

Understanding Convolution for Semantic Segmentation

Panqu Wang¹, Pengfei Chen¹, Ye Yuan², Ding Liu³, Zehua Huang¹, Xiaodi Hou¹, Garrison Cottrell⁴

¹TuSimple, ²Carnegie Mellon University, ³University of Illinois Urbana-Champaign, ⁴UC San Diego
¹{panqu.wang,pengfei.chen,zehua.huang,xiaodi.hou}@tusimple.ai, ²yey1@andrew.cmu.edu,
³dingliu2@illinois.edu, ⁴gary@ucsd.edu

Abstract

Recent advances in deep learning, especially deep convolutional neural networks (CNNs), have led to significant improvement over previous semantic segmentation systems. Here we show how to improve pixel-wise semantic segmentation by manipulating convolution-related operations that are of both theoretical and practical value. First, we design dense upsampling convolution (DUC) to generate pixel-level prediction, which is able to capture and decode more detailed information that is generally missing in bilinear upsampling. Second, we propose a hybrid dilated convolution (HDC) framework in the encoding phase. This framework 1) effectively enlarges the receptive fields (RF) of the network to aggregate global information; 2) alleviates what we call the “gridding issue” caused by the standard dilated convolution operation. We evaluate our approaches thoroughly on the Cityscapes dataset, and achieve a state-of-art result of 80.1% mIoU in the test set at the time of submission. We also have achieved state-of-the-art overall on the KITTI road estimation benchmark and the PASCAL VOC2012 segmentation task. Our source code can be found at <https://github.com/TuSimple/TuSimple-DUC>.

1. Introduction

Semantic segmentation aims to assign a categorical label to every pixel in an image, which plays an important role in image understanding and self-driving systems. The recent success of deep convolutional neural network (CNN) models [17] [26] [13] has enabled remarkable progress in pixel-wise semantic segmentation tasks due to rich hierarchical features and an end-to-end trainable framework [21] [31] [29] [20] [18] [3]. Most state-of-the-art semantic segmentation systems have three key components: 1) a fully-convolutional network (FCN), first introduced in [21], re-

placing the last few fully connected layers by convolutional layers to make efficient end-to-end learning and inference that can take arbitrary input size; 2) Conditional Random Fields (CRFs), to capture both local and long-range dependencies within an image to refine the prediction map; 3) dilated convolution (or Atrous convolution), which is used to increase the resolution of intermediate feature maps in order to generate more accurate predictions while maintaining the same computational cost.

Since the introduction of FCN in [21], improvements on fully-supervised semantic segmentation systems are generally focused on two perspectives: First, applying deeper FCN models. Significant gains in mean Intersection-over-Union (mIoU) scores on PASCAL VOC2012 dataset [8] were reported when the 16-layer VGG-16 model [26] was replaced by a 101-layer ResNet-101 [13] model [3]; using 152 layer ResNet-152 model yields further improvements [28]. This trend is consistent with the performance of these models on ILSVRC [23] object classification tasks, as deeper networks generally can model more complex representations and learn more discriminative features that better distinguish among categories. Second, making CRFs more powerful. This includes applying fully connected pairwise CRFs [16] as a post-processing step [3], integrating CRFs into the network by approximating its mean-field inference steps [31] [20] [18] to enable end-to-end training, and incorporating additional information into CRFs such as edges [15] and object detections [1].

We are pursuing further improvements on semantic segmentation from another perspective: the *convolutional* operations for both decoding (from intermediate feature map to output label map) and encoding (from input image to feature map) counterparts. In decoding, most state-of-the-art semantic segmentation systems simply use bilinear upsampling (before the CRF stage) to get the output label map [18] [20] [3]. Bilinear upsampling is not learnable and may lose fine details. Inspired by work in image super-resolution [25], we propose a method called *dense upsampling convo-*

lution (DUC), which is extremely easy to implement and can achieve pixel-level accuracy: instead of trying to recover the full-resolution label map at once, we learn an array of upscaling filters to upscale the downsized feature maps into the final dense feature map of the desired size. DUC naturally fits the FCN framework by enabling end-to-end training, and it increases the mIOU of pixel-level semantic segmentation on the Cityscapes dataset [5] significantly, especially on objects that are relatively small.

For the encoding part, dilated convolution recently became popular [3, 29, 28, 32], as it maintains the resolution and receptive field of the network by inserting “holes” in the convolution kernels, thus eliminating the need for down-sampling (by max-pooling or strided convolution). However, an inherent problem exists in the current dilated convolution framework, which we identify as “gridding”: as zeros are padded between two pixels in a convolutional kernel, the receptive field of this kernel only covers an area with checkerboard patterns - only locations with non-zero values are sampled, losing some neighboring information. The problem gets worse when the rate of dilation increases, generally in higher layers when the receptive field is large: the convolutional kernel is too sparse to cover any local information, since the non-zero values are too far apart. Information that contributes to a fixed pixel always comes from its predefined gridding pattern, thus losing a huge portion of information. Here we propose a simple *hybrid dilation convolution (HDC)* framework as a first attempt to address this problem: instead of using the same rate of dilation for the same spatial resolution, we use a range of dilation rates and concatenate them serially the same way as “blocks” in ResNet-101 [13]. We show that HDC helps the network to alleviate the gridding problem. Moreover, choosing proper rates can effectively increase the receptive field size and improves the accuracy for objects that are relatively big.

We design DUC and HDC to make *convolution* operations better serve the need of pixel-level semantic segmentation. The technical details are described in Section 3 below. Combined with post-processing by Conditional Random Fields (CRFs), we show that this approach achieves state-of-the-art performance on the Cityscapes pixel-level semantic labeling task, KITTI road estimation benchmark, and PASCAL VOC2012 segmentation task.

2. Related Work

Decoding of Feature Representation: In the pixel-wise semantic segmentation task, the output label map has the same size as the input image. Because of the operation of max-pooling or strided convolution in CNNs, the size of feature maps of the last few layers of the network are inevitably downsampled. Multiple approaches have been proposed to decode accurate information from the downsampled feature map to label maps. Bilinear interpolation

is commonly used [18, 20, 3], as it is fast and memory-efficient. Another popular method is called deconvolution, in which the unpooling operation, using stored pooling switches from the pooling step, recovers the information necessary for feature visualization [30]. In [21], a single deconvolutional layer is added in the decoding stage to produce the prediction result using stacked feature maps from intermediate layers. In [7], multiple deconvolutional layers are applied to generate chairs, tables, or cars from several attributes. Noh et al. [22] employ deconvolutional layers as mirrored version of convolutional layers by using stored pooled location in unpooling step. [22] show that coarse-to-fine object structures, which are crucial to recover fine-detailed information, can be reconstructed along the propagation of the deconvolutional layers. Fischer et al. [9] use a similar mirrored structure, but combine information from multiple deconvolutional layers and perform upsampling to make the final prediction.

Dilated Convolution: Dilated Convolution (or Atrous convolution) was originally developed in *algorithme à trous* for wavelet decomposition [14]. The main idea of dilated convolution is to insert “holes” (zeros) between pixels in convolutional kernels to increase image resolution, thus enabling dense feature extraction in deep CNNs. In the semantic segmentation framework, dilated convolution is also used to enlarge the field of convolutional kernels. Yu & Koltun [29] use serialized layers with increasing rates of dilation to enable context aggregation, while [3] design an “atrous spatial pyramid pooling (ASPP)” scheme to capture multi-scale objects and context information by placing multiple dilated convolution layers in parallel. More recently, dilated convolution has been applied to a broader range of tasks, such as object detection [6], optical flow [24], and audio generation [27].

3. Our Approach

3.1. Dense Upsampling Convolution (DUC)

Suppose an input image has height H , width W , and color channels C , and the goal of pixel-level semantic segmentation is to generate a label map with size $H \times W$ where each pixel is labeled with a category label. After feeding the image into a deep FCN, a feature map with dimension $h \times w \times c$ is obtained at the final layer before making predictions, where $h = H/d$, $w = W/d$, and d is the down-sampling factor. Instead of performing bilinear upsampling, which is not learnable, or using deconvolution network (as in [22]), in which zeros have to be padded in the unpooling step before the convolution operation, DUC applies convolutional operations directly on the feature maps to get the dense pixel-wise prediction map. Figure I depicts the architecture of our network with a DUC layer.

The DUC operation is all about *convolution*, which is

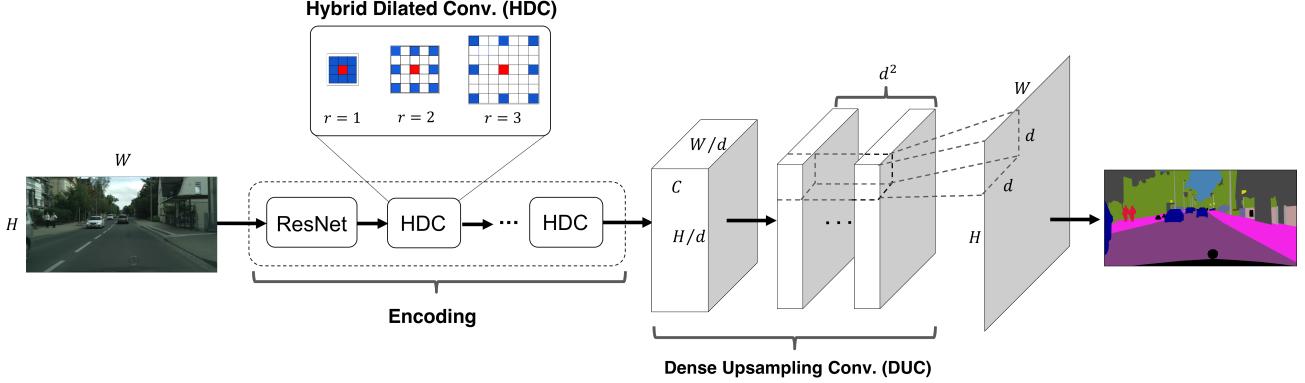


Figure 1. Illustration of the architecture of ResNet-101 network with Hybrid Dilated Convolution (HDC) and Dense Upsampling Convolution (DUC) layer. HDC is applied within ResNet blocks, and DUC is applied on top of network and is used for decoding purpose.

performed on the feature map from ResNet of dimension $h \times w \times c$ to get the output feature map of dimension $h \times w \times (d^2 \times L)$, where L is the total number of classes in the semantic segmentation task. Thus each layer of the dense convolution is learning the prediction for each pixel. The output feature map is then reshaped to $H \times W \times L$ with a softmax layer, and an elementwise *argmax* operator is applied to get the final label map. In practice, the “reshape” operation may not be necessary, as the feature map can be collapsed directly to a vector to be fed into the softmax layer. The key idea of DUC is to divide the whole label map into equal d^2 subparts which have the same height and width as the incoming feature map. This is to say, we transform the whole label map into a smaller label map with multiple channels. This transformation allows us to apply the convolution operation directly between the input feature map and the output label maps without the need of inserting extra values in deconvolutional networks (the “unpooling”operation).

Since DUC is learnable, it is capable of capturing and recovering fine-detailed information that is generally missing in the bilinear interpolation operation. For example, if a network has a downsample rate of $1/16$, and an object has a length or width less than 16 pixels (such as a pole or a person far away), then it is more than likely that bilinear up-sampling will not be able to recover this object. Meanwhile, the corresponding training labels have to be downsampled to correspond with the output dimension, which will already cause information loss for fine details. The prediction of DUC, on the other hand, is performed at the original resolution, thus enabling pixel-level decoding. In addition, the DUC operation can be naturally integrated into the FCN framework, and makes the whole encoding and decoding process end-to-end trainable.

3.2. Hybrid Dilated Convolution (HDC)

In 1-D, dilated convolution is defined as:

$$g[i] = \sum_{l=1}^L f[i + r \cdot l] h[l], \quad (1)$$

where $f[i]$ is the input signal, $g[i]$ is the output signal, $h[l]$ denotes the filter of length L , and r corresponds to the dilation rate we use to sample $f[i]$. In standard convolution, $r = 1$.

In a semantic segmentation system, 2-D dilated convolution is constructed by inserting “holes”(zeros) between each pixel in the convolutional kernel. For a convolution kernel with size $k \times k$, the size of resulting dilated filter is $k_d \times k_d$, where $k_d = k + (k - 1) \cdot (r - 1)$. Dilated convolution is used to maintain high resolution of feature maps in FCN through replacing the max-pooling operation or strided convolution layer while maintaining the receptive field (or “field of view”in [3]) of the corresponding layer. For example, if a convolution layer in ResNet-101 has a stride $s = 2$, then the stride is reset to 1 to remove downsampling, and the dilation rate r is set to 2 for all convolution kernels of subsequent layers. This process is applied iteratively through all layers that have a downsampling operation, thus the feature map in the output layer can maintain the same resolution as the input layer. In practice, however, dilated convolution is generally applied on feature maps that are already downsampled to achieve a reasonable efficiency/accuracy trade-off [3].

However, one theoretical issue exists in the above dilated convolution framework, and we call it “**gridding**”(Figure 2): For a pixel p in a dilated convolutional layer l , the information that contributes to pixel p comes from a nearby $k_d \times k_d$ region in layer $l - 1$ centered at p . Since dilated convolution introduces zeros in the convolutional kernel, the actual pixels that participate in

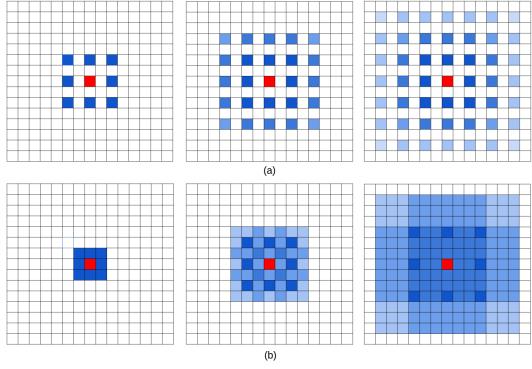


Figure 2. Illustration of the gridding problem. Left to right: the pixels (marked in blue) contributes to the calculation of the center pixel (marked in red) through three convolution layers with kernel size 3×3 . (a) all convolutional layers have a dilation rate $r = 2$. (b) subsequent convolutional layers have dilation rates of $r = 1, 2, 3$, respectively.

the computation from the $k_d \times k_d$ region are just $k \times k$, with a gap of $r - 1$ between them. If $k = 3, r = 2$, only 9 out of 25 pixels in the region are used for the computation (Figure 2 (a)). Since all layers have equal dilation rates r , then for pixel p in the top dilated convolution layer l_{top} , the maximum possible number of locations that contribute to the calculation of the value of p is $(w' \times h')/r^2$ where w', h' are the width and height of the bottom dilated convolution layer, respectively. As a result, pixel p can only view information in a checkerboard fashion, and lose a large portion (at least 75% when $r = 2$) of information. When r becomes large in higher layers due to additional downsampling operations, the sample from the input can be very sparse, which may not be good for learning because 1) local information is completely missing; 2) the information can be irrelevant across large distances. Another outcome of the gridding effect is that pixels in nearby $r \times r$ regions at layer l receive information from completely different set of “grids” which may impair the consistency of local information.

Here we propose a simple solution- *hybrid dilated convolution (HDC)*, to address this theoretical issue. Suppose we have N convolutional layers with kernel size $K \times K$ that have dilation rates of $[r_1, \dots, r_i, \dots, r_n]$, the goal of HDC is to let the final size of the RF of a series of convolutional operations fully covers a square region without any holes or missing edges. We define the “maximum distance between two nonzero values” as

$$M_i = \max[M_{i+1} - 2r_i, M_{i+1} - 2(M_{i+1} - r_i), r_i], \quad (2)$$

with $M_n = r_n$. The design goal is to let $M_2 \leq K$. For example, for kernel size $K = 3$, an $r = [1, 2, 5]$ pattern works as $M_2 = 2$; however, an $r = [1, 2, 9]$ pattern does not work as $M_2 = 5$. Practically, instead of using the same dilation rate for all layers after the downsampling occurs, we

use a different dilation rate for each layer. In our network, the assignment of dilation rate follows a sawtooth wave-like heuristic: a number of layers are grouped together to form the “rising edge” of the wave that has an increasing dilation rate, and the next group repeats the same pattern. For example, for all layers that have dilation rate $r = 2$, we form 3 succeeding layers as a group, and change their dilation rates to be 1, 2, and 3, respectively. By doing this, the top layer can access information from a broader range of pixels, in the same region as the original configuration (Figure 2 (b)). This process is repeated through all layers, thus making the receptive field unchanged at the top layer.

Another benefit of HDC is that it can use arbitrary dilation rates through the process, thus naturally enlarging the receptive fields of the network without adding extra modules [29], which is important for recognizing objects that are relatively big. One important thing to note, however, is that the dilation rate within a group should not have a common factor relationship (like 2,4,8, etc.), otherwise the gridding problem will still hold for the top layer. This is a key difference between our HDC approach and the atrous spatial pyramid pooling (ASPP) module in [3], or the context aggregation module in [29], where dilation factors that have common factor relationships are used. In addition, HDC is naturally integrated with the original layers of the network, without any need to add extra modules as in [29, 3].

4. Experiments and Results

We report our experiments and results on three challenging semantic segmentation datasets: Cityscapes [5], KITTI dataset [10] for road estimation, and PASCAL VOC2012 [8]. We use ResNet-101 or ResNet-152 networks that have been pretrained on the ImageNet dataset as a starting point for all of our models. The output layer contains the number of semantic categories to be classified depending on the dataset (including background, if applicable). We use the cross-entropy error at each pixel over the categories. This is then summed over all pixel locations of the output map, and we optimize this objective function using standard Stochastic Gradient Descent (SGD). We use MXNet [4] to train and evaluate all of our models on NVIDIA TITAN X GPUs.

4.1. Cityscapes Dataset

The Cityscapes Dataset is a large dataset that focuses on semantic understanding of urban street scenes. The dataset contains 5000 images with fine annotations across 50 cities, different seasons, varying scene layout and background. The dataset is annotated with 30 categories, of which 19 categories are included for training and evaluation (others are ignored). The training, validation, and test set contains 2975, 500, and 1525 fine images, respectively. An additional 20000 images with coarse (polygonal) annotations are also provided, but are only used for training.

4.1.1 Baseline Model

We use the DeepLab-V2 [3] ResNet-101 framework to train our baseline model. Specifically, the network has a down-sampling rate of 8, and dilated convolution with rate of 2 and 4 are applied to *res4b* and *res5b* blocks, respectively. An ASPP module with dilation rate of 6, 12, 18, and 24 is added on top of the network to extract multiscale context information. The prediction maps and training labels are downsampled by a factor of 8 compared to the size of original images, and bilinear upsampling is used to get the final prediction. Since the image size in the Cityscapes dataset is 1024×2048 , which is too big to fit in the GPU memory, we partition each image into twelve 800×800 patches with partial overlapping, thus augmenting the training set to have 35700 images. This data augmentation strategy is to make sure all regions in an image can be visited. This is an improvement over random cropping, in which nearby regions may be visited repeatedly.

We train the network using mini-batch SGD with patch size 544×544 (randomly cropped from the 800×800 patch) and batch size 12, using multiple GPUs. The initial learning rate is set to 2.5×10^{-4} , and a “*poly*” learning rate (as in [3]) with *power* = 0.9 is applied. Weight decay is set to 5×10^{-4} , and momentum is 0.9. The network is trained for 20 epochs and achieves mIoU of 72.3% on the validation set.

4.1.2 Dense Upsampling Convolution (DUC)

We examine the effect of DUC on the baseline network. In DUC the only thing we change is the shape of the top convolutional layer. For example, if the dimension of the top convolutional layer is $68 \times 68 \times 19$ in the baseline model (19 is the number of classes), then the dimension of the same layer for a network with DUC will be $68 \times 68 \times (r^2 \times 19)$ where r is the total downsampling rate of the network ($r = 8$ in this case). The prediction map is then reshaped to size $544 \times 544 \times 19$. DUC will introduce extra parameters compared to the baseline model, but only at the top convolutional layer. We train the ResNet-DUC network the same way as the baseline model for 20 epochs, and achieve a mean IOU of 74.3% on the validation set, a 2% increase compared to the baseline model. Visualization of the result of ResNet-DUC and comparison with the baseline model is shown in Figure 3

From Figure 3, we can clearly see that DUC is very helpful for identifying small objects, such as poles, traffic lights, and traffic signs. Consistent with our intuition, pixel-level dense upsampling can recover detailed information that is generally missed by bilinear interpolation.

Ablation Studies We examine the effect of different settings of the network on the performance. Specifically, we examine: 1) the downsampling rate of the network, which

controls the resolution of the intermediate feature map; 2) whether to apply the ASPP module, and the number of parallel paths in the module; 3) whether to perform 12-fold data augmentation; and 4) cell size, which determines the size of neighborhood region ($cell \times cell$) that one predicted pixel projects to. Pixel-level DUC should use $cell = 1$; however, since the ground-truth label generally cannot reach pixel-level precision, we also try $cell = 2$ in the experiments. From Table 1 we can see that making the downsampling rate smaller decreases the accuracy. Also it significantly raises the computational cost due to the increasing resolution of the feature maps. ASPP generally helps to improve the performance, and increasing ASPP channels from 4 to 6 (dilation rate 6 to 36 with interval 6) yields a 0.2% boost. Data augmentation helps to achieve another 1.5% improvement. Using $cell = 2$ yields slightly better performance when compared with $cell = 1$, and it helps to reduce computational cost by decreasing the channels of the last convolutional layer by a factor of 4.

Network	DS	ASPP	Augmentation	Cell	mIoU
<i>Baseline</i>	8	4	<i>yes</i>	<i>n/a</i>	72.3
<i>Baseline</i>	4	4	<i>yes</i>	<i>n/a</i>	70.9
<i>DUC</i>	8	<i>no</i>	<i>no</i>	1	71.9
<i>DUC</i>	8	4	<i>no</i>	1	72.8
<i>DUC</i>	8	4	<i>yes</i>	1	74.3
<i>DUC</i>	4	4	<i>yes</i>	1	73.7
<i>DUC</i>	8	6	<i>yes</i>	1	74.5
<i>DUC</i>	8	6	<i>yes</i>	2	74.7

Table 1. Ablation studies for applying ResNet-101 on the Cityscapes dataset. **DS**: Downsampling rate of the network. **Cell**: neighborhood region that one predicted pixel represents.

Bigger Patch Size Since setting $cell = 2$ reduces GPU memory cost for network training, we explore the effect of patch size on the performance. Our assumption is that, since the original images are all 1024×2048 , the network should be trained using patches as big as possible in order to aggregate both local detail and global context information that may help learning. As such, we make the patch size to be 880×880 , and set the batch size to be 1 on each of the 4 GPUs used in training. Since the patch size exceeds the maximum dimension (800×800) in the previous 12-fold data augmentation framework, we adopt a new 7-fold data augmentation strategy: seven center locations with $x = 512$, $y = \{256, 512, \dots, 1792\}$ are set in the original image; for each center location, a 880×880 patch is obtained by randomly setting its center within a 160×160 rectangle area centered at each center. This strategy makes sure that we can sample all areas in the image, including edges. Training with a bigger patch size boosts the performance to **75.7%**, a 1% improvement over the previous best

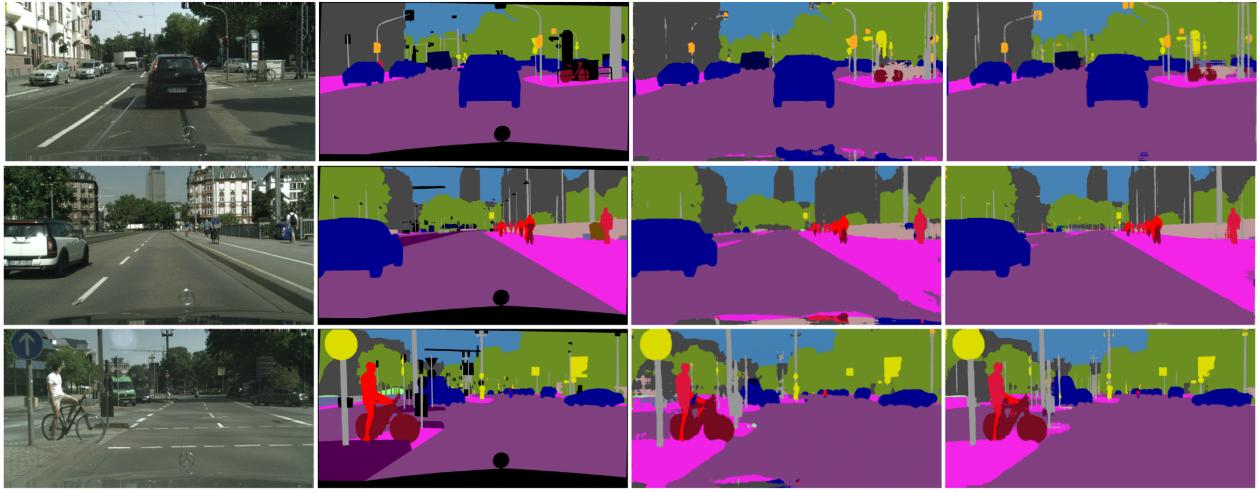


Figure 3. Effect of Dense Upsampling Convolution (DUC) on the Cityscapes validation set. From left to right: input image, ground truth (areas with black color are ignored in evaluation), baseline model, and our ResNet-DUC model.

result.

Compared with Deconvolution We compare our DUC model with deconvolution, which also involves learning for upsampling. Particularly, we compare with 1) direct deconvolution from the prediction map (dowsampled by 8) to the original resolution; 2) deconvolution with an upsampling factor of 2 first, followed by an upsampling factor of 4. We design the deconv network to have approximately the same number of parameters as DUC. We use the ResNet-DUC bigger patch model to train the networks. The above two models achieve mIoU of 75.1% and 75.0%, respectively, lower than the ResNet-DUC model (75.7% mIoU).

Conditional Random Fields (CRFs) Fully-connected CRFs [16] are widely used for improving semantic segmentation quality as a post-processing step of an FCN [3]. We follow the formation of CRFs as shown in [3]. We perform a grid search on parameters on the validation set, and use $\sigma_\alpha = 15$, $\sigma_\beta = 3$, $\sigma_\gamma = 1$, $w_1 = 3$, and $w_2 = 3$ for all of our models. Applying CRFs to our best ResNet-DUC model yields an mIoU of 76.7%, a 1% improvement over the model does not use CRFs.

4.1.3 Hybrid Dilated Convolution (HDC)

We use the best 101 layer ResNet-DUC model as a starting point of applying HDC. Specifically, we experiment with several variants of the HDC module:

1. No dilation: For all ResNet blocks containing dilation, we make their dilation rate $r = 1$ (no dilation).
2. Dilation-conv: For all blocks contain dilation, we group every 2 blocks together and make $r = 2$ for the first block, and $r = 1$ for the second block.

3. Dilation-RF: For the *res4b* module that contains 23 blocks with dilation rate $r = 2$, we group every 3 blocks together and change their dilation rates to be 1, 2, and 3, respectively. For the last two blocks, we keep $r = 2$. For the *res5b* module which contains 3 blocks with dilation rate $r = 4$, we change them to 3, 4, and 5, respectively.
4. Dilation-bigger: For *res4b* module, we group every 4 blocks together and change their dilation rates to be 1, 2, 5, and 9, respectively. The rates for the last 3 blocks are 1, 2, and 5. For *res5b* module, we set the dilation rates to be 5, 9, and 17.

The result is summarized in Table 2. We can see that increasing receptive field size generally yields higher accuracy. Figure 5 illustrates the effectiveness of the ResNet-DUC-HDC model in eliminating the gridding effect. A visualization result is shown in Figure 4. We can see our best ResNet-DUC-HDC model performs particularly well on objects that are relatively big.

Network	RF increased	mIoU (without CRF)
No dilation	54	72.9
Dilation-conv	88	75.0
Dilation-RF	116	75.4
Dilation-bigger	256	76.4

Table 2. Result of different variations of the HDC module. “RF increased” is the total size of receptive field increase along a single dimension compared to the layer before the dilation operation.

Deeper Networks We have also tried replacing our ResNet-101 based model with the ResNet-152 network, which is deeper and achieves better performance on the ILSVRC image classification task than ResNet-101 [13].

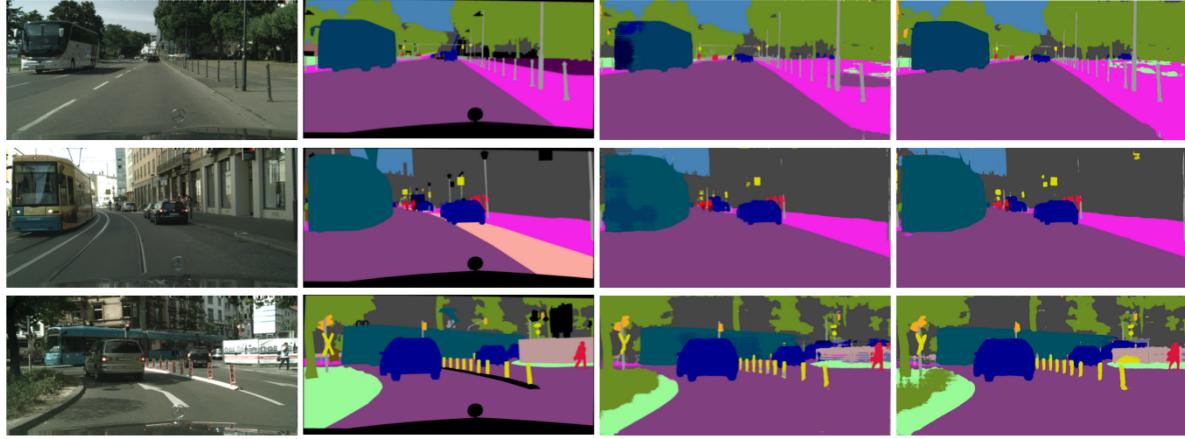


Figure 4. Effect of Hybrid Dilated Convolution (HDC) on the Cityscapes validation set. From left to right: input image, ground truth, result of the ResNet-DUC model, result of the ResNet-DUC-HDC model (Dilation-bigger).



Figure 5. Effectiveness of HDC in eliminating the gridding effect. First row: ground truth patch. Second row: prediction of the ResNet-DUC model. A strong gridding effect is observed. Third row: prediction of the ResNet-DUC-HDC (Dilation-RF) model.

Due to the network difference, we first train the ResNet-152 network to learn the parameters in all batch normalization (BN) layers for 10 epochs, and continue fine-tuning the network by fixing these BN parameters for another 20 epochs. The results are summarized in Table 3. We can see that using the deeper ResNet-152 model generally yields better performance than the ResNet-101 model.

Network	Method	data	mIoU
ResNet-101	<i>Deconv</i>	fine	75.1
ResNet-101	<i>DUC+HDC</i>	fine	76.4
ResNet-101	<i>DUC+HDC</i>	fine+coarse	76.2
ResNet-152	<i>Deconv</i>	fine	76.4
ResNet-152	<i>DUC+HDC</i>	fine	76.7
ResNet-152	<i>DUC+HDC</i>	fine+coarse	77.1

Table 3. Effect of depth of the network and upsampling method for Cityscapes validation set (without CRF).

4.1.4 Test Set Results

Our results on the Cityscapes test set are summarized in Table 4. There are separate entries for models trained using fine-labels only, and using a combination of fine and coarse labels. Our ResNet-DUC-HDC model achieves 77.6% mIoU using fine data only. Adding coarse data help us achieve **78.5%** mIoU.

In addition, inspired by the design of the VGG network [26], in that a single 5×5 convolutional layer can be decomposed into two adjacent 3×3 convolutional layers to increase the expressiveness of the network while maintaining the receptive field size, we replaced the 7×7 convolutional layer in the original ResNet-101 network by three 3×3 convolutional layers. By retraining the updated network, we achieve a mIoU of **80.1%** on the test set using a single model without CRF post-processing. Our result achieves the state-of-the-art performance on the Cityscapes dataset at the time of submission. Compared with the strong baseline of Chen et al. [3], we improve the mIoU by a significant margin (9.7%), which demonstrates the effectiveness of our approach.

Method	mIoU
<i>fine</i>	
FCN 8s [21]	65.3%
Dilation10 [29]	67.1%
DeepLabv2-CRF [3]	70.4%
Adelaide_context [18]	71.6%
ResNet-DUC-HDC (ours)	77.6%
<i>coarse</i>	
LRR-4x [11]	71.8%
ResNet-DUC-HDC-Coarse	78.5%
ResNet-DUC-HDC-Coarse (better network)	80.1%

Table 4. Performance on Cityscapes test set.

4.2. KITTI Road Segmentation

Dataset The KITTI road segmentation task contains images of three various categories of road scenes, including 289 training images and 290 test images. The goal is to decide if each pixel in images is road or not. It is challenging to use neural network based methods due to the limited number of training images. In order to avoid overfitting, we crop patches of 320×320 pixels with a stride of 100 pixels from the training images, and use the ResNet-101-DUC model pretrained from ImageNet during training. Other training settings are the same as Cityscapes experiment. We did not apply CRFs for post-processing.

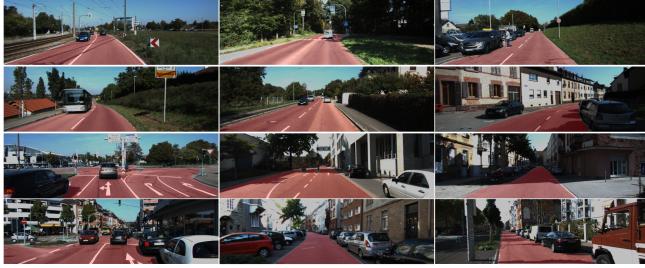


Figure 6. Examples of visualization on Kitti road segmentation test set. The road is marked in red.

Results We achieve the state-of-the-art results at the time of submission without using any additional information of stereo, laser points and GPS. Specifically, our model attains the highest maximum F1-measure in the sub-categories of urban unmarked (UU_ROAD), urban multiple marked (UMM_ROAD) and the overall category URBAN_ROAD of all sub-categories, the highest average precision across all three sub-categories and the overall category by the time of submission of this paper. Examples of visualization results are shown in Figure 6. The detailed results are displayed in Table 5¹.

	MaxF	AP
UM_ROAD	95.64%	93.50%
UMM_ROAD	97.62%	95.53%
UM_ROAD	95.17%	92.73%
URBAN_ROAD	96.41%	93.88%

Table 5. Performance on different road scenes in KITTI test set. MaxF: Maximum F1-measure, AP: Average precision.

4.3. PASCAL VOC2012 dataset

Dataset The PASCAL VOC2012 segmentation benchmark contains 1464 training images, 1449 validation images, and 1456 test images. Using the extra annotations provided by [12], the training set is augmented to have 10582

¹For thorough comparison with other methods, please check http://www.cvlabs.net/datasets/kitti/eval_road.php.

images. The dataset has 20 foreground object categories and 1 background class with pixel-level annotation.

Results We first pretrain our 152 layer ResNet-DUC model using a combination of augmented VOC2012 training set and MS-COCO dataset [19], and then finetune the pretrained network using augmented VOC2012 trainval set. We use patch size 512×512 (zero-padded) throughout training. All other training strategies are the same as Cityscapes experiment. We achieve mIoU of **83.1%** on the test set using a single model without any model ensemble or multiscale testing, which is the best-performing method at the time of submission². The detailed results are displayed in Table 6, and the visualizations are shown in Figure 7.

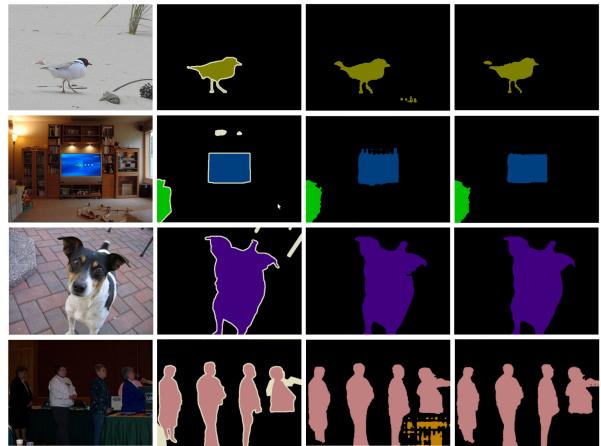


Figure 7. Examples of visualization on the PASCAL VOC2012 segmentation validation set. Left to right: input image, ground truth, our result before CRF, and after CRF.

Method	mIoU
DeepLabv2-CRF[3]	79.7%
CentraleSupelec Deep G-CRF[2]	80.2%
ResNet-DUC (ours)	83.1%

Table 6. Performance on the Pascal VOC2012 test set.

5. Conclusion

We propose simple yet effective convolutional operations for improving semantic segmentation systems. We designed a new dense upsampling convolution (DUC) operation to enable pixel-level prediction on feature maps, and hybrid dilated convolution (HDC) to solve the gridding problem, effectively enlarging the receptive fields of the network. Experimental results demonstrate the effectiveness of our framework on various semantic segmentation tasks.

²Result link: <http://host.robots.ox.ac.uk:8080/anonymous/LQ2ACW.html>

6. Acknowledgments

We thank the members of TuSimple and Gary’s Unbelievable Research Unit (GURU) for comments on this work. GWC was supported in part by Guangzhou Science and Technology Planning Project (201704030051) and NSF cooperative agreement SMA 1041755 to the Temporal Dynamics of Learning Center, an NSF Science of Learning Center.

References

- [1] A. Arnab, S. Jayasumana, S. Zheng, and P. Torr. Higher order potentials in end-to-end trainable conditional random fields. *arXiv preprint arXiv:1511.08119*, 2015.
- [2] S. Chandra and I. Kokkinos. Fast, exact and multi-scale inference for semantic image segmentation with deep gaussian crfs. *arXiv preprint arXiv:1603.08358*, 2016.
- [3] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv preprint arXiv:1606.00915*, 2016.
- [4] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015.
- [5] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [6] J. Dai, Y. Li, K. He, and J. Sun. R-fcn: Object detection via region-based fully convolutional networks. *arXiv preprint arXiv:1605.06409*, 2016.
- [7] A. Dosovitskiy, J. Tobias Springenberg, and T. Brox. Learning to generate chairs with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1538–1546, 2015.
- [8] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.
- [9] P. Fischer, A. Dosovitskiy, E. Ilg, P. Häusser, C. Hazırbaş, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox. Flownet: Learning optical flow with convolutional networks. *arXiv preprint arXiv:1504.06852*, 2015.
- [10] J. Fritsch, T. Kuehnl, and A. Geiger. A new performance measure and evaluation benchmark for road detection algorithms. In *International Conference on Intelligent Transportation Systems (ITSC)*, 2013.
- [11] G. Ghiasi and C. Fowlkes. Laplacian reconstruction and refinement for semantic segmentation. *arXiv preprint arXiv:1605.02264*, 2016.
- [12] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik. Semantic contours from inverse detectors. In *2011 International Conference on Computer Vision*, pages 991–998. IEEE, 2011.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- [14] M. Holschneider, R. Kronland-Martinet, J. Morlet, and P. Tchamitchian. A real-time algorithm for signal analysis with the help of the wavelet transform. In *Wavelets*, pages 286–297. Springer, 1990.
- [15] I. Kokkinos. Pushing the boundaries of boundary detection using deep learning. *arXiv preprint arXiv:1511.07386*, 2015.
- [16] P. Krähenbühl and V. Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *Advances in neural information processing systems*, 2011.
- [17] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [18] G. Lin, C. Shen, I. Reid, et al. Efficient piecewise training of deep structured models for semantic segmentation. *arXiv preprint arXiv:1504.01013*, 2015.
- [19] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014.
- [20] Z. Liu, X. Li, P. Luo, C.-C. Loy, and X. Tang. Semantic image segmentation via deep parsing network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1377–1385, 2015.
- [21] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [22] H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1520–1528, 2015.
- [23] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [24] L. Sevilla-Lara, D. Sun, V. Jampani, and M. J. Black. Optical flow with semantic segmentation and localized layers. *arXiv preprint arXiv:1603.03911*, 2016.
- [25] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1874–1883, 2016.
- [26] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [27] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.

- [28] Z. Wu, C. Shen, and A. v. d. Hengel. High-performance semantic segmentation using very deep fully convolutional networks. *arXiv preprint arXiv:1604.04339*, 2016.
- [29] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [30] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833. Springer, 2014.
- [31] 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. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1529–1537, 2015.
- [32] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Semantic understanding of scenes through the ade20k dataset. *arXiv preprint arXiv:1608.05442*, 2016.