# Learning Sparse High Dimensional Filters:
# Image Filtering, Dense CRFs and Bilateral Neural Networks

Varun Jampani[1], Martin Kiefel[1,2] and Peter V. Gehler[1,2]
[1]Max Planck Institute for Intelligent Systems, Tübingen, 72076, Germany
[2]Bernstein Center for Computational Neuroscience, Tübingen, 72076, Germany
{varun.jampani, martin.kiefel, peter.gehler}@tuebingen.mpg.de

## Abstract

*Bilateral filters have wide spread use due to their edge-preserving properties. The common use case is to manually choose a parametric filter type, usually a Gaussian filter. In this paper, we will generalize the parametrization and in particular derive a gradient descent algorithm so the filter parameters can be learned from data. This derivation allows to learn high dimensional linear filters that operate in sparsely populated feature spaces. We build on the permutohedral lattice construction for efficient filtering. The ability to learn more general forms of high-dimensional filters can be used in several diverse applications. First, we demonstrate the use in applications where single filter applications are desired for runtime reasons. Further, we show how this algorithm can be used to learn the pairwise potentials in densely connected conditional random fields and apply these to different image segmentation tasks. Finally, we introduce layers of bilateral filters in CNNs and propose* bilateral neural networks *for the use of high-dimensional sparse data. This view provides new ways to encode model structure into network architectures. A diverse set of experiments empirically validates the usage of general forms of filters.*

## 1. Introduction

Image convolutions are basic operations for many image processing and computer vision applications. In this paper we will study the class of bilateral filter convolutions and propose a general image adaptive convolution that can be learned from data. The bilateral filter [5, 45, 50] was originally introduced for the task of image denoising as an edge preserving filter. Since the bilateral filter contains the spatial convolution as a special case, we will in the following directly state the general case. Given an image $\mathbf{v} = (\mathbf{v}_1, \ldots, \mathbf{v}_n), \mathbf{v}_i \in \mathbb{R}^c$ with $n$ pixels and $c$ channels, and for every pixel $i$, a $d$ dimensional feature vector $\mathbf{f}_i \in \mathbb{R}^d$ (*e.g.*, the $(x, y)$ position in the image $\mathbf{f}_i = (x_i, y_i)^\top$). The

bilateral filter then computes

$$\mathbf{v}'_i = \sum_{j=1}^{n} \mathbf{w}_{\mathbf{f}_i, \mathbf{f}_j} \mathbf{v}_j. \tag{1}$$

for all $i$. Almost the entire literature refers to the bilateral filter as a synonym of the Gaussian parametric form $\mathbf{w}_{\mathbf{f}_i, \mathbf{f}_j} = \exp\left(-\frac{1}{2}(\mathbf{f}_i - \mathbf{f}_j)^\top \Sigma^{-1}(\mathbf{f}_i - \mathbf{f}_j)\right)$. The features $\mathbf{f}_i$ are most commonly chosen to be position $(x_i, y_i)$ and color $(r_i, g_i, b_i)$ or pixel intensity. To appreciate the edge-preserving effect of the bilateral filter, consider the five-dimensional feature $\mathbf{f} = (x, y, r, g, b)^\top$. Two pixels $i, j$ have a strong influence $\mathbf{w}_{\mathbf{f}_i, \mathbf{f}_j}$ on each other only if they are close in position *and* color. At edges the color changes, therefore pixels lying on opposite sides have low influence and thus this filter does not blur across edges. This behaviour is sometimes referred to as "image adaptive", since the filter has a different shape when evaluated at different locations in the image. More precisely, it is the projection of the filter to the two-dimensional image plane that changes, the filter values $\mathbf{w}_{\mathbf{f}, \mathbf{f}'}$ do not change. The filter itself can be of $c$ dimensions $\mathbf{w}_{\mathbf{f}_i, \mathbf{f}_j} \in \mathbb{R}^c$, in which case the multiplication in Eq. (1) becomes an inner product. For the Gaussian case the filter can be applied independently per channel. For an excellent review of image filtering we refer to [39].

The filter operation of Eq. (1) is a sparse high-dimensional convolution, a view first developed in [40]. An image $v$ is not sparse in the spatial domain, we observe pixels values for all locations $(x, y)$. However, when pixels are understood in a higher dimensional feature space, *e.g.*, $(x, y, r, g, b)$, the image becomes a sparse signal, since the $r, g, b$ values lie scattered in this five-dimensional space. This view on filtering is the key difference of the bilateral filter compared to the common spatial convolution. An image edge is not "visible" for a filter in the spatial domain alone, whereas in the 5D space it is. The edge-preserving behaviour is possible due to the higher dimensional operation. Other data can naturally be understood as sparse signals, *e.g.*, 3D surface points.

The contribution of this paper is to propose a general and learnable sparse high dimensional convolution. Our technique builds on efficient algorithms that have been developed to approximate the Gaussian bilateral filter and re-uses them for more general high-dimensional filter operations. Due to its practical importance (see related work in Sec. 2) several efficient algorithms for computing Eq. (1) have been developed, including the bilateral grid [40], Gaussian KD-trees [3], and the permutohedral lattice [2]. The design goal for these algorithms was to provide a) fast runtimes and b) small approximation errors for the Gaussian filter case. The key insight of this paper is to use the permutohedral lattice and use it not as an approximation of a predefined kernel but to freely parametrize its values. We relax the separable Gaussian filter case from [2] and show how to compute gradients of the convolution (Sec. 3) in lattice space. This enables learning the filter from data.

This insight has several useful consequences. We discuss applications where the bilateral filter has been used before: image filtering (Sec. 4) and CRF inference (Sec. 5). Further we will demonstrate how the free parametrization of the filters enables us to use them in Deep convolutional neural networks (CNN) and allow convolutions that go beyond the regular spatially connected receptive fields (Sec. 6). For all domains, we present various empirical evaluations with a wide range of applications.

## 2. Related Work

We categorize the related work according to the three different generalizations of this work.

**Image Adaptive Filtering:** The literature in this area is rich and we can only provide a brief overview. Important classes of image adaptive filters include the bilateral filters [5, 50, 45], non-local means [13, 6], locally adaptive regressive kernels [47], guided image filters [26] and propagation filters [42]. The kernel least-squares regression problem can serve as a unified view of many of them [39]. In contrast to the present work that learns the filter kernel using supervised learning, all these filtering schemes use a predefined kernel. Because of the importance of the bilateral filtering to many applications in image processing, much effort has been devoted to derive fast algorithms; most notably [40, 2, 3, 22]. Surprisingly, the only attempt to learn the bilateral filter we found is [27] that casts the learning problem in the spatial domain by rearranging pixels. However, the learned filter does not necessarily obey the full region of influence of a pixel as in the case of a bilateral filter. The bilateral filter also has been proposed to regularize a large set of applications in [9, 8] and the respective optimization problems are parametrized in a bilateral space. They provide gradients to use the filters as part of a learning system but unlike this work restrict the filters to be Gaussian.

**Dense CRF:** The key observation of [33] is that mean-field inference update steps in densely connected CRFs with Gaussian edge potentials require Gaussian bilateral filtering operations. This enables tractable inference through the application of a fast filter implementation from [2]. This quickly found wide-spread use, *e.g.*, the combination of CNNs with a dense CRF is among the best performing segmentation models [16, 55, 11]. These works combine structured prediction frameworks on top of CNNs, to model the relationship between the desired output variables thereby significantly improving upon the CNN result. Bilateral neural networks, that are presented in this work, provide a principled framework for encoding the output relationship, using the feature transformation inside the network itself thereby alleviating some of the need for later processing. Several works [34, 18, 31] demonstrate how to learn free parameters of the dense CRF model. However, the parametric form of the pairwise term always remains a Gaussian. Campbell *et al.* [15] embed complex pixel dependencies into an Euclidean space and use a Gaussian filter for pairwise connections. This embedding is a pre-processing step and can not directly be learned. In Sec. 5 we will discuss how to learn the pairwise potentials, while retaining the efficient inference strategy of [33].

**Neural Networks:** In recent years, tremendous progress in a wide range of computer vision applications has been observed due to the use of CNNs. Most CNN architectures use spatial convolution layers, which have fixed local receptive fields. This work suggests to replace these layers with bilateral filters, which have a varying spatial receptive field depending on the image content. The equivalent representation of the filter in a higher dimensional space leads to sparse samples that are handled by a permutohedral lattice data structure. Similarly, Bruna *et al.* [12] propose convolutions on irregularly sampled data. Their graph construction is closely related to the high-dimensional convolution that we propose and defines weights on local neighborhoods of nodes. However, the structure of the graph is bound to be fixed and it is not straightforward to add new samples. Furthermore, re-using the same filter among neighborhoods is only possible with their costly spectral construction. Both cases are handled naturally by our sparse convolution. Jaderberg *et al.* [29] propose a spatial transformation of signals within the neural network to learn invariances for a given task. The work of [28] propose matrix backpropagation techniques which can be used to build specialized structrual layers such as normalized-cuts. Graham *et al.* [24] propose extensions from 2D CNNs to 3D sparse signals. Our work enables sparse 3D filtering as a special case, since we use an algorithm that allows for even higher dimensional data.
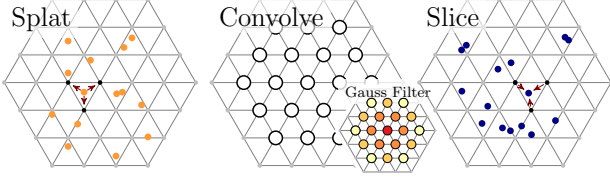
Figure 1. **Schematic of the permutohedral convolution.** Left: splatting the input points (orange) onto the lattice corners (black); Middle: The extent of a filter on the lattice with a $s = 2$ neighborhood (white circles), for reference we show a Gaussian filter, with its values color coded. The general case has a free scalar/vector parameter per circle. Right: The result of the convolution at the lattice corners (black) is projected back to the output points (blue). Note that in general the output and input points may be different.
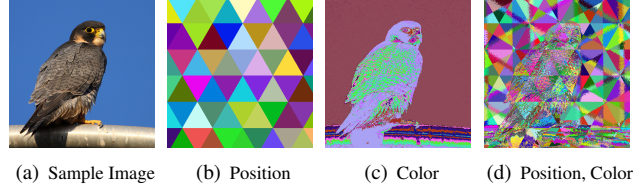


Figure 2. **Visualization of the Permutohedral Lattice.** (a) Input image; Lattice visualizations for different feature spaces: (b) 2D position features: $0.01(x, y)$, (c) color features: $0.01(r, g, b)$ and (d) position and color features: $0.01(x, y, r, g, b)$. All pixels falling in the same simplex cell are shown with the same color.

## 3. Learning Sparse High Dimensional Filters

In this section, we describe the main technical contribution of this work, we generalize the permutohedral convolution [2] and show how the filter can be learned from data.

Recall the form of the bilateral convolution from Eq. (1). A naive implementation would compute for every pixel $i$ all associated filter values $\mathbf{w}_{\mathbf{f}_i, \mathbf{f}_j}$ and perform the summation independently. The view of $\mathbf{w}$ as a linear filter in a higher dimensional space, as proposed by [40], opened the way for new algorithms. Here, we will build on the permutohedral lattice convolution developed in Adams *et al.* [2] for approximate Gaussian filtering.

### 3.1. Permutohedral Lattice Convolutions

We first review the permutohedral lattice convolution for Gaussian bilateral filters from Adams *et al.* [2] and describe its most general case.

As before, we assume that every image pixel $i$ is associated with a $d$-dimensional feature vector $\mathbf{f}_i$. Gaussian bilateral filtering using a permutohedral lattice approximation involves 3 steps. We begin with an overview of the algorithm, then discuss each step in more detail in the next paragraphs. Fig. 1 schematically shows the three operations for 2D features. First, interpolate the image signal on the $d$-dimensional grid plane of the permutohedral lattice, which is called *splatting*. A permutohedral lattice is the tessellation of space into permutohedral simplices. We refer to [2] for details of the lattice construction and its properties. In Fig. 2, we visualize the permutohedral lattice in the image plane, where every simplex cell receives a different color. All pixels of the same lattice cell have the same color. Second, *convolve* the signal on the lattice. And third, retrieve the result by interpolating the signal at the original $d$-dimensional feature locations, called *slicing*. For example, if the features used are a combination of position and color $\mathbf{f}_i = (x_i, y_i, r_i, g_i, b_i)^\top$, the input signal is mapped into the 5D cross product space of position and color and then convolved with a 5D tensor. The result is then mapped back to the original space. In practice we use a feature scal-

ing $D\mathbf{f}$ with a diagonal matrix $D$ and use separate scales for position and color features. The scale determines the distance of points and thus the size of the lattice cells. More formally, the computation is written by $\mathbf{v}' = S_{\text{slice}} B S_{\text{splat}} \mathbf{v}$ and all involved matrices are defined below. For notational convenience we will assume scalar input signals $v_i$, the vector valued case is analogous, the lattice convolution changes from scalar multiplications to inner products.

**Splat:** The splat operation (cf. left-most image in Fig. 1) finds the enclosing simplex in $\mathcal{O}(d^2)$ on the lattice of a given pixel feature $\mathbf{f}_i$ and distributes its value $v_i$ onto the corners of the simplex. How strong a pixel contributes to a corner $j$ is defined by its barycentric coordinate $b_{i,j} \in \mathbb{R}$ inside the simplex. Thus, the value $l_j \in \mathbb{R}$ at a lattice point $j$ is computed by summing over all enclosed input points; more precisely, we define an index set $J_i$ for a pixel $i$, which contains all the lattice points $j$ of the enclosing simplex

$$\ell = S_{\text{splat}} \mathbf{v}; (S_{\text{splat}})_{j,i} = b_{i,j}, \text{ if } j \in J_i, \text{ otherwise } 0. \quad (2)$$

**Convolve:** The permutohedral convolution is defined on the lattice neighborhood $N_s(j)$ of lattice point $j$, *e.g.*, only $s$ grid hops away. More formally

$$\ell' = B\ell; (B)_{j',j} = w_{j,j'}, \text{ if } j' \in N_s(j), \text{ otherwise } 0. \quad (3)$$

An illustration of a two-dimensional permutohedral filter is shown in Fig. 1 (middle). Note that we already presented the convolution in the general form that we will make use of. The work of [2] chooses the filter weights such that the resulting operation approximates a Gaussian blur, which is illustrated in Fig. 1. Further, the algorithm of [2] takes advantage of the separability of the Gaussian kernel. Since we are interested in the most general case, we extended the convolution to include non-separable filters $B$.

**Slice:** The slice operation (cf. right-most image in Fig. 1) computes an output value $v'_{i'}$ for an output pixel $i'$ again based on its barycentric coordinates $b_{i,j}$ and sums over the corner points $j$ of its lattice simplex

$$\mathbf{v}' = S_{\text{slice}} \ell'; (S_{\text{slice}})_{i,j} = b_{i,j}, \text{ if } j \in J_i, \text{ otherwise } 0 \quad (4)$$

The splat and slice operations take a role of an interpolation between the different signal representations: the irregular and sparse distribution of pixels with their associated feature vectors and the regular structure of the permutohedral lattice points. Since high-dimensional spaces are

usually sparse, performing the convolution densely on all lattice points is inefficient. So, for speed reasons, we keep track of the populated lattice points using a hash table and only convolve at those locations.

## 3.2. Learning Permutohedral Filters

The *fixed* set of filter weights $\mathbf{w}$ from [2] in Eq. 3 are designed to approximate a Gaussian filter. However, the convolution kernel $\mathbf{w}$ can naturally be understood as a general filtering operation in the permutohedral lattice space with free parameters. In the exposition above we already presented this general case. As we will show in more detail later, this modification has non-trivial consequences for bilateral filters, CNNs and probabilistic graphical models.

The size of the neighborhood $N_s(k)$ for the blur in Eq. 3 compares to the filter size of a spatial convolution. The filtering kernel of a common spatial convolution that considers $s$ points to either side in all dimensions has $(2s + 1)^d \in \mathcal{O}(s^d)$ parameters. A comparable filter on the permutohedral lattice with an $s$ neighborhood is specified by $(s + 1)^{d+1} - s^{d+1} \in \mathcal{O}(s^d)$ elements (cf. Sec. A). Thus, both share the same asymptotic size.

By computing the gradients of the filter elements we enable the use of gradient based optimizers, *e.g.*, backpropagation for CNN in the same way that spatial filters in a CNN are learned. The gradients with respect to $\mathbf{v}$ and the filter weights in $B$ of a scalar loss $L$ are:

$$\frac{\partial L}{\partial \mathbf{v}} = S'_{\text{splat}} B' S'_{\text{slice}} \frac{\partial L}{\partial \mathbf{v}'}, \tag{5}$$

$$\frac{\partial L}{\partial (B)_{i,j}} = \left( S'_{\text{slice}} \frac{\partial L}{\partial \mathbf{v}} \right)_i (S_{\text{splat}} \mathbf{v})_j. \tag{6}$$

Both gradients are needed during backpropagation and in experiments, we use stochastic backpropagation for learning the filter kernel. The permutohedral lattice convolution is parallelizable, and scales linearly with the filter size. Specialized implementations run at interactive speeds in image processing applications [2]. Our implementation allows arbitrary filter parameters and the computation of the gradients on both CPU and GPU. It is also integrated into the caffe deep learning framework [30].

## 4. Single Bilateral Filter Applications

In this section we will consider the problem of joint bilateral upsampling [32] as a prominent instance of a single bilateral filter application. See [41] for a recent overview of other bilateral filter applications. Further experiments on image denoising and 3D body mesh denoising are included in Sec. B.3, together with details about exact experimental protocols and more visualizations.

| Upsampling factor | Bicubic | Gaussian | Learned |
|---|---|---|---|
| Color Upsampling (PSNR) | | | |
| 2x | 24.19 / 30.59 | 33.46 / 37.93 | **34.05 / 38.74** |
| 4x | 20.34 / 25.28 | 31.87 / 35.66 | **32.28 / 36.38** |
| 8x | 17.99 / 22.12 | 30.51 / 33.92 | **30.81 / 34.41** |
| 16x | 16.10 / 19.80 | 29.19 / 32.24 | **29.52 / 32.75** |
| Depth Upsampling (RMSE) | | | |
| 8x | 0.753 | 0.753 | **0.748** |

Table 1. **Joint bilateral upsampling.** (top) PSNR values corresponding to various upsampling factors and upsampling strategies on the test images of the Pascal VOC12 segmentation / high-resolution 2MP dataset; (bottom) RMSE error values corresponding to upsampling depth images estimated using [19] computed on the test images from the NYU depth dataset [43].

## 4.1. Joint Bilateral Upsampling

A typical technique to speed up computer vision algorithms is to compute results on a lower scale and upsample the result to the full resolution. This upsampling step may use the original resolution image as a guidance image. A joint bilateral upsampling approach for this problem setting was developed in [32]. We describe the procedure for the example of upsampling a color image. Given a high resolution gray scale image (the guidance image) and the same image on a lower resolution but in colors, the task is to upsample the color image to the same resolution as the guidance image. Using the permutohedral lattice, joint bilateral upsampling proceeds by splatting the color image into the lattice, using 2D position and 1D intensity as features and the 3D RGB values as the signal. A convolution is applied in the lattice and the result is read out at the pixels of the high resolution image, that is using the 2D position and intensity of the guidance image. The possibility of reading out (slicing) points that are not necessarily the input points is an appealing feature of the permutohedral lattice convolution.

### 4.1.1 Color Upsampling

For the task of color upsampling, we compare the Gaussian bilateral filter [32] against a learned generalized filter. We experimented with two different datasets: One is Pascal VOC2012 segmentation dataset [20] using the pre-defined train, validation and test splits, and other is 200 high-resolution (2MP) image dataset collected using Google image search [1] with 100 train, 50 validation and 50 test images. For training we use the mean squared error (MSE) criterion and perform stochastic gradient descent with a momentum term of $0.9$, and weight decay of $0.0005$. The learning rate and feature scales $D$ have been cross-validated using the validation split. The Gaussian and the learned filter have the same support. In Table 1 we report result in terms of PSNR for the upsampling factors $2\times, 4\times, 8\times$ and $16\times$. We compare a standard bicubic interpolation, that does not use a guidance image, the Gaussian bilateral filter case (with feature scales optimized on the validation set),

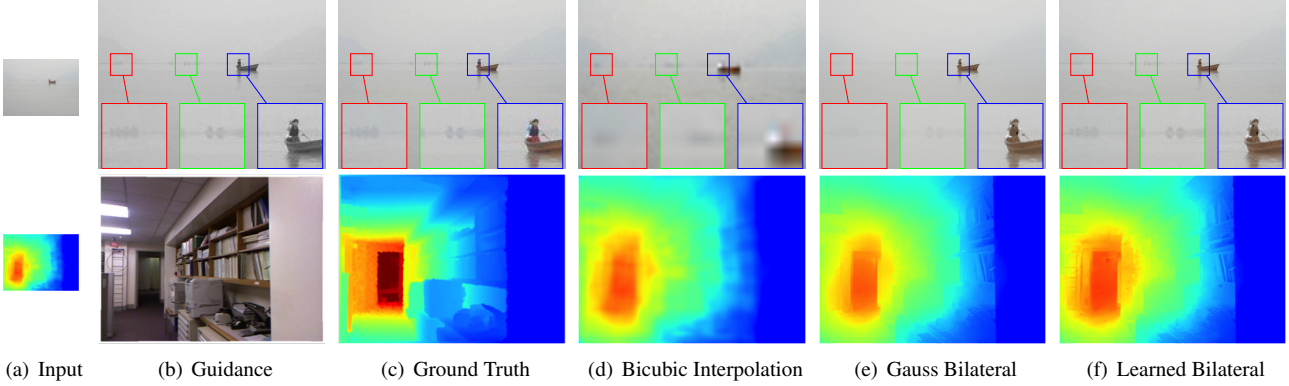| (a) Input | (b) Guidance | (c) Ground Truth | (d) Bicubic Interpolation | (e) Gauss Bilateral | (f) Learned Bilateral |

Figure 3. **Guided Upsampling.** Color (top) and depth (bottom) $8\times$ upsampling results using different methods (best viewed on screen).

and the learned filter. For all upsampling factors, we observe that joint bilateral Gaussian upsampling outperforms bicubic interpolation. Learning a general convolution further improves over the Gaussian filter case. A result of the upsampling is shown in Fig. 3 and more results are included in Sec. C.2. The learned filter recovers finer details in the images

### 4.1.2 Depth Upsampling

We use the work on recovering depth estimates from single RGB images [19] to validate the learned filter on another domain for the task of joint upsampling. The dataset of [43] provides a benchmark for this task and comes with pre-defined train, validation and test splits that we re-use. The approach of [19] is a CNN model that produces a result at 1/4th of the input resolution due to down-sampling operations in max-pooling layers. Furthermore, the authors downsample the $640 \times 480$ images to $320 \times 240$ as a pre-processing step before CNN convolutions. The final depth result is bicubic interpolated to the original resolution. It is this interpolation that we replace with a Gaussian and learned joint bilateral upsampling. The features are five-dimensional position and color information from the high resolution input image. The filter is learned using the same protocol as for color upsampling minimizing MSE prediction error. The quantitative results are shown in Table 1, the Gaussian filter performs equal to the bicubic interpolation ($p$-value 0.311), the learned filter is better ($p$-value 0.015). Qualitative results are shown in Fig 3, both joint bilateral upsampling respect image edges in the result. For this [21] and other tasks specialized interpolation algorithms exist, *e.g.*, deconvolution networks [54]. Part of future work is to equip these approaches with bilateral filters.

## 5. Learning Pairwise Potentials in Dense CRFs

The bilateral convolution from Sec. 3 generalizes the class of dense CRF models for which the mean-field inference from [33] applies. The dense CRF models have

found wide-spread use in various computer vision applications [46, 10, 55, 52, 51, 11]. Let us briefly review the dense CRF and then discuss its generalization.

### 5.1. Dense CRF Review

Consider a fully pairwise connected CRF with discrete variables over each pixel in the image. For an image with $n$ pixels, we have random variables $\mathbf{X} = \{X_1, X_2, \ldots, X_n\}$ with discrete values $\{x_1, \ldots, x_n\}$ in the label space $x_i \in \{1, \ldots, \mathcal{L}\}$. The dense CRF model has unary potentials $\psi_u(x) \in \mathbb{R}^{\mathcal{L}}$, *e.g.*, these can be the output of CNNs. The pairwise potentials in [33] are of the form $\psi_p(x_i, x_j) = \mu(x_i, x_j)k(\mathbf{f}_i, \mathbf{f}_j)$ where $\mu$ is a label compatibility matrix, $k$ is a Gaussian kernel $k(\mathbf{f}_i, \mathbf{f}_j) = \exp(-(\mathbf{f}_i - \mathbf{f}_j)^\top \Sigma^{-1}(\mathbf{f}_i - \mathbf{f}_j))$ and the vectors $\mathbf{f}_i$ are feature vectors, just alike the ones used for the bilateral filtering, *e.g.*, $(x, y, r, g, b)$. The Gibbs distribution for an image $v$ thus reads $p(x|v) \propto \exp(-\sum_i \psi_u(x_i) - \sum_{i>j} \psi_p(x_i, x_j))$. Because of dense connectivity, exact MAP or marginal inference is intractable. The main result of [33] is to derive the mean-field approximation of this model and to relate it to bilateral filtering which enables tractable approximate inference. Mean-field approximates the model $p$ with a fully factorized distribution $q = \prod_i q_i(x_i)$ and solves for $q$ by minimizing their KL divergence $KL(q||p)$. This results in a fixed point equation which can be solved iteratively $t = 0, 1, \ldots$ to update the marginal distributions $q_i$,

$$q_i^{t+1}(x_i) = \frac{1}{Z_i}\exp\{-\psi_u(x_i) - \underbrace{\sum_{j \neq i}\sum_{x_j \in \mathcal{L}}\psi_p(x_i, x_j)q_j^t(x_j)}_{\text{bilateral filtering}}\}.$$

(7)

for all $i$. Here, $Z_i$ ensures normalizations of the distribution $q_i$.

### 5.2. Learning Pairwise Potentials

The proposed bilateral convolution generalizes the class of potential functions $\psi_p$, since they allow a richer class of kernels $k(\mathbf{f}_i, \mathbf{f}_j)$ that furthermore can be learned from data.

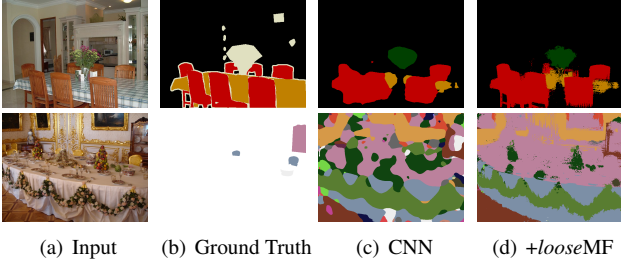| | (a) Input | (b) Ground Truth | (c) CNN | (d) +*loose*MF |

Figure 4. **Segmentation results.** An example result for semantic (top) and material (bottom) segmentation. (c) depicts the unary results before application of MF, (d) after two steps of *loose*-MF with a learned CRF. More examples with comparisons to Gaussian pairwise potentials are in Sec. C.5 and Sec. C.6.

So far, all dense CRF models have used Gaussian potential functions $k$, we replace it with the general bilateral convolution and learn the parameters of kernel $k$, thus in effect learn the pairwise potentials of the dense CRF. This retains the desirable properties of this model class – efficient inference through mean-field and the feature dependency of the pairwise potential. The benefit of the Gaussian pairwise potentials is to encode similarity between pixels through distance and appearance, in other words, pixels with similar $(r, g, b)$ values are likely to be of the same class. In order to learn the form of the pairwise potentials $k$ we make use of the gradients for filter parameters in $k$ and use back-propagation through the mean-field iterations [18, 37] to learn them.

The work of [34] derived gradients to learn the feature scaling $D$ but not the form of the kernel $k$, which still was Gaussian. In [15], the features $f_i$ were derived using a non-parametric embedding into a Euclidean space and again a Gaussian kernel was used. The computation of the embedding was a pre-processing step, not integrated in a end-to-end learning framework. Both aforementioned works are generalizations that are orthogonal to our development and can be used in conjunction.

## 5.3. Experimental Evaluation

We evaluate the effect of learning more general forms of potential functions on two pixel labeling tasks, semantic segmentation of VOC data [20] and material classification [11]. We use pre-trained models from the literature and empirically evaluate the relative change when learning the pairwise potentials. For both the experiments, we use multinomial logistic classification loss and learn the filters via back-propagation.

### 5.3.1 Semantic Segmentation

Semantic segmentation is the task of assigning semantic label (*e.g.*, car, person etc.) to every pixel. We choose the DeepLab network [16], a variant of the VGGnet [44] for obtaining unaries. The DeepLab architecture runs a CNN

| | + MF-1step | + MF-2 step | + *loose* MF-2 step |
|---|---|---|---|
| Semantic segmentation (IoU) - CNN [16]: 72.08 / 66.95 | | | |
| Gauss CRF | +2.48 | +3.38 | +3.38 / +3.00 |
| Learned CRF | +2.93 | +3.71 | **+3.85 / +3.37** |
| Material segmentation (Pixel Accuracy) - CNN [11]: 67.21 / 69.23 | | | |
| Gauss CRF | +7.91 / +6.28 | +9.68 / +7.35 | +9.68 / +7.35 |
| Learned CRF | +9.48 / +6.23 | +11.89 / +6.93 | **+11.91 / +6.93** |

Table 2. **Improved mean-field inference with learned potentials.** (top) Average IoU score on Pascal VOC12 validation/test data [20] for semantic segmentation; (bottom) Accuracy for all pixels / averaged over classes on the MINC test data [11] for material segmentation.

model on the input image to obtain a result that is down-sampled by a factor of 8. The result is then bilinear interpolated to the desired resolution and serves as unaries $\psi_u(x_i)$ in a dense CRF. We use the same Pott's label compatibility function $\mu$, and also use two kernels $k^1(\mathbf{f}_i, \mathbf{f}_j) + k^2(p_i, p_j)$ with the same features $\mathbf{f}_i = (x_i, y_i, r_i, g_i)^\top$ and $\mathbf{p}_i = (x_i, y_i)^\top$ as in [16]. Thus, the two filters operate in parallel on color & position, and spatial domain respectively. We also initialize the mean-field update equations with the CNN unaries. The only change in the model is the type of the pairwise potential function from Gauss to a generalized form.

We evaluate the result after 1 step and 2 steps of mean-field inference and compare the Gaussian filter versus the learned version (cf. Tab. 2). First, as in [16] we observe that one step of mean field improves the performance by 2.48% in Intersection over Union (IoU) score. However, a learned potential increases the score by 2.93%. The same behaviour is observed for 2 steps: the learned result again adds on top of the raised Gaussian mean field performance. Further, we tested a variant of the mean-field model that learns a separate kernel for the first and second step [37]. This "loose" mean-field model leads to further improvement of the performance. It is not obvious how to take advantage of a loose model in the case of Gaussian potentials.

### 5.3.2 Material Segmentation

We adopt the method and dataset from [11] for material segmentation task. Their approach proposes the same architecture as in the previous section; a CNN to predict the material labels (*e.g.*, wool, glass, sky, etc.) followed by a densely connected CRF using Gaussian potentials and mean-field inference. We re-use their pre-trained CNN that is publicly available and further choose the CRF parameters and Lab color/position features as in [11]. We compare the Gaussian pairwise potentials versus learned potentials with the same training procedure as for semantic segmentation. Results for pixel accuracy and class-averaged pixel accuracy are shown in Table 2. Following the CRF validation in [11], we ignored 'other' label for both the training and evaluation. For this dataset, the availability of training data is

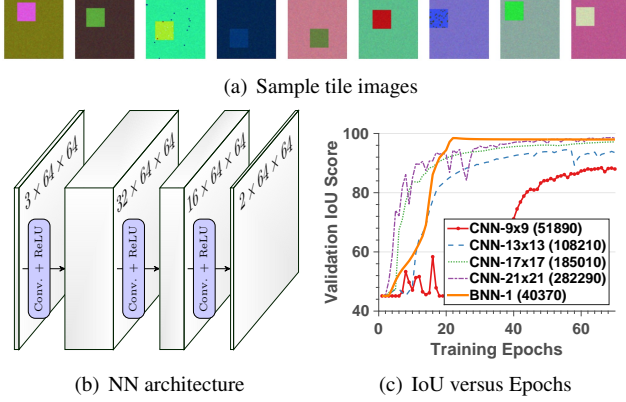(a) Sample tile images



(b) NN architecture

(c) IoU versus Epochs

Figure 5. **Segmenting Tiles.** (a) Example tile input images; (b) the 3-layer NN architecture used in experiments. "Conv" stands for spatial convolutions, resp. bilateral convolutions; (c) Training progress in terms of Validation IoU versus training epochs.

small, since there are 928 images that are only sparsely annotated with segments. While that is enough to cross-validate parameters for Gaussian CRFs, we would expect the general bilateral convolution to benefit from more training data. Especially for class-averaged results we anticipate a higher increase in performance, which we believe is hindered through size of the training set. Visual results with annotations are shown in Fig. 4 and more are included in Sec. C.6.

## 6. Bilateral Neural Networks

Probably the most promising opportunity for the generalized bilateral filter is its use in Convolutional Neural Networks. Since we are not restricted to the Gaussian case anymore, we can run several filters sequentially in the same way as filters are ordered in layers in typical spatial CNN architectures. Having the gradients available allows for end-to-end training with backpropagation, without the need for any change in CNN training protocols. We refer to the layers of bilateral filters as "bilateral convolution layers" (BCL). As discussed in the introduction, these can be understood as either linear filters in a high dimensional space or a filter with an image adaptive receptive field. In the remainder we will refer to CNNs that include at least one bilateral convolutional layer as a bilateral neural network (BNN).

What are the possibilities of a BCL compared to a standard spatial layer? First, we can define a feature space $\mathbf{f}_i \in \mathbb{R}^d$ to define proximity between elements to perform the convolution. This can include color or intensity as in the previous example. Second, since the permutohedral lattice convolution takes advantage of hashing it is advantageous for data in a sparse format. We performed a runtime comparison between our implementation of a BCL[1] and the

---

[1] We have identified several possible speed-ups that are not taken advantage of yet.

caffe [30] implementation of a $d$-dimensional convolution. For 2D positional features (first row), the standard layer is faster since the permutohedral algorithm comes with an overhead. For higher dimensions $d > 2$, the runtime depends on the sparsity; but ignoring the sparsity is quickly leading to intractable runtimes. The permutohedral lattice convolution was designed especially for this purpose, although with Gaussian filtering in mind.

| Dim.-Features | d-dim caffe | BCL |
|---|---|---|
| 2D-$(x, y)$ | **3.3 ± 0.3 / 0.5± 0.1** | 4.8 ± 0.5 / 2.8 ± 0.4 |
| 3D-$(r, g, b)$ | 364.5 ± 43.2 / 12.1 ± 0.4 | **5.1 ± 0.7 / 3.2 ± 0.4** |
| 4D-$(x, r, g, b)$ | 30741.8 ± 9170.9 / 1446.2 ± 304.7 | **6.2 ± 0.7 / 3.8 ± 0.5** |
| 5D-$(x, y, r, g, b)$ | out of memory | **7.6 ± 0.4 / 4.5 ± 0.4** |

Table 3. **Runtime comparison: BCL vs. spatial convolution.** Average CPU/GPU runtime (in ms) of 50 1-neighborhood filters averaged over 1000 images from Pascal VOC. All scaled features $(x, y, r, g, b) \in [0, 50]$. BCL includes splatting and splicing operations which in layered networks can be re-used.

Next we illustrate two use cases of BNNs and compare against spatial CNNs. Section B.1 contains further explanatory experiments with examples on MNIST digit recognition.

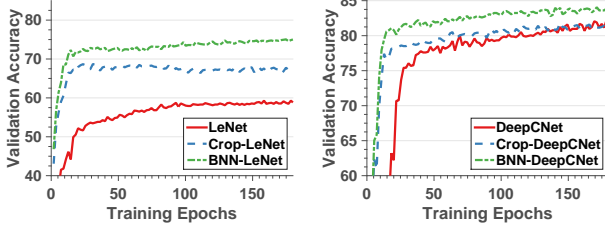### 6.1. An Illustrative Example: Segmenting Tiles

In order to highlight the model possibilities of using higher dimensional sparse feature spaces for convolutions through BCLs, we constructed the following illustrative problem. A randomly colored foreground tile with size $20 \times 20$ is placed on a random colored background of size $64 \times 64$ with additional Gaussian noise with std. dev. $0.02$ added with color values normalized to $[0, 1]$. The task is to segment out the smaller tile. Some example images are shown in Fig. 5(a). A pixel classifier can not distinguish foreground from background since the color is random. We use a three layer CNN with convolutions interleaved with ReLU functions. The schematic of the architecture is shown in Fig 5(b) ($32, 16, 2$ filters). We train standard CNNs, with filters of size $n \times n, n \in \{9, 13, 17, 21\}$. We create 10k training, 1k validation and 1k test images and, use the validation set to choose learning rates. In Fig. 5(c) we plot the validation IoU as a function of training epochs. All CNNs (and BNNs) converge to almost perfect test predictions.

Now, we replace all spatial convolutions with bilateral convolutions for a full BNN. The features are $\mathbf{f}_i = (x_i, y_i, r_i, g_i, b_i)^\top$ and the filter has a neighborhood of 1. The total number of parameters in this network is around $40k$ compared to $52k$ for $9 \times 9$ up to $282k$ for a $21 \times 21$ CNN. With the same training protocol and optimizer, the convergence rate of BNN is much faster. In this example as in semantic segmentation discussed in the last section, color is a discriminative feature for the label. The bilateral convolutions "see" the color difference, the points are al-
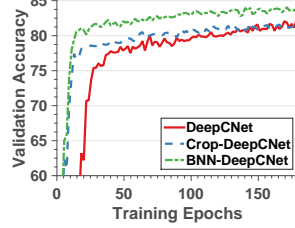
(a) Sample Assamese character images (9 classes, 2 samples each)



(b) LeNet training



(c) DeepCNet training

Figure 6. **Character recognition.** (a) Sample Assamese character images [7]; and training progression of various models with (b) LeNet and (c) DeepCNet base networks.

| | LeNet | Crop-LeNet | BNN-LeNet | DeepC-Net | Crop-DeepCNet | BNN-DeepCNet |
|---|---|---|---|---|---|---|
| Validation | 59.29 | 68.67 | **75.05** | 82.24 | 81.88 | **84.15** |
| Test | 55.74 | 69.10 | **74.98** | 79.78 | 80.02 | **84.21** |

Table 4. **Results on Assamese character images.** Total recognition accuracy for the different models.

ready pre-grouped in the permutohedral lattice and the task remains to assign a label to the two groups.

## 6.2. Character Recognition

The results for tile, semantic, and material segmentations when using general bilateral filters mainly improved because the feature space was used to encode useful prior information about the problem (similar RGB of close-by pixels have the same label). Such prior knowledge is often available when structured predictions are to be made but the input signal may also be in a sparse format to begin with. Let us consider handwritten character recognition, one of the prime cases for CNN use.

The Assamese character dataset [7] contains 183 different Indo-Aryan symbols with 45 writing samples per class. Some sample character images are shown in Fig. 6(a). This dataset has been collected on a tablet PC using a pen input device and has been pre-processed to binary images of size $96 \times 96$. Only about $3\%$ of the pixels contain a pen stroke, which we will denote by $v_i = 1$.

A CNN is a natural choice to approach this classification task. We experiment with two CNN architectures for this experiment that have been used for this task, LeNet-7 from [35] and DeepCNet [17, 23]. The LeNet is a shallower network with bigger filter sizes whereas DeepCNet is deeper with smaller convolutions. Both networks are fully specified in Sec. C.4. In order to simplify the task for the networks we cropped the characters by placing a tight bounding box around them and providing the bounding boxes as input to the networks. We will call these networks "Crop-"LeNet, resp. DeepCNet. For training, we randomly divided the data into 30 writers for training, 6 for validation and the remaining 9 for test. Fig.6(b) and Fig. 6(c) show the training progress for various LeNet and DeepCNet models respectively. DeepCNet is a better choice for this problem

and for both cases pre-processing the data by cropping improves convergence.

The input is spatially sparse and the BCL provides a natural way to take advantage of this. For both networks we create a BNN variant (BNN-LeNet and BNN-DeepCNet) by replacing the first layer with a bilateral convolutions using the features $\mathbf{f}_i = (x_i, y_i)^\top$ and we *only* consider the foreground points $v_i = 1$. The values $(x_i, y_i)$ denote the position of the pixel with respect to the top-left corner of the bounding box around the character. In effect the lattice is very sparse which reduces runtime because the convolutions are only performed on $3\%$ of the points that are actually observed. A bilateral filter has 7 parameters compared to a receptive field of $3 \times 3$ for the first DeepCNet layer and $5 \times 5$ for the first LeNet layer. Thus, a BCL with the same number of filters has fewer parameters. The result of the BCL convolution is then splatted at all points $(x_i, y_i)$ and passed on to the remaining spatial layers. The convergence behaviour is shown in Fig.6 and again we find faster convergence and also better validation accuracy. The empirical results of this experiment for all tested architectures are summarized in Table 4, with BNN variants clearly outperforming their spatial counterparts.

The absolute results can be vastly improved by making use of virtual examples, *e.g.*, by affine transformations [23]. The purpose of these experiments is to compare the networks on equal grounds while we believe that additional data will be beneficial for both networks. We have no reason to believe that a particular network benefits more.

## 7. Conclusion

We proposed to learn bilateral filters from data. In hindsight, it may appear obvious that this leads to performance improvements compared to a particular chosen parametric form, *e.g.*, the Gaussian. To understand the algorithms that facilitate fast approximative computation of Eq. (1) as a parameterized implementation of a bilateral filter that has free parameters is the key insight and enables gradient descent based learning. We relaxed the non-separability in the algorithm from [2] to allow for more general filter functions. There is a wide range of possible applications for learned bilateral filters [41]. In this paper we discussed some examples that generalize previous work. These include joint bilateral upsampling, inference in dense CRFs, and the use of bilateral convolutions in neural networks. On several and diverse experiments, we observed empirical improvements

when learning bilateral filters.

The bilateral convolutional layer allows for filters whose receptive field change given the input image. The feature space view provides a canonical way to encode similarity between any kind of objects, not only pixels, but *e.g.*, bounding boxes, segmentations, surfaces, etc. The proposed filtering operation is then a natural candidate to define a filter convolutions on these objects, it takes advantage of sparsity and scales to higher dimensions. Therefore, we believe that this view will be useful for several problems where CNNs can be applied. An open research problem is whether the sprase higher dimensional structure also allows for efficient or compact representations for intermediate layers inside CNN architectures.

# References

[1] Google Images. https://images.google.com/. [Online; accessed 1-March-2015]. 4

[2] A. Adams, J. Baek, and M. A. Davis. Fast high-dimensional filtering using the permutohedral lattice. In *Computer Graphics Forum*, volume 29, pages 753–762. Wiley Online Library, 2010. 2, 3, 4, 8, 10, 11

[3] A. Adams, N. Gelfand, J. Dolson, and M. Levoy. Gaussian kd-trees for fast high-dimensional filtering. In *ACM Transactions on Graphics (TOG)*, volume 28, page 21. ACM, 2009. 2

[4] P. Arbeláez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(5):898–916, May 2011. 12

[5] V. Aurich and J. Weule. Non-linear gaussian filters performing edge preserving diffusion. In *DAGM*, pages 538–545. Springer, 1995. 1, 2, 12

[6] S. P. Awate and R. T. Whitaker. Higher-order image statistics for unsupervised, information-theoretic, adaptive, image filtering. In *Computer Vision and Pattern Recognition, 2005. IEEE Computer Society Conference on*, volume 2, pages 44–51. 2

[7] K. Bache and M. Lichman. UCI machine learning repository, 2013. 8

[8] J. T. Barron, A. Adams, Y. Shih, and C. Hernández. Fast bilateral-space stereo for synthetic defocus. *CVPR*, 2015. 2

[9] J. T. Barron and B. Poole. The Fast Bilateral Solver. *ArXiv e-prints*, Nov. 2015. 2

[10] S. Bell, K. Bala, and N. Snavely. Intrinsic images in the wild. *ACM Transactions on Graphics (TOG)*, 33(4):159, 2014. 5

[11] S. Bell, P. Upchurch, N. Snavely, and K. Bala. Material recognition in the wild with the materials in context database. *Computer Vision and Pattern Recognition (CVPR)*, 2015. 2, 5, 6

[12] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013. 2

[13] A. Buades, B. Coll, and J.-M. Morel. A non-local algorithm for image denoising. In *Computer Vision and Pattern Recognition, 2005. IEEE Computer Society Conference on*, volume 2, pages 60–65. 2

[14] H. C. Burger, C. J. Schuler, and S. Harmeling. Image denoising: Can plain neural networks compete with BM3D? In *Computer Vision and Pattern Recognition (CVPR)*, 2012. 12

[15] N. D. Campbell, K. Subr, and J. Kautz. Fully-connected CRFs with non-parametric pairwise potential. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 1658–1665. 2, 6

[16] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected CRFs. *arXiv preprint arXiv:1412.7062*, 2014. 2, 6, 14

[17] D. Ciresan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 3642–3649. 8, 14, 16

[18] J. Domke. Learning graphical model parameters with approximate marginal inference. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(10):2454–2467, 2013. 2, 6

[19] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in Neural Information Processing Systems*, pages 2366–2374, 2014. 4, 5

[20] M. Everingham, L. V. Gool, C. Williams, J. Winn, and A. Zisserman. The PASCAL VOC2012 challenge results. 2012. 4, 6, 19

[21] D. Ferstl, M. Ruether, and H. Bischof. Variational depth superresolution using example-based edge representations. In *Proceedings International Conference on Computer Vision (ICCV), IEEE*, December 2015. 5

[22] E. S. Gastal and M. M. Oliveira. Adaptive manifolds for real-time high-dimensional filtering. *ACM Transactions on Graphics (TOG)*, 31(4):33, 2012. 2

[23] B. Graham. Spatially-sparse convolutional neural networks. *arXiv preprint arXiv:1409.6070*, 2014. 8, 14, 16

[24] B. Graham. Sparse 3D convolutional neural networks. *arXiv preprint arXiv:1505.02890*, 2015. 2

[25] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik. Semantic contours from inverse detectors. In *International Conference on Computer Vision (ICCV)*, 2011. 19

[26] K. He, J. Sun, and X. Tang. Guided image filtering. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(6):1397–1409, 2013. 2

[27] H. Hu and G. de Haan. Trained bilateral filters and applications to coding artifacts reduction. In *Image Processing, 2007. IEEE International Conference on*, volume 1, pages I–325. 2

[28] C. Ionescu, O. Vantzos, and C. Sminchisescu. Matrix backpropagation for deep networks with structured layers. *ICCV*, 2015. 2

[29] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial transformer networks. *arXiv preprint arXiv:1506.02025*, 2015. 2

[30] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the ACM International Conference on Multimedia*, pages 675–678. ACM, 2014. 4, 7, 11, 12

[31] M. Kiefel and P. V. Gehler. Human pose estimation with fields of parts. In *Computer Vision–ECCV 2014*, pages 331–346. Springer, 2014. 2

[32] J. Kopf, M. F. Cohen, D. Lischinski, and M. Uyttendaele. Joint bilateral upsampling. *ACM Transactions on Graphics (TOG)*, 26(3):96, 2007. 4

[33] P. Krähenbühl and V. Koltun. Efficient inference in fully connected CRFs with Gaussian edge potentials. *arXiv preprint arXiv:1210.5644*, 2012. 2, 5

[34] P. Krähenbühl and V. Koltun. Parameter learning and convergent inference for dense random fields. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pages 513–521, 2013. 2, 6

[35] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 8, 11, 12

[36] Y. LeCun, C. Cortes, and C. J. Burges. The mnist database of handwritten digits, 1998. 11, 14, 16

[37] Y. Li and R. Zemel. Mean-field networks. *arXiv preprint arXiv:1410.5884*, 2014. 6

[38] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015. 12, 13

[39] P. Milanfar. A tour of modern image filtering. *IEEE Signal Processing Magazine*, 2, 2011. 1, 2

[40] S. Paris and F. Durand. A fast approximation of the bilateral filter using a signal processing approach. In *Computer Vision–ECCV 2006*, pages 568–580. Springer, 2006. 1, 2, 3

[41] S. Paris, P. Kornprobst, J. Tumblin, and F. Durand. *Bilateral filtering: Theory and applications*. Now Publishers Inc, 2009. 4, 8

[42] J.-H. Rick Chang and Y.-C. Frank Wang. Propagated image filtering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10–18, 2015. 2

[43] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from rgbd images. In *Computer Vision–ECCV 2012*, pages 746–760. Springer, 2012. 4, 5, 14

[44] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 6

[45] S. M. Smith and J. M. Brady. SUSAN – a new approach to low level image processing. *Int. J. Comput. Vision*, 23(1):45–78, May 1997. 1, 2

[46] D. Sun, J. Wulff, E. B. Sudderth, H. Pfister, and M. J. Black. A fully-connected layered model of foreground and background flow. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 2451–2458. 5

[47] H. Takeda, S. Farsiu, and P. Milanfar. Kernel regression for image processing and reconstruction. *Image Processing, IEEE Transactions on*, 16(2):349–366, 2007. 2

[48] J. B. Tenenbaum. A Global Geometric Framework for Nonlinear Dimensionality Reduction. http://isomap.stanford.edu/, 2000. [Online; accessed 12-October-2015]. 13

[49] J. B. Tenenbaum, V. De Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000. 13

[50] C. Tomasi and R. Manduchi. Bilateral filtering for gray and color images. In *Computer Vision, 1998. Sixth International Conference on*, pages 839–846. IEEE, 1998. 1, 2

[51] V. Vineet, G. Sheasby, J. Warrell, and P. H. Torr. Posefield: An efficient mean-field based method for joint estimation of human pose, segmentation and depth. In *EMMCVPR*, 2013. 5

[52] V. Vineet, J. Warrell, and P. H. Torr. Filter-based mean-field inference for random fields with higher order terms and product label-spaces. In *ECCV*, 2012. 5

[53] H. Yagou, Y. Ohtake, and A. Belyaev. Mesh smoothing via mean and median filtering applied to face normals. In *Geometric Modeling and Processing, 2002. Proceedings*, pages 124–131. IEEE, 2002. 13

[54] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus. Deconvolutional networks. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 2528–2535. 5

[55] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr. Conditional random fields as recurrent neural networks. *IEEE International Conference on Computer Vision*, 2015. 2, 5

# Appendix

The appendix contains a more detailed overview of the permutohedral lattice convolution in Section A, more experiments in Section B and additional results with experiment protocols for the experiments presented before in Section C.

## A. General Permutohedral Convolutions

A core technical contribution of this work is the generalization of the Gaussian permutohedral lattice convolution proposed in [2] to the full non-separable case with the ability to perform backpropagation. Although, conceptually, there are minor difference between non-Gaussian and general parameterized filters, there are non-trivial practical differences in terms of the algorithmic implementation. The Gauss filters belong to the separable class and can thus be decomposed into multiple sequential one dimensional convolutions. We are interested in the general filter convolutions, which can not be decomposed. Thus, performing a general permutohedral convolution at a lattice point requires the computation of the inner product with the neighboring elements in all the directions in the high-dimensional space.

Here, we give more details of the implementation differences of separable and non-separable filters. In the following we will explain the scalar case first. Recall, that
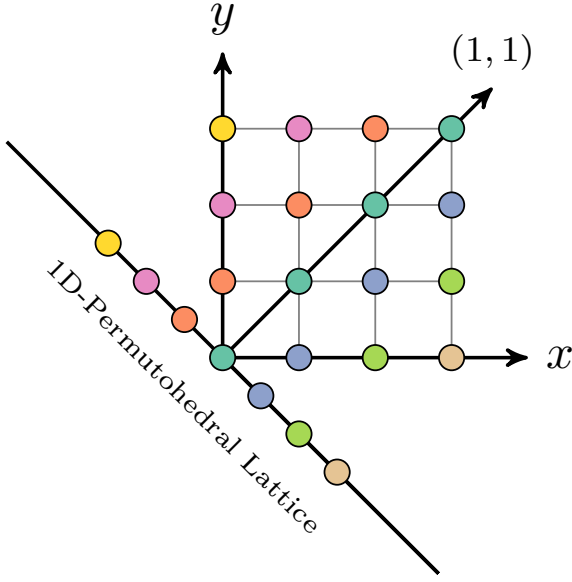
Figure 7. **Illustration of 1D permutohedral lattice construction.** A $4 \times 4$ $(x, y)$ grid lattice is projected onto the plane defined by the normal vector $(1, 1)^\top$. This grid has $s + 1 = 4$ and $d = 2$ $(s + 1)^d = 4^2 = 16$ elements. In the projection, all points of the same color are projected onto the same points in the plane. The number of elements of the projected lattice is $t = (s + 1)^d - s^d = 4^2 - 3^2 = 7$, that is the $(4 \times 4)$ grid minus the size of lattice that is 1 smaller at each size, in this case a $(3 \times 3)$ lattice (the upper right $(3 \times 3)$ elements).

the forward pass of general permutohedral convolution involves 3 steps: *splatting*, *convolving* and *slicing*. We follow the same splatting and slicing strategies as in [2] since these operations do not depend on the filter kernel. The main difference between our work and the existing implementation of [2] is the way that the convolution operation is executed. This proceeds by constructing a *blur neighbor* matrix $K$ that stores for every lattice point all values of the lattice neighbors that are needed to compute the filter output.

The blur neighbor matrix is constructed by traversing through all the populated lattice points and their neighboring elements. This is done recursively to share computations. For any lattice point, the neighbors that are $n$ hops away are the direct neighbors of the points that are $n - 1$ hops away. The size of a $d$ dimensional spatial filter with width $s + 1$ is $(s + 1)^d$ (*e.g.*, a $3 \times 3$ filter, $s = 2$ in $d = 2$ has $3^2 = 9$ elements) and this size grows exponentially in the number of dimensions $d$. The permutohedral lattice is constructed by projecting a regular grid onto the plane spanned by the $d$ dimensional normal vector $(1, \ldots, 1)^\top$. See Fig. 7 for an illustration of 1D lattice construction. Many corners of a grid filter are projected onto the same point, in total $t = (s + 1)^d - s^d$ elements remain in the permuto-

hedral filter with $s$ neighborhood in $d - 1$ dimensions. If the lattice has $m$ populated elements, the matrix $K$ has size $t \times m$. Note that, since the input signal is typically sparse, only a few lattice corners are being populated in the *slicing* step. We use a hash-table to keep track of these points and traverse only through the populated lattice points for this neighborhood matrix construction.

Once the blur neighbor matrix $K$ is constructed, we can perform the convolution by the matrix vector multiplication

$$l' = BK, \tag{8}$$

where $B$ is the $1 \times t$ filter kernel (whose values we will learn) and $l' \in \mathbb{R}^{1 \times m}$ is the result of the filtering at the $m$ lattice points. In practice, we found that the matrix $K$ is sometimes too large to fit into GPU memory and we divided the matrix $K$ into smaller pieces to compute Eq. 8 sequentially.

In the general multi-dimensional case, the signal $l$ is of $c$ dimensions. Then the kernel is of size $B \times t$ and $K$ stores the $c$ dimensional vectors accordingly. When the input and output points are different, we slice only the input points and splat only at the output points.

## B. Additional Experiments

In this section we discuss more use-cases for the learned bilateral filters, one use-case of BNNs and two single filter applications for image and 3D mesh denoising.

### B.1. Recognition of subsampled MNIST

One of the strengths of the proposed filter convolution is that it does not require the input to lie on a regular grid. The only requirement is to define a distance between features of the input signal. We highlight this feature with the following experiment using the classical MNIST ten class classification problem [36]. We sample a sparse set of $N$ points $(x, y) \in [0, 1] \times [0, 1]$ uniformly at random in the input image, use their interpolated values as signal and the *continuous* $(x, y)$ positions as features. This mimics subsampling of a high-dimensional signal. To compare against a spatial convolution, we interpolate the sparse set of values at the grid positions.

We take a reference implementation of LeNet [35] that is part of the Caffe project [30] and compare it against the same architecture but replacing the first convolutional layer with a bilateral convolution layer (BCL). The filter size and numbers are adjusted to get a comparable number of parameters ($5 \times 5$ for LeNet, 2-neighborhood for BCL).

The results are shown in Table 5. We see that training on the original MNIST data (column Original, LeNet vs. BNN) leads to a slight decrease in performance of the BNN (99.03%) compared to LeNet (99.19%). The BNN can be trained and evaluated on sparse signals, and we resample the

| Method | | Test Subsampling | | |
| --- | --- | --- | --- | --- |
| | Original | 100% | 60% | 20% |
| LeNet | **0.9919** | 0.9660 | 0.9348 | **0.6434** |
| BNN | 0.9903 | **0.9844** | **0.9534** | 0.5767 |
| LeNet 100% | 0.9856 | 0.9809 | 0.9678 | **0.7386** |
| BNN 100% | **0.9900** | 0.9863 | 0.9699 | 0.6910 |
| LeNet 60% | 0.9848 | 0.9821 | 0.9740 | 0.8151 |
| BNN 60% | **0.9885** | **0.9864** | **0.9771** | **0.8214** |
| LeNet 20% | **0.9763** | **0.9754** | 0.9695 | 0.8928 |
| BNN 20% | 0.9728 | 0.9735 | **0.9701** | **0.9042** |

Table 5. Classification accuracy on MNIST. We compare the LeNet [35] implementation that is part of Caffe [30] to the network with the first layer replaced by a bilateral convolution layer (BCL). Both are trained on the original image resolution (first two rows). Three more BNN and CNN models are trained with randomly subsampled images (100%, 60% and 20% of the pixels). An additional bilinear interpolation layer samples the input signal on a spatial grid for the CNN model.

image as described above for $N = 100\%$, 60% and 20% of the total number of pixels. The methods are also evaluated on test images that are subsampled in the same way. Note that we can train and test with different subsampling rates. We introduce an additional bilinear interpolation layer for the LeNet architecture to train on the same data. In essence, both models perform a spatial interpolation and thus we expect them to yield a similar classification accuracy. Once the data is of higher dimensions the permutohedral convolution will be faster due to hashing the sparse input points, as well as less memory demanding in comparison to naive application of a spatial convolution with interpolated values.

## B.2. Image Denoising

The main application that inspired the development of the bilateral filtering operation is image denoising [5], there using a single Gaussian kernel. Our development allows to learn this kernel function from data and we explore how to improve using a *single* but more general bilateral filter.

We use the Berkeley segmentation dataset (BSDS500) [4] as a test bed. The color images in the dataset are converted to gray-scale, and corrupted with Gaussian noise with a standard deviation of $\frac{25}{255}$.

We compare the performance of four different filter models on a denoising task. The first baseline model ("Spatial" in Table 6, 25 weights) uses a single spatial filter with a kernel size of 5 and predicts the scalar gray-scale value at the center pixel. The next model ("Gauss Bilateral") applies a bilateral *Gaussian* filter to the noisy input, using position and intensity features $\mathbf{f} = (x, y, v)^{\top}$. The third setup ("Learned Bilateral", 65 weights) takes a Gauss kernel as initialization and fits all filter weights on the "train" image set to minimize the mean squared error with respect to the clean images. We run a combination of spatial and permutohedral convolutions on spatial and bilateral features ("Spa-



Figure 8. **Sample data for 3D mesh denoising.** (top) Some 3D body meshes sampled from [38] and (bottom) the corresponding noisy meshes used in denoising experiments.

tial + Bilateral (Learned)") to check for a complementary performance of the two convolutions.

| Method | PSNR |
| --- | --- |
| Noisy Input | 20.17 |
| Spatial | 26.27 |
| Gauss Bilateral | 26.51 |
| Learned Bilateral | 26.58 |
| Spatial + Bilateral (Learned) | 26.65 |

Table 6. PSNR results of a denoising task using the BSDS500 dataset [4]

The PSNR scores evaluated on full images of the "test" image set are shown in Table 6. We find that an untrained bilateral filter already performs better than a trained spatial convolution (26.27 to 26.51). A learned convolution further improve the performance slightly. We chose this simple one-kernel setup to validate an advantage of the generalized bilateral filter. A competitive denoising system would employ RGB color information and also needs to be properly adjusted in network size. Multi-layer perceptrons have obtained state-of-the-art denoising results [14] and the permutohedral lattice layer can readily be used in such an architecture, which is intended future work.

## B.3. 3D Mesh Denoising

Permutohedral convolutions can naturally be extended to higher ($> 2$) dimensional data. To highlight this, we use the proposed convolution for the task of denoising 3D meshes.

Figure 9. **4D isomap features for 3D human bodies.** Visualization of 4D isomap features for a sample 3D mesh. Isomap feature values are overlaid onto mesh vertices.

We sample 3D human body meshes using a generative 3D body model from [38]. To the clean meshes, we add Gaussian random noise displacements along the surface normal at each vertex location. Figure 8 shows some sample 3D meshes sampled from [38] and corresponding noisy meshes. The task is to take the noisy meshes as inputs and recover the original 3D body meshes. We create 1000 training, 200 validation and another 500 testing examples for the experiments.

**Mesh Representation:** The 3D human body meshes from [38] are represented with 3D vertex locations and the edge connections between the vertices. We found that this signal representation using global 3D coordinates is not suitable for denoising with bilateral filtering. Therefore, we first smooth the noisy mesh using mean smoothing applied to the face normals [53] and represent the noisy mesh vertices as 3D vector displacements with respect to the corresponding smoothed mesh. Thus, the task becomes denoising the 3D vector displacements with respect to the smoothed mesh.

**Isomap Features:** To apply permutohedral convolution, we need to define features at each input vertex point. We use 4 dimensional isomap embedding [49] of the given 3D mesh graph as features. The given 3D mesh is converted into weighted edge graph with edge weights set to Euclidean distance between the connected vertices and to infinity between the non-connected vertices. Then 4 dimensional isomap embedding is computed for this weighted edge graph using a publicly available implementation [48]. Fig. 9 shows the visualization of isomap features on a sample 3D mesh.

**Experimental Results:** Mesh denoising with a bilateral filter proceeds by splatting the input 3D mesh vectors (displacements with respect to smoothed mesh) into the 4D isomap feature space, filtering the signal in this 4D space

| | Noisy Mesh | Normal Smoothing | Gauss Bilateral | Learned Bilateral |
|---|---|---|---|---|
| Vertex Distance (RMSE) | 5.774 | 3.183 | 2.872 | **2.825** |
| Normal Angle Error | 19.680 | 19.707 | 19.357 | **19.207** |

Table 7. **Body Denoising.** Vertex distance RMSE values and normal angle error (in degrees) corresponding to different denoising strategies averaged over 500 test meshes.
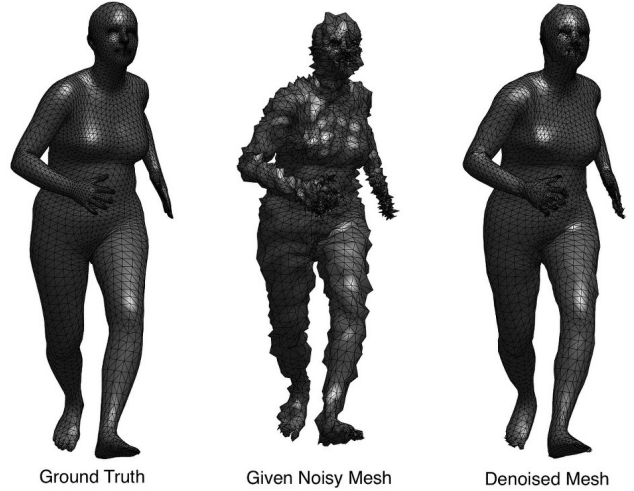


Figure 10. **Sample Denoising Result.** Ground truth mesh (left), corresponding given noisy mesh (middle) and the denoised result (right) using learned bilateral filter.

and then slicing back into original 3D input space. The Table 7 shows quantitative results as RMSE for different denoising strategies. The normal smoothing [53] already reduces the RMSE. The Gauss bilateral filter results in significant improvement over normal smoothing with and learning the filter weights again improves the result. A visual result is shown in Figure 10.

## C. Additional results

This section contains more qualitative results for the experiments of the main paper.

### C.1. Lattice Visualization

Figure 11 shows sample lattice visualizations for different feature spaces.

### C.2. Color Upsampling

In addition to the experiments discussed in the main paper, we performed the cross-factor analysis of training and testing at different upsampling factors. Table 8 shows the PSNR results for this analysis. Although, in terms of PSNR, it is optimal to train and test at the same upsampling factor, the differences are small when training and testing upsam-
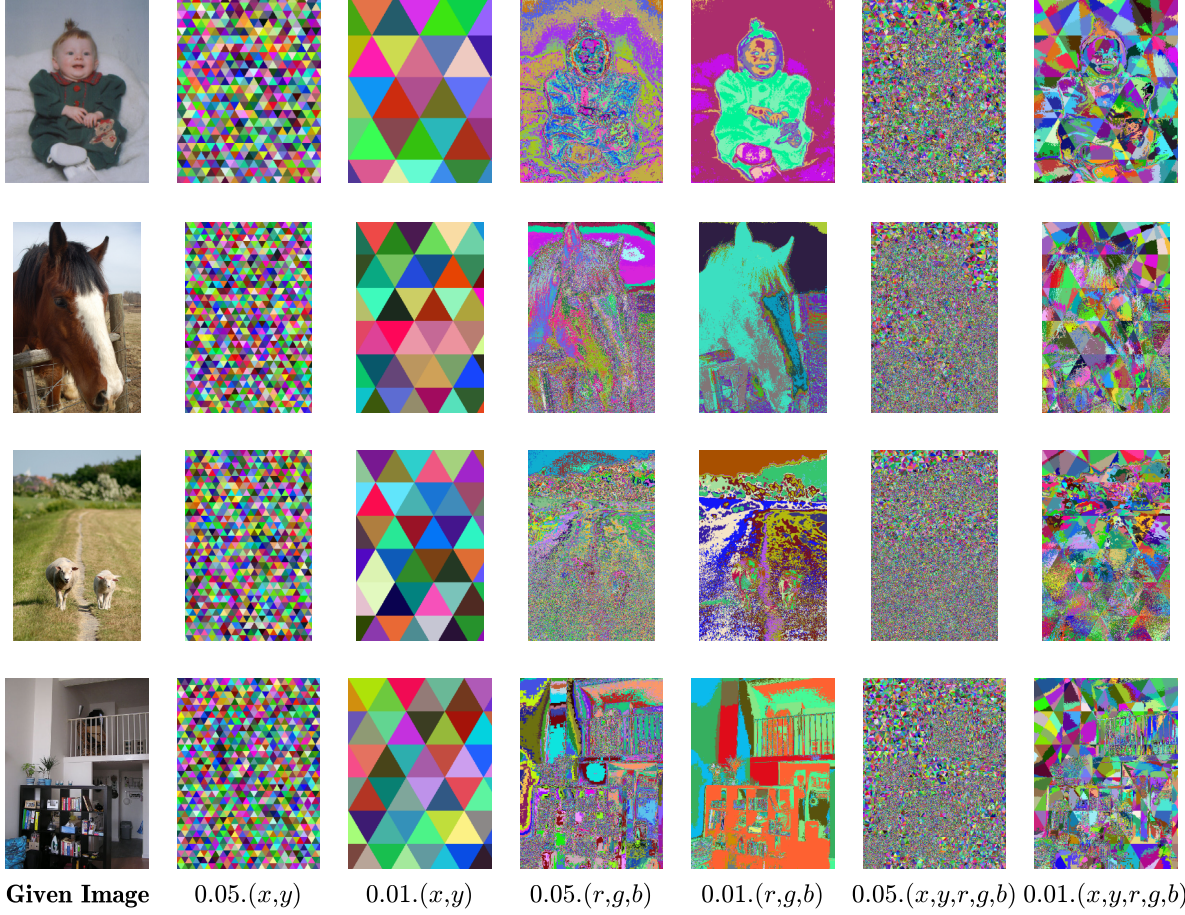
Figure 11. **Visualization of the Permutohedral Lattice.** Sample lattice visualizations for different feature spaces. All pixels falling in the same simplex cell are shown with the same color. $(x, y)$ features correspond to image pixel positions, and $(r, g, b) \in [0, 255]$ correspond to the red, green and blue color values.

|  |  | Test Factor | | |
|  | | $2\times$ | $4\times$ | $8\times$ | $16\times$ |
|---|---|---|---|---|---|
| | $2\times$ | **38.45** | 36.12 | 34.06 | 32.43 |
| Train Factor | $4\times$ | 38.40 | **36.16** | **34.08** | 32.47 |
| | $8\times$ | 38.40 | 36.15 | **34.08** | 32.47 |
| | $16\times$ | 38.26 | 36.13 | 34.06 | **32.49** |

Table 8. **Color Upsampling with different train and test upsampling factors.** PSNR values corresponding to different upsampling factors used at train and test times on the 2 megapixel image dataset, using our learned bilateral filters.

pling factors are different. Some images of the upsampling for the Pascal VOC12 dataset are shown in Fig. 12. It is especially the low level image details that are better preserved with a learned bilateral filter compared to the Gaussian case.

### C.3. Depth Upsampling

Figure 13 presents some more qualitative results comparing bicubic interpolation, Gauss bilateral and learned bilateral upsampling on NYU depth dataset image [43].

### C.4. Character Recognition

Figure 14 shows the schematic of different layers of the network architecture for LeNet-7 [36] and DeepCNet(5, 50) [17, 23]. For the BNN variants, the first layer filters are replaced with learned bilateral filters and are learned end-to-end.

### C.5. Semantic Segmentation

Some more visual results for semantic segmentation are shown in Figure 15. These include the underlying DeepLab CNN[16] result (DeepLab), the 2 step mean-field result with Gaussian edge potentials (+2stepMF-GaussCRF) and also corresponding results with learned edge potentials (+2stepMF-LearndCRF). In general, we observe that mean-field in learned CRF leads to slightly dilated classification regions in comparison to using Gaussian CRF thereby filling-in the false negative pixels and also correcting some mis-classified regions.
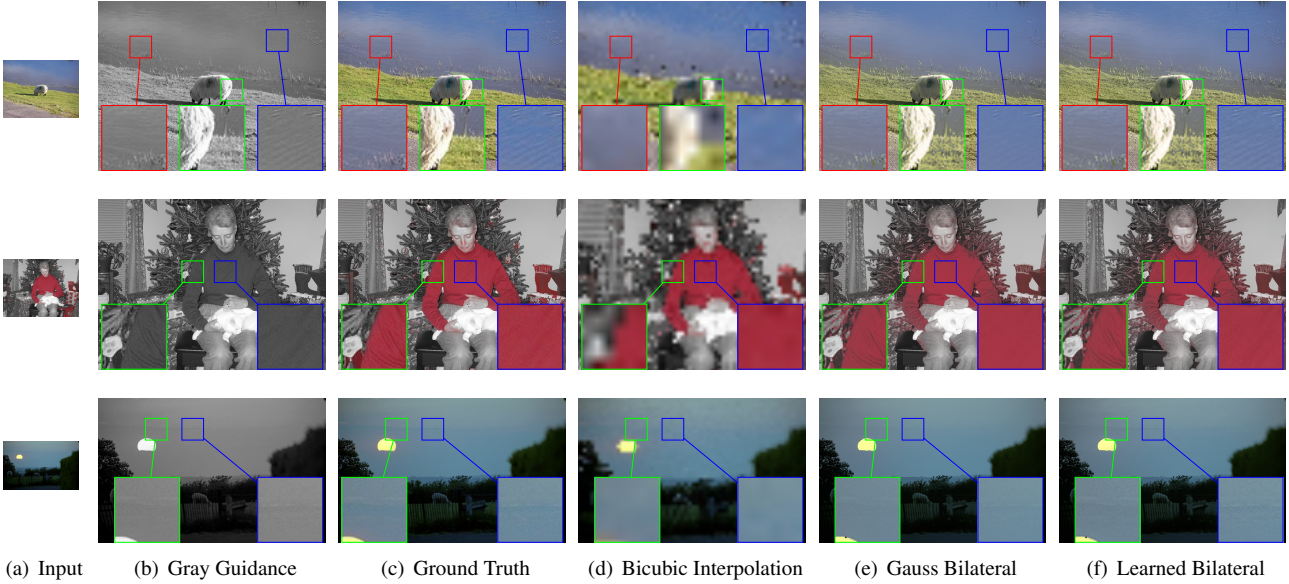
| (a) Input | (b) Gray Guidance | (c) Ground Truth | (d) Bicubic Interpolation | (e) Gauss Bilateral | (f) Learned Bilateral |

Figure 12. **Color Upsampling.** Color 8× upsampling results using different methods (best viewed on screen).



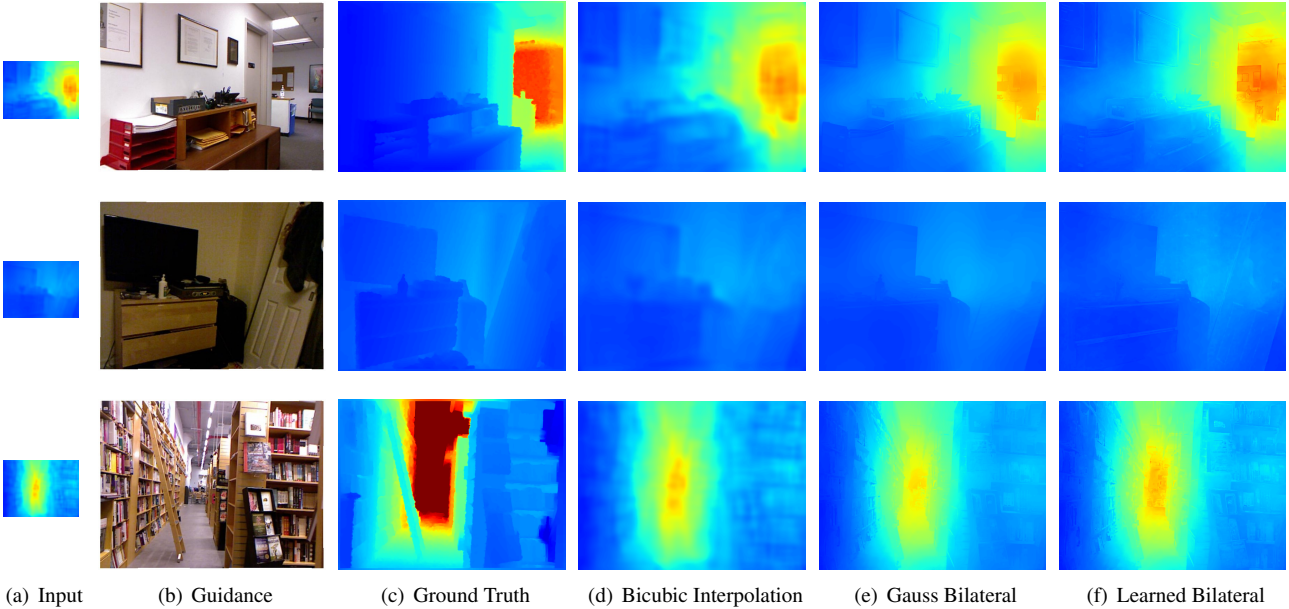| (a) Input | (b) Guidance | (c) Ground Truth | (d) Bicubic Interpolation | (e) Gauss Bilateral | (f) Learned Bilateral |

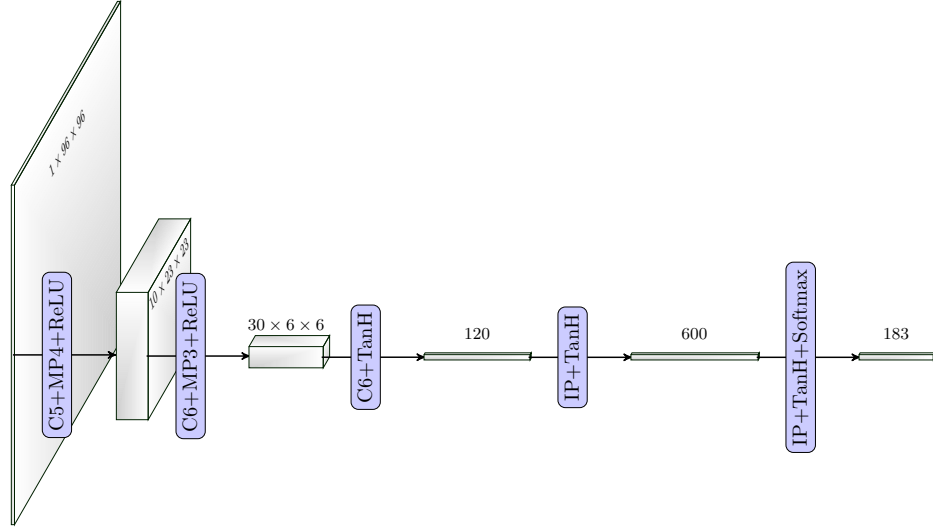Figure 13. **Depth Upsampling.** Depth 8× upsampling results using different upsampling strategies.

## C.6. Material Segmentation

In Fig. 16, we present visual results comparing 2 step mean-field inference with Gaussian and learned pairwise CRF potentials. In general, we observe that the pixels belonging to dominant classes in the training data are being more accurately classified with learned CRF. This leads to a significant improvements in overall pixel accuracy. This also results in a slight decrease of the accuracy from less frequent class pixels thereby slightly reducing the average class accuracy with learn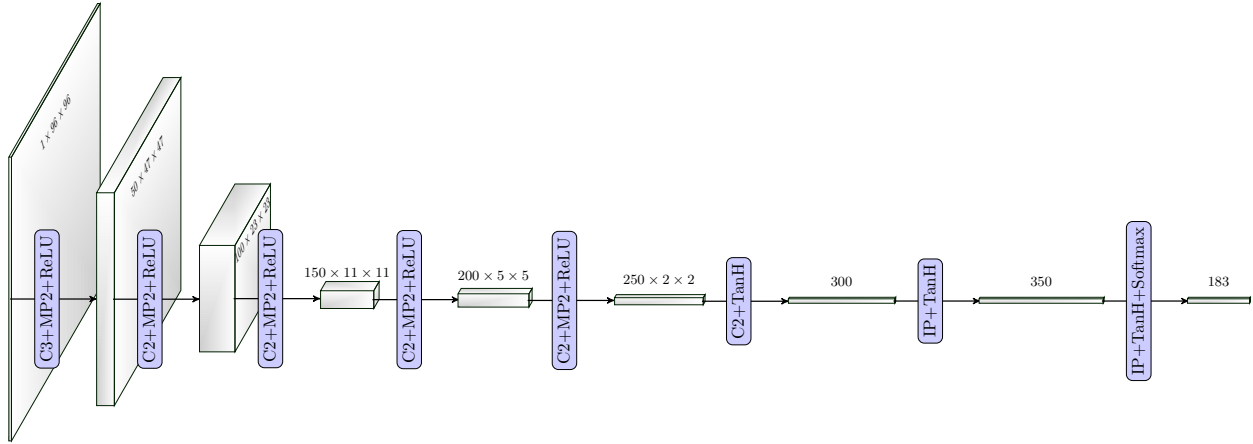ing. We attribute this to the type of annotation that is available for this dataset, which is not for the entire image but for some segments in the image. We have very few images of the infrequent classes to combat this behaviour during training.

## C.7. Experiment Protocols

Table 9 shows experiment protocols of different experiments.

(a) LeNet-7



(b) DeepCNet

Figure 14. **CNNs for Character Recognition.** Schematic of (top) LeNet-7 [36] and (bottom) DeepCNet(5,50) [17, 23] architectures used in Assamese character recognition experiments.

| (a) Input | (b) Ground Truth | (c) DeepLab | (d) +2stepMF-GaussCRF | (e) +2stepMF-LearnedCRF |

Figure 15. **Semantic Segmentation.** Example results of semantic segmentation. (c) depicts the unary results before application of MF, (d) after two steps of MF with Gaussian edge CRF potentials, (e) after two steps of MF with learned edge CRF potentials.

| Brick | Carpet | Ceramic | Fabric | Foliage | Food | Glass | Hair |
| Leather | Metal | Mirror | Other | Painted | Paper | Plastic | Polished Stone |
| Skin | Sky | Stone | Tile | Wallpaper | Water | Wood |

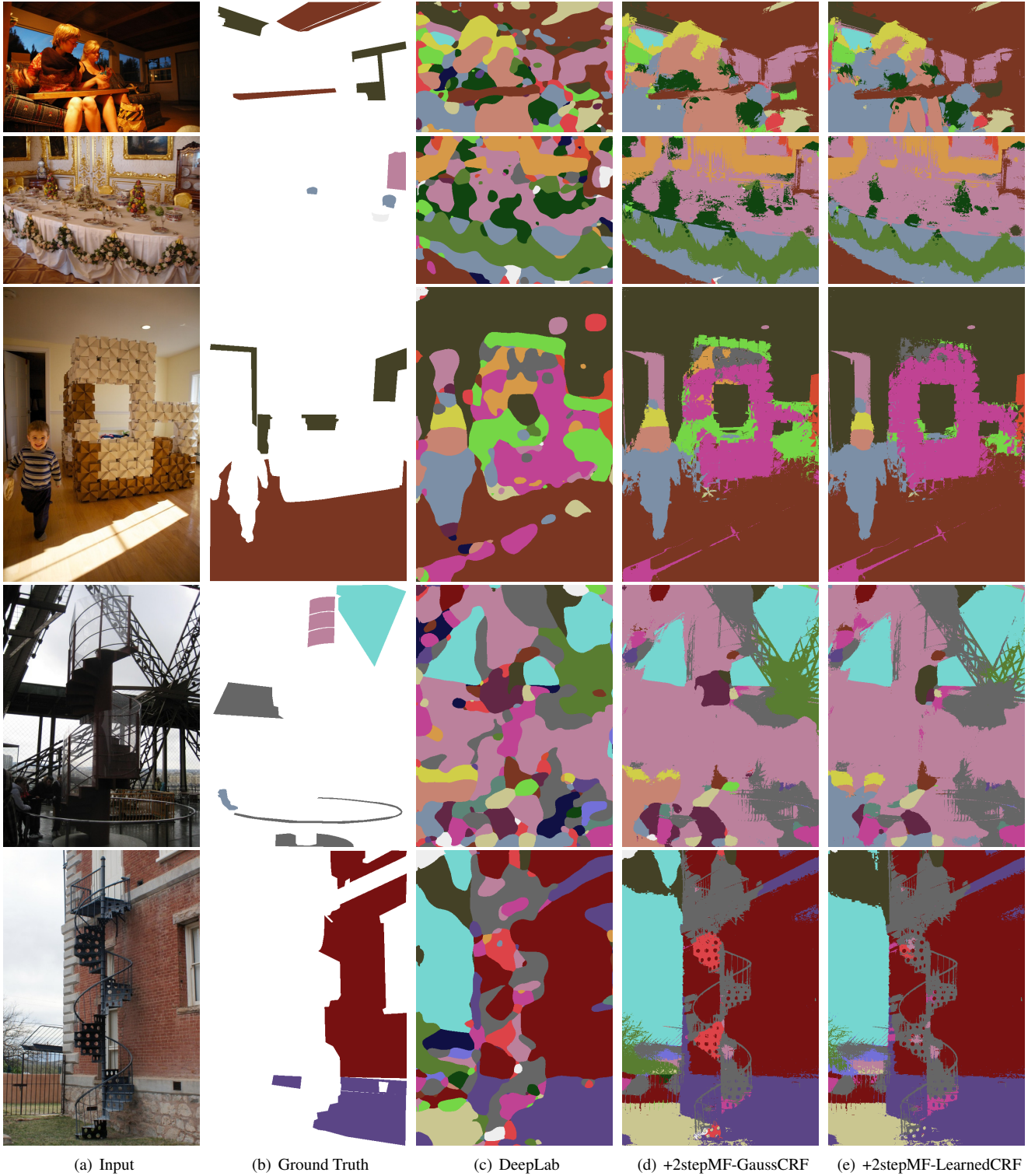(a) Input    (b) Ground Truth    (c) DeepLab    (d) +2stepMF-GaussCRF    (e) +2stepMF-LearnedCRF

Figure 16. **Material Segmentation.** Example results of material segmentation. (c) depicts the unary results before application of MF, (d) after two steps of MF with Gaussian edge CRF potentials, (e) after two steps of MF with learned edge CRF potentials.

| | | | Filter | Filter | Data Statistics | | | Training Protocol | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Experiment | Feature Types | Feature Scales | Size | Nbr. | Train | Val. | Test | Loss Type | LR | Batch | Epochs |
| **Single Bilateral Filter Applications** | | | | | | | | | | | |
| **2× Color Upsampling** | $Position_1$, Intensity (3D) | 0.13, 0.17 | 65 | 2 | 10581 | 1449 | 1456 | MSE | 1e-06 | 200 | 94.5 |
| **4× Color Upsampling** | $Position_1$, Intensity (3D) | 0.06, 0.17 | 65 | 2 | 10581 | 1449 | 1456 | MSE | 1e-06 | 200 | 94.5 |
| **8× Color Upsampling** | $Position_1$, Intensity (3D) | 0.03, 0.17 | 65 | 2 | 10581 | 1449 | 1456 | MSE | 1e-06 | 200 | 94.5 |
| **16× Color Upsampling** | $Position_1$, Intensity (3D) | 0.02, 0.17 | 65 | 2 | 10581 | 1449 | 1456 | MSE | 1e-06 | 200 | 94.5 |
| **Depth Upsampling** | $Position_1$, Color (5D) | 0.05, 0.02 | 665 | 2 | 795 | 100 | 654 | MSE | 1e-07 | 50 | 251.6 |
| **Mesh Denoising** | Isomap (4D) | 46.00 | 63 | 2 | 1000 | 200 | 500 | MSE | 100 | 10 | 100.0 |
| **DenseCRF Applications** | | | | | | | | | | | |
| **Semantic Segmentation** | | | | | | | | | | | |
| **- 1step MF** | $Position_1$, Color (5D); $Position_1$ (2D) | 0.01, 0.34; 0.34 | 665; 19 | 2; 2 | 10581 | 1449 | 1456 | Logistic | 0.1 | 5 | 1.4 |
| **- 2step MF** | $Position_1$, Color (5D); $Position_1$ (2D) | 0.01, 0.34; 0.34 | 665; 19 | 2; 2 | 10581 | 1449 | 1456 | Logistic | 0.1 | 5 | 1.4 |
| **- *loose* 2step MF** | $Position_1$, Color (5D); $Position_1$ (2D) | 0.01, 0.34; 0.34 | 665; 19 | 2; 2 | 10581 | 1449 | 1456 | Logistic | 0.1 | 5 | +1.9 |
| **Material Segmentation** | | | | | | | | | | | |
| **- 1step MF** | $Position_2$, Lab-Color (5D) | 5.00, 0.05, 0.30 | 665 | 2 | 928 | 150 | 1798 | Weighted Logistic | 1e-04 | 24 | 2.6 |
| **- 2step MF** | $Position_2$, Lab-Color (5D) | 5.00, 0.05, 0.30 | 665 | 2 | 928 | 150 | 1798 | Weighted Logistic | 1e-04 | 12 | +0.7 |
| **- *loose* 2step MF** | $Position_2$, Lab-Color (5D) | 5.00, 0.05, 0.30 | 665 | 2 | 928 | 150 | 1798 | Weighted Logistic | 1e-04 | 12 | +0.2 |
| **Neural Network Applications** | | | | | | | | | | | |
| **Tiles: CNN-9×9** | - | - | 81 | 4 | 10000 | 1000 | 1000 | Logistic | 0.01 | 100 | 500.0 |
| **Tiles: CNN-13×13** | - | - | 169 | 6 | 10000 | 1000 | 1000 | Logistic | 0.01 | 100 | 500.0 |
| **Tiles: CNN-17×17** | - | - | 289 | 8 | 10000 | 1000 | 1000 | Logistic | 0.01 | 100 | 500.0 |
| **Tiles: CNN-21×21** | - | - | 441 | 10 | 10000 | 1000 | 1000 | Logistic | 0.01 | 100 | 500.0 |
| **Tiles: BNN** | $Position_1$, Color (5D) | 0.05, 0.04 | 63 | 1 | 10000 | 1000 | 1000 | Logistic | 0.01 | 100 | 30.0 |
| **LeNet** | - | - | 25 | 2 | 5490 | 1098 | 1647 | Logistic | 0.1 | 100 | 182.2 |
| **Crop-LeNet** | - | - | 25 | 2 | 5490 | 1098 | 1647 | Logistic | 0.1 | 100 | 182.2 |
| **BNN-LeNet** | $Position_2$ (2D) | 20.00 | 7 | 1 | 5490 | 1098 | 1647 | Logistic | 0.1 | 100 | 182.2 |
| **DeepCNet** | - | - | 9 | 1 | 5490 | 1098 | 1647 | Logistic | 0.1 | 100 | 182.2 |
| **Crop-DeepCNet** | - | - | 9 | 1 | 5490 | 1098 | 1647 | Logistic | 0.1 | 100 | 182.2 |
| **BNN-DeepCNet** | $Position_2$ (2D) | 40.00 | 7 | 1 | 5490 | 1098 | 1647 | Logistic | 0.1 | 100 | 182.2 |

Table 9. **Experiment Protocols.** Experiment protocols for the different experiments presented in this work. **Feature Types**: Feature spaces used for the bilateral convolutions. $Position_1$ corresponds to un-normalized pixel positions whereas $Position_2$ corresponds to pixel positions normalized to $[0, 1]$ with respect to the given image. **Feature Scales**: Cross-validated scales for the features used. **Filter Size**: Number of elements in the filter that is being learned. **Filter Nbr.**: Half-width of the filter. **Train**, **Val.** and **Test** corresponds to the number of train, validation and test images used in the experiment. **Loss Type**: Type of loss used for back-propagation. "MSE" corresponds to Euclidean mean squared error loss and "Logistic" corresponds to multinomial logistic loss. "Weighted Logistic" is the class-weighted multinomial logistic loss. We weighted the loss with inverse class probability for material segmentation task due to the small availability of training data with class imbalance. **LR**: Fixed learning rate used in stochastic gradient descent. **Batch**: Number of images used in one parameter update step. **Epochs**: Number of training epochs. In all the experiments, we used fixed momentum of 0.9 and weight decay of 0.0005 for stochastic gradient descent. "'Color Upsampling" experiments in this Table corresponds to those performed on Pascal VOC12 dataset images. For all experiments using Pascal VOC12 images, we use extended training segmentation dataset available from [25], and used standard validation and test splits from the main dataset [20].