

Decoders Matter for Semantic Segmentation: Data-Dependent Decoding Enables Flexible Feature Aggregation

Zhi Tian¹ Tong He¹ Chunhua Shen^{1*} Youliang Yan²

¹The University of Adelaide, Australia

²Noah’s Ark Lab, Huawei Technologies

Abstract

Recent semantic segmentation methods exploit encoder-decoder architectures to produce the desired pixel-wise segmentation prediction. The last layer of the decoders is typically a bilinear upsampling procedure to recover the final pixel-wise prediction. We empirically show that this over-simple and data-independent bilinear upsampling may lead to sub-optimal results.

In this work, we propose a data-dependent upsampling (DUpsampling) to replace bilinear, which takes advantages of the redundancy in the label space of semantic segmentation and is able to recover the pixel-wise prediction from low-resolution outputs of CNNs. The main advantage of the new upsampling layer lies in that with a relatively lower-resolution feature map such as $\frac{1}{16}$ or $\frac{1}{32}$ of the input size, we can achieve even better segmentation accuracy, significantly reducing computation complexity. This is made possible by 1) the new upsampling layer’s much improved reconstruction capability; and more importantly 2) the DUpsampling based decoder’s flexibility in leveraging almost arbitrary combinations of the CNN encoders’ features. Experiments demonstrate that our proposed decoder outperforms the state-of-the-art decoder, with only $\sim 20\%$ of computation. Finally, without any post-processing, the framework equipped with our proposed decoder achieves new state-of-the-art performance on two datasets: 88.1% mIOU on PASCAL VOC with 30% computation of the previously best model; and 52.5% mIOU on PASCAL Context.

1. Introduction

Fully convolutional networks (FCNs) [21] have achieved tremendous success in dense pixel prediction applications such as semantic segmentation, for which the algorithm is asked to predict a variable for each pixel of an input image and is a fundamental problem in computer vision. The great achievement of FCNs results from powerful features

*Appearing in IEEE Conf. Computer Vision and Pattern Recognition, 2019. First two authors equally contributed to this work. C. Shen is the corresponding author: chunhua.shen@adelaide.edu.au

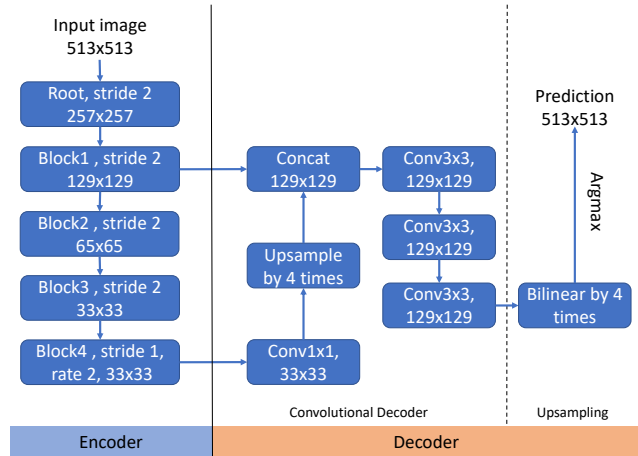


Figure 1: An example of the encoder-decoder architecture used by DeepLabv3+. Its decoder fuses low-level features of downsample ratio = 4 and upsamples high-level features before merging them. Finally, bilinear upsampling is applied to restore the full-resolution prediction. “rate” denotes the atrous rate in atrous convolution.

extracted by CNNs. Importantly, the sharing convolutional computation mechanism makes training and inference computationally very efficient.

In the original FCNs, several stages of strided convolutions and/or spatial pooling reduce the final image prediction typically by a factor of 32, thus losing fine image structure information and leading to inaccurate predictions, especially at the object boundaries. DeepLab [3] applies atrous (a.k.a dilation) convolutions, achieving large receptive fields while maintaining a higher-resolution feature map. Alternatively the encoder-decoder architecture is often used to address this problem. The encoder-decoder architecture views the backbone CNN as an encoder, responsible for encoding a raw input image into lower-resolution feature maps (e.g., $\frac{1}{r}$ of the input image size with $r = 8, 16,$ or 32). Afterwards, a decoder is used to recover the pixel-wise prediction from the lower-resolution feature maps. In previous works [5, 16], a decoder consists of a few convolutional layers and a bilinear upsampling. The light-weight

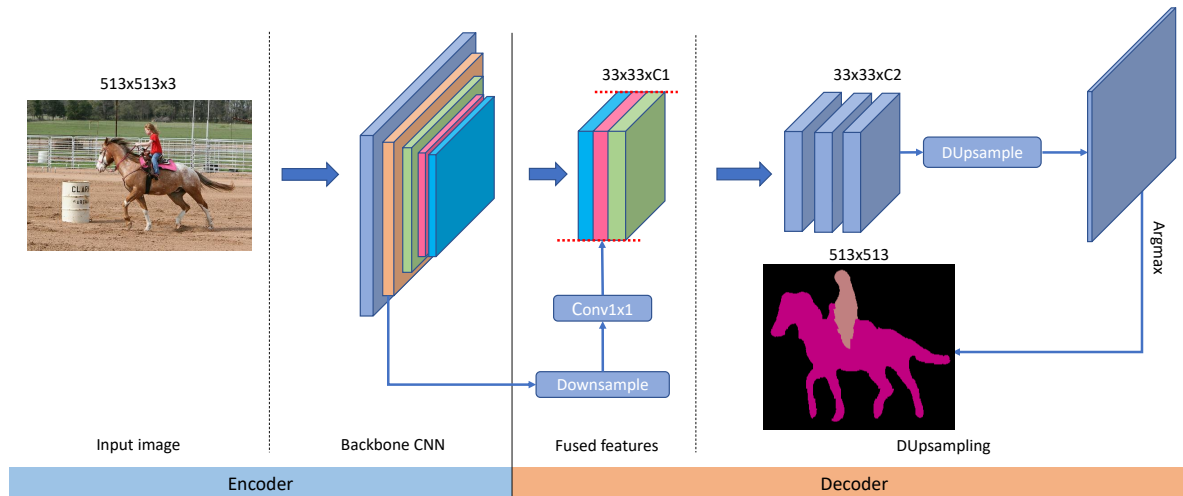


Figure 2: The framework with our proposed decoder. The major differences from the previous framework shown in Fig. 1 are 1) all fused features are downsampled to the lowest features resolution before merging. 2) The incapable bilinear is replaced with our proposed DUpsampling to recover the full-resolution prediction.

convolutional decoder yields high-resolution feature maps and bilinear upsampling is finally applied to the resulting feature maps to obtain the desired pixel-wise prediction. The decoder commonly fuses low-level features to capture the fine-grained information lost by convolution and pooling operations in CNNs. A standard DeepLabv3+ encoder-decoder architecture is illustrated in Fig. 1.

A drawback of the oversimple bilinear upsampling is its limited capability in recovering the pixel-wise prediction accurately. Bilinear upsampling does not take into account the correlation among the prediction of each pixel since it is data independent. As a consequence, the convolutional decoder is required to produce relatively higher-resolution feature maps in order to obtain good final prediction (e.g., $\frac{1}{4}$ or $\frac{1}{8}$ of the input size). This requirement causes two issues for semantic segmentation.

1) The overall strides of the encode must be reduced very aggressively by using multiple atrous convolutions [3, 33]. The price is much heavier computation complexity and memory footprint, hampering the training on large data and deployment for real-time applications.

For example, in order to achieve state-of-the-art performance, the recent DeepLabv3+ [5] reduces the overall strides of its encoder by four times (from 32 to 8). Thus inference of DeepLabv3+ is very slow.

2) The decoder is often needed to fuse features at very low levels. For example, DeepLabv3+ fuses features of downsample ratio = 4^1 in block1 as shown in Fig. 1. It is because that the fineness of the final prediction is actually dominated by the resolution of the fused low-level fea-

tures due to the inability of bilinear. As a result, in order to produce high-resolution prediction, the decoder has to fuse the high-resolution features at a low level. This constraint narrows down the design space of the feature aggregation and therefore is likely to cause a *suboptimal* combination of features to be aggregated in the decoder. In experiments, *we show that a better feature aggregation strategy can be found if the feature aggregation can be designed without the constraint imposed by the resolution of feature maps.*

In order to tackle the aforementioned issues caused by bilinear, here we propose a new data-dependent upsampling method, termed DUPSampling, to recover the pixel-wise prediction from the final outputs of the CNNs, replacing bilinear upsampling used extensively in previous works. Our proposed DUPSampling takes advantages of the redundancy in the segmentation label space and proves to be capable of accurately recovering the pixel-wise prediction from relatively coarse CNNs outputs, alleviating the need for precise responses from the convolutional decoder.

As a result, the encoder is not required to overly reduce its overall strides, dramatically reducing the computation time and memory footprint of the whole segmentation framework. Meanwhile, due to the effectiveness of DUPSampling, it allows the decoder to downsample the fused features to the lowest resolution of feature maps before merging them. This downsampling not only reduces the amount of computation of the decoder, but much more importantly it decouples the resolution of fused features and that of the final prediction. This decoupling allows the decoder to make use of arbitrary feature aggregation and thus a better feature aggregation can be leveraged so as to boost the segmentation performance as much as possible.

Finally, DUPSampling can be seamlessly incorporated

¹downsample ratio denotes the ratio of the resolution of the feature maps to that of the input image.

into the network with a standard 1×1 convolution and thus needs no ad-hoc coding. Our overall framework is shown in Fig. 2.

We summarize our main contributions as follows.

- We propose a simple yet effective Data-dependent Upsampling (DUpsampling) method to recover the pixel-wise segmentation prediction from the coarse outputs of the convolutional decoder, replacing the incapable bilinear used extensively in previous methods.
- Taking advantages of our proposed DUpsampling, we can avoid overly reducing the overall strides of the encoder, significantly reducing the computation time and memory footprint of the semantic segmentation method by a factor of 3 or so.
- DUpsampling also allows the decoder to downsample the fused features to the lowest resolution of feature maps before merging them. The downsampling not only reduces the amount of computation of the decoder dramatically but also enlarges the design space of feature aggregation, allowing a better feature aggregation to be exploited in the decoder.
- Together with the above contributions, we propose a new decoder scheme, which compares favourably with state-of-the-art decoders while using $\sim 20\%$ amount of computation. With the proposed decoder, the framework illustrated in Fig. 2 achieves new state-of-the-art performance: mIOU of $88.1\%^2$ on PASCAL VOC [8] with only 30% computation of the previous best framework of DeepLabv3+ [5]. We also set a new mIOU record of 52.5% on the PASCAL Context dataset [23].

2. Related Work

Efforts have been devoted to improve pixel-wise predictions with FCNs. They can be roughly divided into two groups: atrous convolution [3, 33] and encoder-decoder architectures [16, 5, 21, 2, 24].

Atrous convolution. A straightforward approach is to reduce the overall strides of backbone CNNs by dropping some strided convolutions or pooling layers. However, simply reducing these strides would diminish the receptive field of convolution networks substantially, which proves to be crucial to semantic segmentation [3, 25, 19]. Atrous convolutions [4, 3, 5, 33] can be used to keep the receptive field unchanged, meanwhile not downsampling the feature map resolution too much. The major drawback of atrous convolutions is much heavier computation complexity and larger memory requirement as the size of those atrous convolutional kernels, as well as the resulted feature maps, become much larger [11, 6].

²The results on PASCAL VOC *test* set can be found at <http://host.robots.ox.ac.uk:8080/anonymous/UYT221.html>

Encoder-decoder architectures. Encoder-decoder architectures are proposed to overcome the drawback of atrous convolutions and are widely used for semantic segmentation. DeconvNet [24] uses stacked deconvolutional layers to recover the full-resolution prediction gradually. The method has the potential to produce high-resolution prediction but is difficult to train due to many parameters introduced by the decoder. SegNet [2] shares a similar idea with DeconvNet but uses indices in pooling layers to guide the recovery process, resulting in better performance. RefineNet [16] further fuse low-level features to improve the performance. Recently, DeepLabv3+ [5] takes advantages of both encoder-decoder architectures and atrous convolution, achieving best reported performance on a few datasets to date. Although efforts have been spent on designing a better decoder, so far almost none of them can bypass the restriction on the resolutions of the fused features and exploit better feature aggregation.

3. Our Approach

In this section, we firstly reformulate semantic segmentation with our proposed DUpsampling and then present the adaptive-temperature softmax function which makes the training with DUpsampling much easier. Finally, we show how the framework can be largely improved with the fusion of downsampled low-level features.

3.1. Beyond Bilinear: Data-dependent Upsampling

In this section, we firstly consider the simplest decoder, which is only composed of upsampling. Let $\mathbf{F} \in \mathbb{R}^{\tilde{H} \times \tilde{W} \times \tilde{C}}$ denote the final outputs of the encoder CNNs and $\mathbf{Y} \in \{0, 1, 2, \dots, C\}^{H \times W}$ be the ground truth label map, where C and \tilde{C} denotes the number of classes of segmentation and the number of channels of the final outputs, respectively. \mathbf{Y} is commonly encoded with one-hot encoding, i.e., $\mathbf{Y} \in \{0, 1\}^{H \times W \times C}$. Note that \mathbf{F} is typically of a factor of 16 or 32 in spatial size of the ground-truth \mathbf{Y} . In other words, $\frac{\tilde{H}}{H} = \frac{\tilde{W}}{W} = \frac{1}{16}$ or $\frac{1}{32}$. Since semantic segmentation requires per-pixel prediction, \mathbf{F} needs to be upsampled to the spatial size of \mathbf{Y} before computing the training loss.

Typically in semantic segmentation [4, 5, 21, 34, 12], the training loss function is formulated as:

$$\mathcal{L}(\mathbf{F}, \mathbf{Y}) = \text{Loss}(\text{softmax}(\text{bilinear}(\mathbf{F})), \mathbf{Y}). \quad (1)$$

Here Loss is often the cross-entropy loss, and bilinear is used to upsample \mathbf{F} to the spatial size of \mathbf{Y} . We argue that bilinear unsampling may not be the optimal choice here. As we show in the experiments (Sec. 4.1.1), bilinear is oversimple and has an inferior upper bound in terms of reconstructing (best possible reconstruction quality). In order to compensate the loss caused by bilinear, the employed deep network is consequently required to output higher-resolution

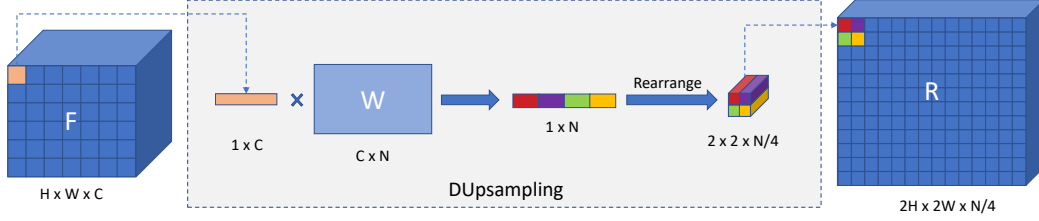


Figure 3: The proposed DUpsampling. In the figure, DUpsampling is used to upsample the CNNs outputs \mathbf{F} by twice. \mathbf{R} denotes the resulting maps. \mathbf{W} , computed with the method described in Sec. 3.1, is the inverse projection matrix of DUpsampling. In practice, the upsampling ratio is typically 16 or 32.

feature maps, which are input to the bilinear operator. As mentioned above, the solution is to apply atrous convolutions, with the price of high computation complexity. For example, *reducing the overall strides from 16 to 8 incurs more than 3 times computation.*

An important observation is that the semantic segmentation label \mathbf{Y} of an image is not i.i.d. and there contains structure information so that \mathbf{Y} can be compressed considerably, with almost no loss. Therefore, unlike previous methods, which upsample \mathbf{F} to the spatial size of \mathbf{Y} , we instead attempt to compress \mathbf{Y} into $\tilde{\mathbf{Y}} \in \mathbb{R}^{\tilde{H} \times \tilde{W} \times \tilde{C}}$ and then compute the training loss between \mathbf{F} and $\tilde{\mathbf{Y}}$. Note that \mathbf{F} and $\tilde{\mathbf{Y}}$ are of the same size.

In order to compress \mathbf{Y} into $\tilde{\mathbf{Y}}$, we seek a transform under some metric to minimize the reconstruction error between \mathbf{Y} and $\tilde{\mathbf{Y}}$. Specifically, let r indicate the ratio of H to \tilde{H} , which is usually 16 or 32. Next, \mathbf{Y} is divided into an $\frac{H}{r} \times \frac{W}{r}$ grid of sub-windows of size $r \times r$ (if H or W is not dividable by r , a padding is applied). For each sub-window $\mathbf{S} \in \{0, 1\}^{r \times r \times C}$, we reshape \mathbf{S} into a vector $\mathbf{v} \in \{0, 1\}^N$, with $N = r \times r \times C$. Finally, we compress the vector \mathbf{v} to a lower-dimensional vector $\mathbf{x} \in \mathbb{R}^{\tilde{C}}$ and then vertically and horizontally stack all \mathbf{x} 's to form $\tilde{\mathbf{Y}}$.

Although a variety of ways can be used to achieve the compression, we find that simply using linear projecting, i.e., multiplying \mathbf{v} by a matrix $\mathbf{P} \in \mathbb{R}^{\tilde{C} \times N}$ works well in this case. Formally, we have,

$$\mathbf{x} = \mathbf{P}\mathbf{v}; \quad \tilde{\mathbf{v}} = \mathbf{W}\mathbf{x}, \quad (2)$$

where $\mathbf{P} \in \mathbb{R}^{\tilde{C} \times N}$ is used to compress \mathbf{v} into \mathbf{x} . $\mathbf{W} \in \mathbb{R}^{N \times \tilde{C}}$ is the inverse projection matrix (a.k.a. reconstruction matrix) and used to reconstruct \mathbf{x} back to \mathbf{v} . $\tilde{\mathbf{v}}$ is the reconstructed \mathbf{v} . We have omitted the offset term here. In practice prior to the compression, \mathbf{v} is centered by subtracting its mean over the training set.

The matrices \mathbf{P} and \mathbf{W} can be found by minimizing the reconstruction error between \mathbf{v} and $\tilde{\mathbf{v}}$ over the training set.

Formally,

$$\begin{aligned} \mathbf{P}^*, \mathbf{W}^* &= \arg \min_{\mathbf{P}, \mathbf{W}} \sum_{\mathbf{v}} \|\mathbf{v} - \tilde{\mathbf{v}}\|^2 \\ &= \arg \min_{\mathbf{P}, \mathbf{W}} \sum_{\mathbf{v}} \|\mathbf{v} - \mathbf{W}\mathbf{P}\mathbf{v}\|^2. \end{aligned} \quad (3)$$

This objective can be iteratively optimized with standard stochastic gradient descent (SGD). With an orthogonality constraint, we can simply use principal component analysis (PCA) [29] to achieve a closed-form solution for the objective.

Using $\tilde{\mathbf{Y}}$ as the target, we may pre-train the network with a regression loss by observing that the compressed labels $\tilde{\mathbf{Y}}$ is real-valued

$$\mathcal{L}(\mathbf{F}, \mathbf{Y}) = \|\mathbf{F} - \tilde{\mathbf{Y}}\|^2. \quad (4)$$

Thus any regression loss, ℓ_2 being a typical example as in Eq. (4), can be employed here. Alternatively, a more direct approach is to compute the loss in the space of \mathbf{Y} . Therefore, instead of compressing \mathbf{Y} into $\tilde{\mathbf{Y}}$, we up-sample \mathbf{F} with the learned reconstruction matrix \mathbf{W} and then compute the pixel classification loss between the decompressed \mathbf{F} and \mathbf{Y} :

$$\mathcal{L}(\mathbf{F}, \mathbf{Y}) = \text{Loss}(\text{softmax}(\text{DUpsample}(\mathbf{F})), \mathbf{Y}). \quad (5)$$

With linear reconstruction, $\text{DUpsample}(\mathbf{F})$ applies linear upsampling of $\mathbf{W}\mathbf{f}$ to each feature $\mathbf{f} \in \mathbb{R}^{\tilde{C}}$ in the tensor \mathbf{F} . Comparing with Eq. (1), we have replaced the bilinear upsampling with a data-dependent upsampling, learned from the ground-truth labels. This upsampling procedure is essentially the same as applying a 1×1 convolution along the spatial dimensions, with convolutional kernels stored in \mathbf{W} . The decompression is illustrated in Fig. 3.

Note that, besides the linear upsampling presented above, we have also conducted experiments using a non-linear auto-encoder for upsampling. Training of the auto-encoder is also to minimize the reconstruction loss, and is more general than the linear case. Empirically, we observe that the final semantic prediction accuracy is almost the same as using the much simpler linear reconstruction.

Therefore we focus on using the linear reconstruction in the sequel.

Discussion with Depth-to-Space and Sub-pixel. The simplest linear form of DUpsample can be viewed as an improved version of Depth-to-Space in [28] or Sub-pixel in [26] with pre-computed upsampling filters. Depth-to-Space and Sub-pixel are typically used to upsample the inputs by a modest upsample ratio (e.g., ≤ 4), in order to avoid incurring too many trainable parameters resulting in difficulties in optimization. In contrast, as the upsampling filters in our method are pre-computed, the upsample ratio of DUpsampling can be very large (e.g., 16 or 32) if needed.

3.2. Incorporating DUpsampling with Adaptive-temperature Softmax

So far, we have shown that DUpsampling can be used to replace the incapable bilinear upsampling in semantic segmentation. The next step is to incorporate the DUpsampling into the encoder-decoder framework, resulting in an end-to-end trainable system. While DUpsampling can be realized with a 1×1 convolution operation, incorporating directly into the framework encounters difficulties in optimization.

Probably due to the \mathbf{W} is computed with one-hot encoded \mathbf{Y} , we find that the combination of vanilla softmax and DUpsampling has difficulty in producing sharp enough activation. As a result, the cross-entropy loss is stuck during the training process (as shown in experiment 4.1.4), which makes the training process slow to converge.

In order to tackle the issue, we instead employ the softmax function with temperature [13], which adds a temperature T into vanilla softmax function to sharpen/soften the activation of softmax.

$$\text{softmax}(z_i) = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}. \quad (6)$$

We find that T can be learned automatically using the standard back-propagation algorithm, eliminating the need for tuning. We show in experiments that this adaptive-temperature softmax makes training converge much faster without introducing extra hyper-parameters.

3.3. Flexible Aggregation of Convolutional Features

The extremely deep CNNs [11, 6, 14] lead to the success in computer vision. However, the depth also causes the loss of fine-grained information essential to semantic segmentation. It has been shown by a number of works [16, 5] that combining the low-level convolutional features can improve the segmentation performance significantly.

Let \mathbf{F} be the eventual CNNs feature maps used to produce the final pixel-wise prediction by bilinear or aforementioned DUpsampling. \mathbf{F}_i and \mathbf{F}_{last} represent the feature maps at level i of the backbone and last convolutional feature maps of the backbone, respectively. For simplicity

we focus on fusing one level of low-level features, but it is straightforward to extend it to multi-level fusion, which perhaps boosts the performance further. The feature aggregation in previous decoders shown in Fig. 1 can be formulated as,

$$\mathbf{F} = f(\text{concat}(\mathbf{F}_i, \text{upsample}(\mathbf{F}_{\text{last}}))), \quad (7)$$

where f denotes a CNN and upsample is usually bilinear. concat is a concatenation operator along the channel. As described above, this arrangement comes with two problems. 1) f is applied after upsampling. Since f is a CNN, whose amount of computation depends on the spatial size of inputs, this arrangement would render the decoder inefficient computationally. Moreover, the computational overhead prevents the decoder from exploiting features at a very low level. 2) The resolution of fused low-level features \mathbf{F}_i is equivalent to that of \mathbf{F} , which is typically around $\frac{1}{4}$ resolution of the final prediction due to the incapable bilinear used to produce the final pixel-wise prediction. In order to obtain high-resolution prediction, the decoder can only choose the feature aggregation with high-resolution low-level features.

In contrast, in our proposed framework, the responsibility to restore the full-resolution prediction has been largely shifted to DUpsampling. Therefore, we can safely down-sample any level of used low-level features to the resolution of last feature maps \mathbf{F}_{last} (the lowest resolution of feature maps) and then fuse these features to produce final prediction, as shown in Fig. 2. Formally, Eq. (7) is changed to,

$$\mathbf{F} = f(\text{concat}(\text{downsample}(\mathbf{F}_i), \mathbf{F}_{\text{last}})), \quad (8)$$

where downsample is bilinear in our case. This rearrangement not only keeps the features always to be computed efficiently at the lowest resolution, but also decouples the resolution of low-level features \mathbf{F}_i and that of the final segmentation prediction, allowing any level of features to be fused. In experiments, we show the flexible feature fusion enables us to exploit a better feature fusion to boost the segmentation performance as much as possible.

Only when cooperating with the aforementioned DUpsampling, the scheme of downsampling low-level features can work. Otherwise, the performance is bounded by the upper bound of the incapable upsampling method of the decoder. This is the reason why previous methods are required to upsample the low-resolution high-level feature maps back to the spatial size of fused low-level feature maps.

4. Experiments

The proposed models are evaluated on the PASCAL VOC 2012 semantic segmentation benchmark [8] and PASCAL Context benchmark [23]. For both benchmarks, we measure the performance in terms of pixel intersection-over-union averaged across the present classes (i.e., mIOU).

PASCAL VOC is the dataset widely used for semantic segmentation. It consists of 21 classes including background. The splits of PASCAL VOC are 1, 464, 1, 449 and 1, 456 for training, validation and test, respectively. The ablation study of our work is conducted over its *val* set. Also, we report our performance over *test* set to compare with other state-of-the-art methods.

PASCAL Context is much larger than PASCAL VOC, including 4, 998 images for training and 5, 105 images for validation. Following previous works [16, 23], we choose the most frequent 59 classes plus one background class (i.e., 60 classes in total) in our experiments. There is not a test server available and therefore we follow previous works [16, 34, 3, 21, 37] to report our result on *val* set.

Cityscapes is a large-scale benchmark for semantic urban scene parsing. It contains 2, 975 images for training, 500 images for validation and 1, 525 images for testing. Additionally, it also provides about 20, 000 weakly annotated images.

Implementation details. For all ablation experiments on PASCAL VOC, we opt for ResNet-50 [11] and Xception-65 [6] as our backbone networks, both of which are modified as in [5]. Following [19, 4, 5], we use “poly” as our learning rate policy for all experiments. The initial learning rate is set as 0.007 and total iteration is $30k$ for ablation experiments on PASCAL VOC. For all ResNet-based experiments, weight decay is set to 0.0001. The batch size is set to 48, but the batch normalization [15] statistics are computed with a batch of 12 images. For all Xception-based experiments, weight decay is 0.00004. We use a batch size of 32 but compute the batch normalization statistics within a batch of 16 images. We follow the practice [5, 4, 35] to use the weights pre-trained on ImageNet [7] to initialize backbone networks. All weights of newly added layers are initialized with Gaussian distribution of variance 0.01 and mean 0. T in adaptive-temperature softmax is initialized to 1. \tilde{C} is set as 64 for ResNet-50 based experiments and 128 for Xception-65 based experiments. Finally, following previous works [4, 3, 5], we augment the training data by randomly scaling the images from 0.5 to 2.0 and left-right flipping them.

4.1. Ablation Study

Our work focuses on the decoder part of the segmentation architecture. Therefore, for all ablation experiments, we use the same encoder, as shown in Fig. 1. The encoder yields the final feature maps with the $\frac{1}{16}$ or $\frac{1}{32}$ size of the original image. The decoder aims to decode the low-resolution feature maps into the prediction with the same resolution as the original image. In this section, we will investigate different decoder schemes, and demonstrate our proposed decoder’s advantages. We make use of official *train* set instead of SBD [10] since it provides more consis-

Method	output stride	mIOU (%)	mIOU* (%)
bilinear	32	70.77	94.80
DUpsampling	32	72.09	99.90
bilinear	16	72.15	98.40
DUpsampling	16	73.15	99.95

Table 1: mIOU over the PASCAL VOC *val* set of DUpsampling vs. bilinear upsampling. “output stride” indicates the ratio of input image spatial resolution to final output resolution. mIOU* denotes the upper bound.

tent annotations.

4.1.1 DUpsampling vs. Bilinear

First of all, we design experiments to show that the upper bound of bilinear is much lower than that of DUpsampling, which results in limited performance of bilinear. Specifically, we design a light-weight CNN including five convolutional layers with kernel size being 3 and stride of 2, which is fed with *ground truth labels* instead of raw images. Next, DUpsampling or bilinear is added on top of that to recover the pixel-wise prediction. This is similar to the *decoder* part in the encoder-decoder architecture.

By training the two networks, with DUpsampling or bilinear as decoder respectively, the ability to restore the pixel-wise prediction can be quantitatively measured via their performance over the *val* set, which can be viewed as the upper bound of both methods. We use the training protocol described in implementation details to train the two networks, except that the total iterations and initial learning rate are set as $100k$ and 0.07, respectively. “output stride” indicates the ratio of input image spatial resolution to the final CNN feature maps resolution. As shown in Table 1, the upper bound performance of DUpsampling is well above that of bilinear both when output stride being 32 and 16.

Given the superior upper bound performance of DUpsampling, we further carry out experiments with raw input images. In the experiments, we employ ResNet-50 as the backbone network. Unsurprisingly, by merely replacing the bilinear with DUpsampling, the mIOU on PASCAL VOC *val* set is improved by 1.3 points and 1 point, when the output stride is 32 and 16 respectively, as shown in Table 1. The improvement is significant because mIOU is strict. Interestingly, the DUpsampling of output stride being 32 achieves similar performance to the bilinear case of output stride being 16. This shows that the proposed DUpsampling may eliminate the need for expensive computationally high-resolution feature maps from the CNNs.

4.1.2 Flexible aggregation of convolutional features

Due to the flexibility of our proposed decoder, we can employ any combination of features to improve segmentation performance, regardless of the resolution of fused features.

Used low-level features	mIOU (%)	FLOPS
N/A	73.15	0.80B
conv1_3	72.70	1.13B
b1u2c3	74.03	1.15B
b3u6c3	73.43	1.23B
b1u2c3 + b3u6c3	73.82	1.58B
conv1_3 + b3u6c3	74.20	1.56B

Table 2: mIOU over PASCAL VOC *val* set when using different fusion of features. *bxuycz* denotes low-level features named block_*x*/unit_*y*/conv_*z* in ResNet. "FLOPS" denotes the amount of computation of the decoder including feature aggregation, convolutional decoder and the final upsampling.

For ResNet-50, we experiment with many different combinations of features, as shown in Table 2. The best one is the combination of conv1_3 + b3u6c3, achieving mIOU 74.20% over *val* set. Additionally, as shown in Table 2, the amount of computation changes little when features at different levels are fused, which allows us to choose the best feature fusion without considering the price of computation incurred by the resolution of fused features.

In order to understand how the fusion works, we visualize the segmentation results with and without low-level features in Fig. 4. Intuitively, the one fusing low-level features yields more consistent segmentation, which suggests the downsampled low-level features are still able to refine the segmentation prediction substantially.

4.1.3 Comparison with the vanilla bilinear decoder

We further compare our proposed decoder scheme with the vanilla bilinear decoder shown in Fig. 1, which fuses low-level features b1u2c3 (downsample ratio = 4). As shown in Table 3, it achieves mIOU 73.26% on *val* set with ResNet-50 as the backbone. By replacing vanilla decoder with our proposed decoder in Fig. 2, the performance is improved to 74.03%. Because of the same low-level features used, the improvement should be due to the capable DUpsampling instead of bilinear used to restore the full-resolution prediction. Furthermore, we explore a better feature fusion conv1_3 + b3u6c3 for proposed decoder and improve the overall performance slightly to 74.20%. When the



Figure 4: The prediction results with low-level features and without low-level features. ResNet-50 is used as the backbone.

Decoder	Low-level features / ratio	mIOU (%)	FLOPS
ResNet-50			
Vanilla	b1u2c3 / 4	73.26	5.53B
Proposed	b1u2c3 / 4	74.03	1.15B
Vanilla	conv1_3 / 2 + b3u6c3 / 16	-	22.34B
Proposed	conv1_3 / 2 + b3u6c3 / 16	74.20	1.56B
Xception-65			
Vanilla	efb2u1c2 / 4	78.70	5.53B
Proposed	efb2u1c2 / 4	79.09	1.93B
Vanilla	mfb1u16c3 / 16	78.74	0.41B
Proposed	mfb1u16c3 / 16	79.67	1.98B

Table 3: mIOU over the PASCAL VOC *val* set when using different fusion strategies of features. *bxuycz* denotes low-level features named block_*x*/unit_*y*/conv_*z* in ResNet or Xception. "ef" and "mf" respectively indicate "entry_flow" and "middle_flow" in Xception. "-" means out-of-memory. "ratio" denotes the ratio of the resolution of feature maps to the resolution of the input image (i.e., downsample ratio). "FLOPS" denotes the amount of computation of the decoders.

Decoder	low-level features / ratio	mIOU (%)	FLOPS
Vanilla	efb2u1c2 / 4	79.36	43.65B
Proposed	mfb1u16c3 / 16	79.06	25.14B

Table 4: mIOU on the Cityscapes *val* set. Our proposed decoder with much less computation complexity achieves a similar performance as the vanilla decoder.

vanilla decoder uses the fusion of features, it incurs much heavier computation complexity and runs out of our GPUs memory due to the high resolution of conv1_3, which prevents the vanilla decoder from exploiting the low-level features.

We also experiment our proposed decoder with Xception-65 as the backbone. Similarly, with the same low-level features efb2u1c3 (downsample ratio = 4), our proposed decoder improves the performance from 78.70% to 79.09%, as shown in Table 3. When using a better low-level features mfb1u16c3 (downsample ratio = 16), the vanilla decoder just improves the performance negligibly by 0.04% because its performance is constrained by the incapable bilinear upsampling used to restore the full-resolution prediction. In contrast, our proposed decoder can still benefit a lot from the better feature fusion due to the use of much powerful DUpsampling. As shown in Table 3, with the better feature fusion, the performance of our proposed decoder is improved to 79.67%. Moreover, since we downsample low-level features before fusing, our proposed decoder requires much fewer FLOPS than the vanilla decoder of the best performance, as shown in Table 3.

Finally, we compare our proposed decoder with the vanilla bilinear decoder on the Cityscapes *val* set. Follow-

Method	mIOU (%)
PSPNet [36]	85.4
DeepLabv3 [4]	85.7
EncNet [34]	85.9
DFN [32]	86.2
IDW-CNN [27]	86.3
CASIA_IVA_SDN [9]	86.6
DIS [22]	86.8
DeepLabv3+ [5] (Xception-65)	87.8
Our proposed (Xception-65)	88.1

Table 5: State-of-the-art methods on PASCAL VOC *test* set.

ing [5], Xception-71 is used as our backbone and the number of iterations is increased to $90k$ with a initial learning rate being 0.01. As shown in Table 4, under the same training and testing settings, our proposed decoder achieves a comparable performance with the vanilla one while using much less computation.

4.1.4 Impact of adaptive-temperature softmax

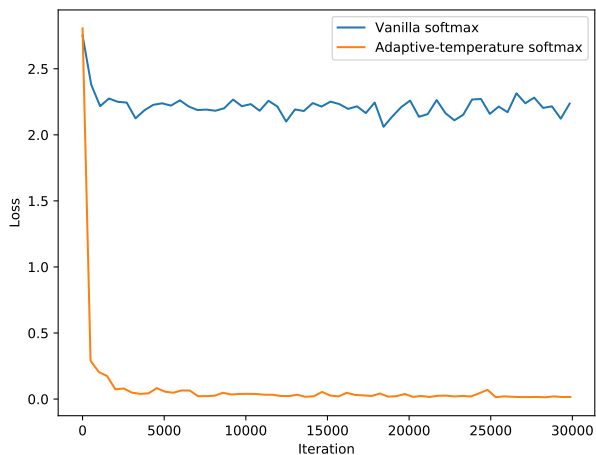


Figure 5: Training losses for vanilla softmax and adaptive-temperature softmax.

As mentioned before, the adaptive-temperature softmax eases the training of the proposed DUpsampling method. When training the framework with vanilla softmax with T being 1, it achieves 69.81% over *val* set, which is significantly lower than 73.15% of the counterpart with adaptive-temperature softmax. We further plot training losses for vanilla softmax and adaptive-temperature softmax in Fig. 5, which shows the advantage of this adaptive-temperature softmax.

Method	mIOU (%)
FCN-8s [21]	37.8
CRF-RNN [37]	39.3
HO_CRF [1]	41.3
Piecewise [17]	43.3
VeryDeep [30]	44.5
DeepLabv2 [3]	45.7
RefineNet [16]	47.3
EncNet [34]	51.7
Our proposed (Xception-65)	51.4
Our proposed (Xception-71)	52.5

Table 6: State-of-the-art methods on PASCAL Context *val* set.

4.2. Comparison with state-of-the-art Methods

Finally, we compare the framework of our proposed decoder with state-of-the-art methods. To compete with these state-of-the-art methods, we choose Xception-65 as the backbone network and the best feature aggregation in the ablation study for our decoder.

Following previous methods, SBD [10] and COCO [18] are used to train the model as well. Specifically, the model is successively trained over COCO, SBD and PASCAL VOC *trainval* set, with the training protocol described in implementation details. Each round is initialized with the last round model and the base learning rate is reduced accordingly (i.e. 0.007 for COCO, 0.001 for SBD and 0.0001 for *trainval*). We use $500k$ iterations when training over COCO and $30k$ iterations for the last two rounds. Additionally, following previous works [4, 5], we make use of multi-scale testing and left-right flipping when inferring over *test* set.

As shown in Table 5, our framework sets the new record on PASCAL VOC and improve the previous method DeepLabv3+ with the same backbone by 0.3%, which is significant due to the benchmark has been very competitive. Meanwhile, since our proposed decoder can eliminate the need for high-resolution feature maps, we employ output stride being 16 instead of 8 in DeepLabv3+ when inferring over *test* set. As a result, our whole framework only takes 30% computation of DeepLabv3+ (897.94B vs. 3055.35B in Multiply-Adds) to achieve the state-of-the-art performance. The performance of our proposed framework on PASCAL Context *val* set is shown in Table 6. With Xception-71 as backbone, our framework sets the new state-of-the-art on this benchmark dataset without pre-training on COCO.

5. Conclusion

We have proposed a flexible and light-weight decoder scheme for semantic image segmentation. This novel decoder employs our proposed DUpsampling to produce the pixel-wise prediction, which eliminates the need for compu-

tationally inefficient high-resolution feature maps from the underlying CNNs and decouples the resolution of the fused low-level features and that of the final prediction. This decoupling expands the design space of feature aggregation of the decoder, allowing almost arbitrary features aggregation to be exploited to boost the segmentation performance as much as possible. Meanwhile, our proposed decoder avoids upsampling low-resolution high-level feature maps back to the spatial size of high-resolution low-level feature maps, reducing the computation of decoder remarkably. Experiments demonstrate that our proposed decoder has advantages of both effectiveness and efficiency over the vanilla decoder extensively used in previous semantic segmentation methods. Finally, the framework with the proposed decoder attains the state-of-the-art performance while requiring much less computation than previous state-of-the-art methods.

Acknowledgments The authors would like to thank Huawei Technologies for the donation of GPU cloud computing resources.

References

- [1] Anurag Arnab, Sadeep Jayasumana, Shuai Zheng, and Philip HS Torr. Higher order conditional random fields in deep neural networks. In *Proc. Eur. Conf. Comp. Vis.*, pages 524–540. Springer, 2016.
- [2] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. SegNet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, (12):2481–2495, 2017.
- [3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(4):834–848, 2018.
- [4] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- [5] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. *Proc. Eur. Conf. Comp. Vis.*, 2018.
- [6] François Chollet. Xception: Deep learning with depthwise separable convolutions. *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 1610–02357, 2017.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 248–255. Ieee, 2009.
- [8] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *Int. J. Comp. Vis.*, 88(2):303–338, 2010.
- [9] Jun Fu, Jing Liu, Yuhang Wang, and Hanqing Lu. Stacked deconvolutional network for semantic segmentation. *arXiv preprint arXiv:1708.04943*, 2017.
- [10] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. 2011.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 770–778, 2016.
- [12] Tong He, Chunhua Shen, Zhi Tian, Dong Gong, Changming Sun, and Youliang Yan. Knowledge adaptation for efficient semantic segmentation. *arXiv preprint arXiv:1903.04688*, 2019.
- [13] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [14] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Weinberger. Densely connected convolutional networks. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*
- [15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [16] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian D Reid. RefineNet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, volume 1, page 5, 2017.
- [17] Guosheng Lin, Chunhua Shen, Anton Van Den Hengel, and Ian Reid. Efficient piecewise training of deep structured models for semantic segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 3194–3203, 2016.
- [18] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *Proc. Eur. Conf. Comp. Vis.*, pages 740–755. Springer, 2014.
- [19] Wei Liu, Andrew Rabinovich, and Alexander C Berg. Parsenet: Looking wider to see better. *arXiv preprint arXiv:1506.04579*, 2015.
- [20] Ziwei Liu, Xiao Xiao Li, Ping Luo, Chen-Change Loy, and Xiaoou Tang. Semantic image segmentation via deep parsing network. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 1377–1385, 2015.
- [21] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 3431–3440, 2015.
- [22] Ping Luo, Guangrun Wang, Liang Lin, and Xiaogang Wang. Deep dual learning for semantic image segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 21–26, 2017.
- [23] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 891–898, 2014.
- [24] Hyeonwoo Noh, Seunghoon Hong, and Bohyung Han. Learning deconvolution network for semantic segmentation. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 1520–1528, 2015.

- [25] Chao Peng, Xiangyu Zhang, Gang Yu, Guiming Luo, and Jian Sun. Large kernel matters improve semantic segmentation by global convolutional network. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 1743–1751. IEEE, 2017.
- [26] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 1874–1883, 2016.
- [27] Guangrun Wang, Ping Luo, Liang Lin, and Xiaogang Wang. Learning object interactions and descriptions for semantic image segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 5859–5867, 2017.
- [28] Zbigniew Wojna, Vittorio Ferrari, Sergio Guadarrama, Nathan Silberman, Liang-Chieh Chen, Alireza Fathi, and Jasper Uijlings. The devil is in the decoder. *arXiv preprint arXiv:1707.05847*, 2017.
- [29] Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, 2(1-3):37–52, 1987.
- [30] Zifeng Wu, Chunhua Shen, and Anton van den Hengel. Bridging category-level and instance-level semantic image segmentation. *arXiv preprint arXiv:1605.06885*, 2016.
- [31] Zifeng Wu, Chunhua Shen, and Anton van den Hengel. Wider or deeper: Revisiting the resnet model for visual recognition. *arXiv preprint arXiv:1611.10080*, 2016.
- [32] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Learning a discriminative feature network for semantic segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018.
- [33] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [34] Hang Zhang, Kristin Dana, Jianping Shi, Zhongyue Zhang, Xiaogang Wang, Amrith Tyagi, and Amit Agrawal. Context encoding for semantic segmentation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018.
- [35] Pingping Zhang, Wei Liu, Hongyu Wang, Yinjie Lei, and Huchuan Lu. Deep gated attention networks for large-scale street-level scene segmentation. *Pattern Recognition*, 88:702–714, 2019.
- [36] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 2881–2890, 2017.
- [37] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 1529–1537, 2015.

Supplementary Material:

In this supplementary material, we 1) provide our result on PASCAL VOC [8] *test* set without COCO [18] pre-training and 2) showcase more visualization results of our proposed method.

6. PASCAL VOC without COCO Pre-training

In this experiment, following previous works [31, 36, 34] without COCO pre-training, we train our model on SBD [10] and then fine-tune it on official *trainval* set. We use the same training protocol as described in the main paper. The multi-scale testing and left-right flipping are employed when our model is evaluated on *test* set. No any post-processing is used. The final performance is obtained by uploading our test results to the official test server. As shown in Table 7, our proposed framework surpasses previous published methods by a large margin.

Method	mIOU (%)
DPN [20]	74.1
Piecewise [17]	75.3
ResNet-38 [31]	82.5
PSPNet [36]	82.6
DFN [32]	82.7
EncNet [34]	82.9
Our proposed (Xception-65)	85.3

Table 7: State-of-the-art methods on PASCAL VOC *test* set without COCO pre-training.

7. Visualization

The visualization results of our method are shown in Fig. 6 and Fig. 7. As shown in Fig. 6, without any post-processing, the proposed method works very well in a lot of challenging cases. Small, distant and incomplete objects can be segmented well. Meanwhile, as shown in Fig. 7, although we employ “output stride” being 16 when evaluating, which results in low-resolution CNNs output feature maps, our model can still yield fine-grained segmentation due to the use of proposed DUpSampling.



Figure 6: Visualization results from *val* set. The proposed method works reliably in a lot of challenging cases including small, distant and incomplete objects.



Figure 7: Visualization results from *val* set. The proposed method can yield fine-grained segmentation, with low-resolution CNNs output feature maps.