. solving the inference problem more than once.

Multiple QO. One straightforward extension of the QO network is to have multiple stacked QO units ($QO_{\times n}$, the subscript n denoting the number of stacked QO units), one after the other, each using the output of the previous one to further enforce pairwise consistencies in the output. We hypothesized this would improve performance because it would allow stricter smoothening/sharpening of unary activations - and can be understood as a variant of recurrent CRF [7], where each layer performs exact inference.

As we described in Sec. 3.2, there are two additional ways to extend the network; the simplest is to have one QO unit per resolution (QO^{res}), thereby enforcing pairwise consistencies individually at each resolution before fusing them, while the more sophisticated one is to have information flow both within and across scales, amount to a joint multi-scale CRF inference task, illustrated in Fig. 4.

When benchmarking these three alternatives, it can be seen that the single 4-connected QO (QO_4) benefits most from stacking, as the performance improves by 0.45% simply by stacking it 3 times. Further, the performance gains due to stacking are not as pronounced in the methods where we have QO units operating at multiple-resolutions (0.07% for QO^{res} , and 0.01% for QO^{mres}). We believe this phenomenon is expected, as the the QO modules operating at different resolutions already do a better job of enforcing pairwise consistencies, therefore not much can be gained by repeatedly solving the inference problem. We also observe that a single QO^{mres} outperforms a single QO by 0.57%, thereby significantly outperforming the baseline by 1.72%.

In conclusion, the introduced multi-resolution architecture effectively alleviates the need for stacking QO's to achieve multi-scale processing.

4.2 Experiments on train+aug+val - test data

In this set of experiments, we train our methods on the *train+aug+val* images, and evaluate them on the *test* images. We use high resolution (865×673) images, and also use the dense CRF post processing as in [5,21]. We note that these results have been obtained from a preliminary network, while the results of the converged network will be in the supplemental material.

Our results are tabulated in tables 7 and 8. We first compare two of our methods QO and QO^{mres} with the basenet, where the relative improvements can be most clearly demonstrated.

Table 7 shows our method brings about a 0.85% improvement in performance over the baseline network. We achieve a further boost in performance upon using the dense CRF post processing strategy, consistently for all methods. We observe that our method yields an improvement, that is entirely complementary to the improvement obtained by combining with Dense-CRF.

Table 8 compares our method with the previously published state of the art methods. We note that we cannot directly compete with [8] and [9], since we do not use their additional (and orthogonal) extensions for higher-resolution networks and global context. When benchmarking against directly comparable techniques, we observe that even though we do not use end-to-end training for the CRF module stacked on top of our QO network, our method outperforms the previous state of the art CRF-RNN system of [7] by a margin of 0.6%. We anticipate further improvements by integrating end-to-end CRF training with our QO.

Method	IoU	IoU after <i>dense CRF</i>
Basenet	72.27	73.47
QO ₄	73.32	74.95
QO ₄ ^{res}	73.75	75.26

Table 7: Performance of our methods on the VOC PASCAL 2012 Image Segmentation Benchmark. Our baseline network (Basenet) is a variant of Deeplab-LargeFOV [5] network. In this table, we demonstrate systematic improvements in performance upon the introduction of our Quadratic Optimization (QO), and multi-resolution (QO^{res}) approaches. Use of denseCRF to post process our results gives us a slight boost in performance.

5 Conclusions and Future Work

In this work we propose a quadratic optimization method for deep networks which can be used for predicting continuous vector-valued variables. The inference is efficient and exact and can be solved in 0.02 seconds on the GPU for each image using the conjugate gradient method. We propose a deep-learning framework which learns features and model parameters simultaneously in an end-to-end FCN training algorithm. Our implementation is fully GPU based, and implemented using the *Caffe* library. Our experimental results indicate that using pairwise terms boosts performance of the network on the task of image segmentation, and our results are competitive with the state of the art methods on the Pascal image segmentation benchmark, while being preliminary and substantially simpler.

Method	mean IoU (%)
DeepLab-CRF [5]	66.4
DeepLab-MSc-CRF [5]	67.1
DeepLab-CRF-7x7 [5]	70.3
DeepLab-CRF-LargeFOV [5]	70.3
DeepLab-MSc-CRF-LargeFOV [5]	71.6
Deeplab-Cross-Joint [21]	73.9
CRFRNN [7]	74.7
Adelaide Context [8]	77.8
Deep Parsing Network [9]	77.4
Ours (QO ₄ ^{res})	75.3

Table 8: Comparison of our method with the previously published state of the art approaches on the VOC PASCAL 2012 image segmentation benchmark.

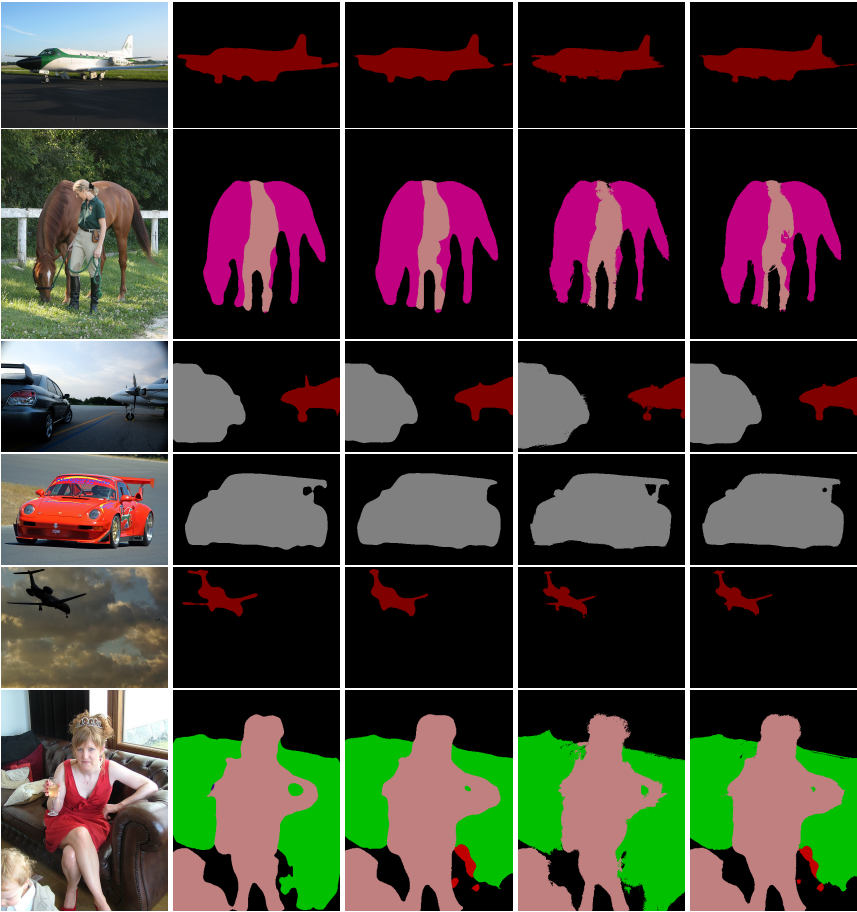


Fig. 5: Visual results on the VOC PASCAL 2012 test set. The first column shows the colour image, the second column shows the QO^{mres} predicted segmentation, the third column shows the basenet predicted segmentation. The fourth column shows the QO^{mres} output after dense CRF post processing, and the final column shows the basenet output after dense CRF. It can be seen that our multi-resolution network captures the finer details better than the basenet: the tail of the airplane in the first image, the horse rider's body in the second image, the aircraft fan in the third image, the road between the car's tail in the fourth image, and the wings of the aircraft in the final image, all indicate this. While dense CRF post-processing quantitatively improves performance, it tends to miss very fine details.

References

1. Farabet, C., Couprie, C., Najman, L., LeCun, Y.: Learning hierarchical features for scene labeling. *IEEE Trans. PAMI* (2013) **1**
2. Mostajabi, M., Yadollahpour, P., Shakhnarovich, G.: Feedforward semantic segmentation with zoom-out features. In: *CVPR*. (2015) **1**
3. Hariharan, B., Arbeláez, P., Girshick, R., Malik, J.: Hypercolumns for object segmentation and fine-grained localization. In: *CVPR*. (2015) **1, 10**
4. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *CVPR*. (2015) 3431–3440 **1**
5. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062* (2014) **1, 2, 3, 10, 12, 13**
6. Krähenbühl, P., Koltun, V.: Efficient inference in fully connected crfs with gaussian edge potentials. In: *NIPS*. (2011) **1, 9**
7. Zheng, S., Jayasumana, S., Romera-Paredes, B., Vineet, V., Su, Z., Du, D., Huang, C., Torr, P.: Conditional random fields as recurrent neural networks. In: *ICCV*. (2015) **1, 2, 12, 13**
8. Lin, G., Shen, C., Reid, I.D., van den Hengel, A.: Efficient piecewise training of deep structured models for semantic segmentation. *CVPR* (2016) **2, 3, 12, 13**
9. Liu, Z., Li, X., Luo, P., Loy, C.C., Tang, X.: Semantic image segmentation via deep parsing network. *CoRR abs/1509.02634* (2015) **2, 3, 12, 13**
10. Tappen, M.F., Liu, C., Adelson, E.H., Freeman, W.T.: Learning gaussian conditional random fields for low-level vision. In: 2007 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2007), 18–23 June 2007, Minneapolis, Minnesota, USA. (2007) **2**
11. Jancsary, J., Nowozin, S., Sharp, T., Rother, C.: Regression tree fields - an efficient, non-parametric approach to image labeling problems. In: *CVPR*. (2012) **2, 4**
12. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093* (2014) **3, 6**
13. Shewchuk, J.R.: An introduction to the conjugate gradient method without the agonizing pain. <https://www.cs.cmu.edu/~quake-papers/painless-conjugate-gradient.pdf> **5**
14. Press, W.H., Teukolsky, S.A., Vetterling, W.T., Flannery, B.P.: *Numerical Recipes in C*, 2nd Edition. Cambridge University Press (1992) **6, 7, 8**
15. Golub, G.H., Loan, C.F.V.: *Matrix computations* (3. ed.). Johns Hopkins University Press (1996) **6**
16. Maire, M., Yu, S.X.: Progressive multigrid eigensolvers for multiscale spectral segmentation. In: *International Conference on Computer Vision (ICCV)*. (2013) **8**
17. Papandreou, G., Maragos, P.: Multigrid geometric active contour models. *IEEE Transactions on Image Processing* **16**(1) (Jan 2007) 229–240 **8**
18. Chen, L., Yang, Y., Wang, J., Xu, W., Yuille, A.L.: Attention to scale: Scale-aware semantic image segmentation. *CoRR abs/1511.03339* (2015) **8**
19. Kokkinos, I.: Pushing the Boundaries of Boundary Detection using Deep Learning. In: *ICLR*. (2016) **8, 10**
20. *et al.*, L.: Microsoft coco: Common objects in context. In: *ECCV*. (2014) **10**
21. Chen, L.C., Papandreou, G., Murphy, K., Yuille, A.L.: Weakly- and semi-supervised learning of a deep convolutional network for semantic image segmentation. *ICCV* (2015) **12, 13**