

# Segmenter: Transformer for Semantic Segmentation

Robin Strudel\*  
Inria<sup>†</sup>

Ricardo Garcia\*  
Inria<sup>†</sup>

Ivan Laptev  
Inria<sup>†</sup>

Cordelia Schmid  
Inria<sup>†</sup>

## Abstract

Image segmentation is often ambiguous at the level of individual image patches and requires contextual information to reach label consensus. In this paper we introduce Segmenter, a transformer model for semantic segmentation. In contrast to convolution-based methods, our approach allows to model global context already at the first layer and throughout the network. We build on the recent Vision Transformer (ViT) and extend it to semantic segmentation. To do so, we rely on the output embeddings corresponding to image patches and obtain class labels from these embeddings with a point-wise linear decoder or a mask transformer decoder. We leverage models pre-trained for image classification and show that we can fine-tune them on moderate sized datasets available for semantic segmentation. The linear decoder allows to obtain excellent results already, but the performance can be further improved by a mask transformer generating class masks. We conduct an extensive ablation study to show the impact of the different parameters, in particular the performance is better for large models and small patch sizes. Segmenter attains excellent results for semantic segmentation. It outperforms the state of the art on the challenging ADE20K dataset and performs on-par on Pascal Context and Cityscapes.

## 1. Introduction

Semantic segmentation is a challenging computer vision problem with a wide range of applications including autonomous driving, robotics, augmented reality, image editing, medical imaging and many others [26, 27, 43]. The goal of semantic segmentation is to assign each image pixel to a category label corresponding to the underlying object and to provide high-level image representations for target tasks, e.g. detecting the boundaries of people and their clothes for virtual try-on applications [28].

\*Equal contribution.

<sup>†</sup>Inria, École normale supérieure, CNRS, PSL Research University, 75005 Paris, France.

Code: <https://github.com/rstrudel/segmenter>

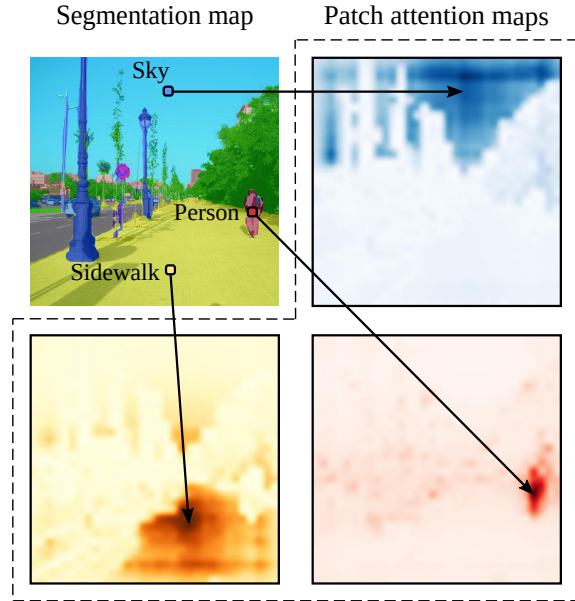


Figure 1: Our approach for semantic segmentation is purely transformer based. It leverages the global image context at every layer of the model. Attention maps from the first Segmenter layer are displayed for three  $8 \times 8$  patches and highlight the early grouping of patches into semantically meaningful categories. The original image (top-left) is overlaid with segmentation masks produced by our method.

Despite much effort and large progress over recent years [9, 20, 30, 35, 45, 60, 61], image segmentation remains a challenging problem due to rich intra-class variation, context variation and ambiguities originating from occlusions and low image resolution.

Recent approaches to semantic segmentation typically rely on convolutional encoder-decoder architectures where the encoder generates low-resolution image features and the decoder upsamples features to segmentation maps with per-pixel class scores. State-of-the-art methods deploy Fully Convolutional Networks (FCN) [42] and achieve impressive results on challenging segmentation benchmarks [9, 21, 53, 54, 56, 58, 61]. These methods rely on learnable stacked convolutions that can capture semantically rich information

and have been highly successful in computer vision. The local nature of convolutional filters, however, limits the access to the global information in the image. Meanwhile, such information is particularly important for segmentation where the labeling of local patches often depends on the global image context. To circumvent this issue, DeepLab methods [7, 8, 9] introduce feature aggregation with dilated convolutions and spatial pyramid pooling. This allows to enlarge the receptive fields of convolutional networks and to obtain multi-scale features. Following recent progresses in NLP [47], several segmentation methods explore alternative aggregation schemes based on channel or spatial [20, 21, 57] attention and point-wise [61] attention to better capture contextual information. Such methods, however, still rely on convolutional backbones and are, hence, biased towards local interactions. An extensive use of specialised layers to remedy this bias [7, 9, 20, 54] suggests limitations of convolutional architectures for segmentation.

To overcome these limitations, we formulate the problem of semantic segmentation as a sequence-to-sequence problem and use a transformer architecture [47] to leverage contextual information at every stage of the model. By design, transformers can capture global interactions between elements of a scene and have no built-in inductive prior, see Figure 1. However, the modeling of global interactions comes at a quadratic cost which makes such methods prohibitively expensive when applied to raw image pixels [10]. Following the recent work on Vision Transformers (ViT) [17, 46], we split the image into patches and treat linear patch embeddings as input tokens for the transformer encoder. The contextualized sequence of tokens produced by the encoder is then upsampled by a transformer decoder to per-pixel class scores. For decoding, we consider either a simple point-wise linear mapping of patch embeddings to class scores or a transformer-based decoding scheme where learnable class embeddings are processed jointly with patch tokens to generate class masks. We conduct an extensive study of transformers for segmentation by ablating model regularization, model size, input patch size and its trade-off between accuracy and performance. Our Segmenter approach attains excellent results while remaining simple, flexible and fast. In particular, when using large models with small input patch size the best model reaches a mean IoU of 50.77% on ADE20K, surpassing all previous state-of-the-art convolutional approaches by a large margin of 4.6%. Such improvement partly stems from the global context captured by our method at every stage of the model as highlighted in Figure 1.

In summary, our work provides the following four contributions: (i) We propose a novel approach to semantic segmentation based on the Vision Transformer (ViT) [17] that does not use convolutions, captures contextual information by design and outperforms FCN based approaches. (ii) We present a family of models with varying levels of resolu-

tion which allows to trade-off between precision and runtime, ranging from *state-of-the-art* performance to models with fast inference and good performances. (iii) We propose a transformer-based decoder generating class masks which outperforms our linear baseline and can be extended to perform more general image segmentation tasks. (iv) We demonstrate that our approach yields *state-of-the-art* results on the challenging ADE20K dataset [63] and is on-par on Pascal Context and Cityscapes benchmarks.

## 2. Related work

**Semantic segmentation.** Methods based on Fully Convolutional Networks (FCN) combined with encoder-decoder architectures have become the dominant approach to semantic segmentation. Initial approaches [19, 34, 37, 38] rely on a stack of consecutive convolutions followed by spatial pooling to perform dense predictions. Consecutive approaches [1, 3, 33, 39, 41] upsample high-level feature maps and combine them with low-level feature maps during decoding to both capture global information and recover sharp object boundaries. To enlarge the receptive field of convolutions in the first layers, several approaches [7, 32, 55] have proposed dilated or atrous convolutions. To capture global information in higher layers, recent work [8, 9, 60] employs spatial pyramid pooling to capture multi-scale contextual information. Combining these enhancements along with atrous spatial pyramid pooling, DeeplabV3+ [9] proposes a simple and effective encoder-decoder FCN architecture. Recent work [20, 21, 53, 54, 57, 61] replace coarse pooling by attention mechanisms on top of the encoder feature maps to better capture long-range dependencies.

While recent segmentation methods are mostly focused on improving FCN, the restriction to local operations imposed by convolutions may imply inefficient processing of global image context and suboptimal segmentation results. Hence, we propose a pure transformer architecture that captures global context at every layer of the model during the encoding and decoding stages.

**Transformers for vision.** Transformers [47] are now state of the art in many Natural Language Processing (NLP) tasks. Such models rely on self-attention mechanisms and capture long-range dependencies among tokens (words) in a sentence. In addition, transformers are well suited for parallelization, facilitating training on large datasets. The success of transformers in NLP has inspired several methods in computer vision combining CNNs with forms of self-attention to address object detection [6], semantic segmentation [49], panoptic segmentation [48], video processing [50] and few-shot classification [16]. Several concurrent works use pure transformer approaches to address computer vision problems, namely 3D point cloud processing [59], segmentation [62] and video classification [5].

The Vision Transformer (ViT) [17] introduces a convolution-free transformer architecture for image classi-

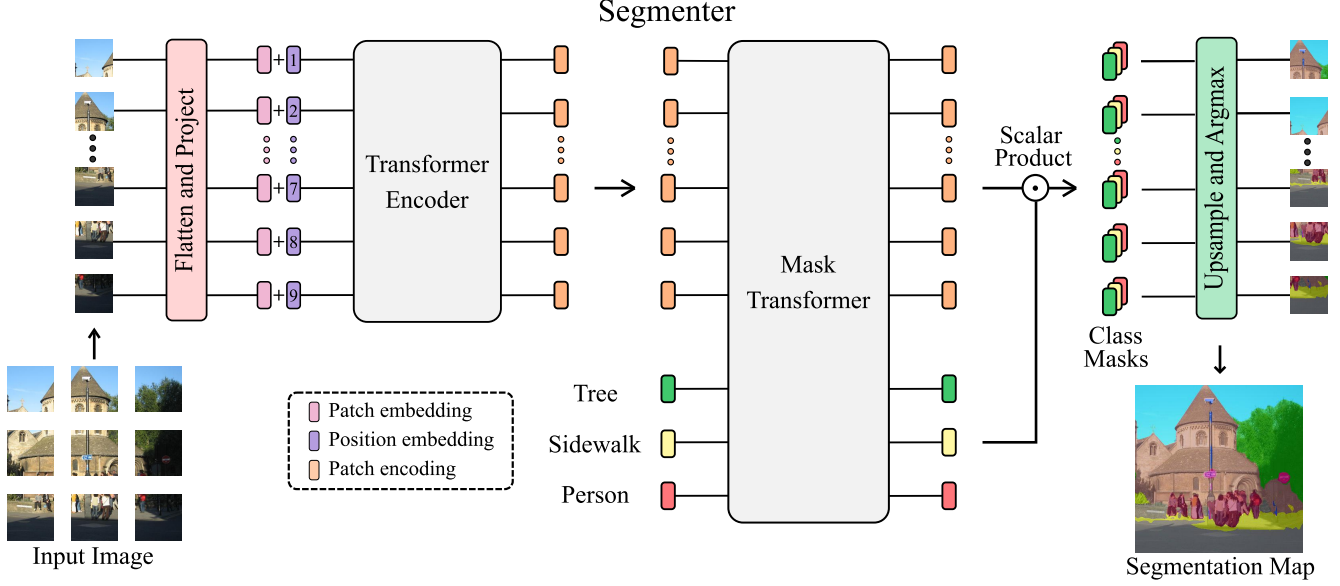


Figure 2: Overview of our approach Segementer. (Left) Encoder: The image patches are projected to a sequence of embeddings and then encoded with a transformer. (Right) Decoder: A mask transformer takes as input the output of the encoder and class embeddings to predict segmentation masks. See text for details.

fication where input images are processed as sequences of patch tokens. While ViT requires training on very large datasets, DeiT [46] proposes a token-based distillation strategy and obtains a competitive vision transformer trained on the ImageNet-1k [14] dataset using a CNN as a teacher. Inspired by this line of work, we propose Segementer, a transformer encoder-decoder architecture for semantic image segmentation. Our architecture does not use convolutions, captures global image context by design and results in competitive performance on standard image segmentation benchmarks.

### 3. Our approach: Segementer

Segementer is based on a fully transformer-based encoder-decoder architecture mapping a sequence of patch embeddings to pixel-level class annotations. An overview of the model is shown in Figure 2. The sequence of patches is encoded by a transformer encoder described in Section 3.1 and decoded by either a point-wise linear mapping or a mask transformer described in Section B. Our model is trained end-to-end with a per-pixel cross-entropy loss. At inference time, argmax is applied after upsampling to obtain a single class per pixel.

#### 3.1. Encoder

An image  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$  is split into a sequence of patches  $\mathbf{x} = [x_1, \dots, x_N] \in \mathbb{R}^{N \times P^2 \times C}$  where  $(P, P)$  is the patch size,  $N = HW/P^2$  is the number of patches and  $C$  is the number of channels. Each patch is flattened

into a 1D vector and then linearly projected to a patch embedding to produce a sequence of patch embeddings  $\mathbf{x}_0 = [\mathbf{E}x_1, \dots, \mathbf{E}x_N] \in \mathbb{R}^{N \times D}$  where  $\mathbf{E} \in \mathbb{R}^{(P^2 C) \times D}$ . To capture positional information, learnable position embeddings  $\mathbf{pos} = [\mathbf{pos}_1, \dots, \mathbf{pos}_N] \in \mathbb{R}^{N \times D}$  are added to the sequence of patches to get the resulting input sequence of tokens  $\mathbf{z}_0 = \mathbf{x}_0 + \mathbf{pos}$ .

A transformer [47] encoder composed of  $L$  layers is applied to the sequence of tokens  $\mathbf{z}_0$  to generate a sequence of contextualized encodings  $\mathbf{z}_L \in \mathbb{R}^{N \times D}$ . A transformer layer consists of a multi-headed self-attention (MSA) block followed by a point-wise MLP block of two layers with layer norm (LN) applied before every block and residual connections added after every block:

$$\mathbf{a}_{i-1} = \text{MSA}(\text{LN}(\mathbf{z}_{i-1})) + \mathbf{z}_{i-1}, \quad (1)$$

$$\mathbf{z}_i = \text{MLP}(\text{LN}(\mathbf{a}_{i-1})) + \mathbf{a}_{i-1}, \quad (2)$$

where  $i \in \{1, \dots, L\}$ . The self-attention mechanism is composed of three point-wise linear layers mapping tokens to intermediate representations, queries  $\mathbf{Q} \in \mathbb{R}^{N \times d_q}$ , keys  $\mathbf{K} \in \mathbb{R}^{N \times d_k}$  and values  $\mathbf{V} \in \mathbb{R}^{N \times d_v}$ . Self-attention is then computed as follows

$$\text{MSA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}. \quad (3)$$

The transformer encoder maps the input sequence  $\mathbf{z}_0 = [z_{0,1}, \dots, z_{0,N}]$  of embedded patches with position encoding to  $\mathbf{z}_L = [z_{L,1}, \dots, z_{L,N}]$ , a contextualized encoding se-

quence containing rich semantic information used by the decoder. In the following section we introduce the decoder.

### 3.2. Decoder

The sequence of patch encodings  $\mathbf{z}_L \in \mathbb{R}^{N \times D}$  is decoded to a segmentation map  $\mathbf{s} \in \mathbb{R}^{H \times W \times K}$  where  $K$  is the number of classes. The decoder learns to map patch-level encodings coming from the encoder to patch-level class scores. Next these patch-level class scores are upsampled by bilinear interpolation to pixel-level scores. We describe in the following a linear decoder, which serves as a baseline, and our approach, a mask transformer, see Figure 2.

**Linear.** A point-wise linear layer is applied to the patch encodings  $\mathbf{z}_L \in \mathbb{R}^{N \times D}$  to produce patch-level class logits  $\mathbf{z}_{lin} \in \mathbb{R}^{N \times K}$ . The sequence is then reshaped into a 2D segmentation map  $\mathbf{s}_{lin} \in \mathbb{R}^{H/P \times W/P \times K}$  and bilinearly upsampled to the original image size to obtain the final segmentation map  $\mathbf{s} \in \mathbb{R}^{H \times P \times K}$ .

**Mask Transformer.** For the transformer-based decoder, we introduce a set of  $K$  learnable class embeddings  $\mathbf{cls} = [\mathbf{cls}_1, \dots, \mathbf{cls}_K] \in \mathbb{R}^{K \times D}$  where  $K$  is the number of classes. Each class embedding is initialized randomly and assigned to a single semantic class. It will be used to generate the class mask. The class embeddings  $\mathbf{cls}$  are processed jointly with patch encodings  $\mathbf{z}_L$  by the decoder as depicted in Figure 2. The decoder is a transformer encoder composed of  $M$  layers. Our mask transformer generates  $K$  masks by computing the scalar product between patch embeddings  $\mathbf{z}_{mask} \in \mathbb{R}^{N \times D}$  and class embeddings  $\mathbf{c} \in \mathbb{R}^{K \times D}$  output by the decoder. The set of class masks is computed as follows

$$\text{Masks}(\mathbf{z}_{mask}, \mathbf{c}) = \text{softmax} \left( \frac{\mathbf{z}_{mask} \mathbf{c}^T}{\sqrt{D}} \right), \quad (4)$$

where the softmax is applied over the class dimension.  $\mathbf{m} = \text{Masks}(\mathbf{z}_{mask}, \mathbf{c}) \in \mathbb{R}^{N \times K}$  is a set of patch sequences, for instance  $m_{i,j}$  denotes the probability that patch  $i$  belongs to class  $j$ . The mask sequences are softly exclusive to each other i.e.  $\sum_{j=1}^K m_{i,j} = 1^N$  for all  $i \in \{1, \dots, N\}$ . Each mask sequence is then reshaped into a 2D mask to form  $\mathbf{s}_{mask} \in \mathbb{R}^{H/P \times W/P \times K}$  and bilinearly upsampled to the original image size to obtain the final segmentation map  $\mathbf{s} \in \mathbb{R}^{H \times P \times K}$ . To enhance decoding and perform fine-grained segmentation, the patch encodings  $\mathbf{z}_L$  are bilinearly upsampled by a factor of two before being passed through the decoder.

Our mask transformer is inspired by Max-DeepLab [48] and SOLO-v2 [51] which introduce object embeddings called memory to produce instance masks. However, unlike our method, MaxDeep-Lab uses an hybrid approach based on CNNs and transformers and splits the pixel and class embeddings into two streams because of computational constraints. Using a pure transformer architecture and leveraging patch level encodings, we propose a simple approach

Model	Backbone	Layers	Token size	MLP size	Heads	Params
Seg-S <sup>†</sup>	DeiT-S	12	384	1536	6	22M
Seg-B <sup>†</sup>	DeiT-B	12	768	3072	12	86M
Seg-B	ViT-B	12	768	3072	12	86M
Seg-L	ViT-L	24	1024	4096	16	307M

Table 1: Details of Transformer variants.

that processes the patch and class embeddings jointly during the decoding phase. While we address semantic segmentation in this work, our mask transformer can be directly adapted to perform panoptic segmentation by replacing the class embeddings by object embeddings.

## 4. Experimental results

### 4.1. Datasets and metrics

**ADE20K** [63]. This dataset contains challenging scenes with fine-grained labels and is one of the most challenging semantic segmentation datasets. The training set contains 20,210 images with 150 semantic classes. The validation and test set contain 2000 and 3352 images respectively.

**Pascal Context** [36]. The training set contains 4998 images with 59 semantic classes plus a background class. The validation set contains 5105 images.

**Cityscapes** [13]. The dataset contains 5,000 images from 50 different cities with 19 semantic classes. There are 2,975 images in the training set, 500 images in the validation set and 1,525 images in the test set.

**Metrics.** We report Intersection over Union (mIoU) averaged over all classes and further report pixel accuracy for ADE20K.

### 4.2. Implementation details

**Transformer models.** For the encoder, we build upon the vision transformer ViT [17] and consider "Small", "Base" and "Large" models described in Table 1. The parameters varying in the transformer encoder are the number of layers and the token size. The head size of a multi-headed self-attention (MSA) block is fixed to 64, the number of heads is the token size divided by the head size and the hidden size of the MLP following MSA is four times the token size. DeiT [46] is a variant of the vision transformer, but trained differently as described in the next section. We consider models representing the image at different resolutions and use input patch sizes  $8 \times 8$ ,  $16 \times 16$  and  $32 \times 32$ . In the following, we use an abbreviation to describe the model variant and patch size, for instance Seg-B/16 denotes the "Base" variant with  $16 \times 16$  input patch size. Models based on DeiT are denoted with a <sup>†</sup>, for instance Seg-B<sup>†</sup>/16.

**ImageNet pre-training.** Our Segmter models are pre-trained on ImageNet, ViT is pre-trained on ImageNet-21k with random cropping, and its data-efficient variant DeiT is pre-trained on ImageNet-1k with strong data augmentation.



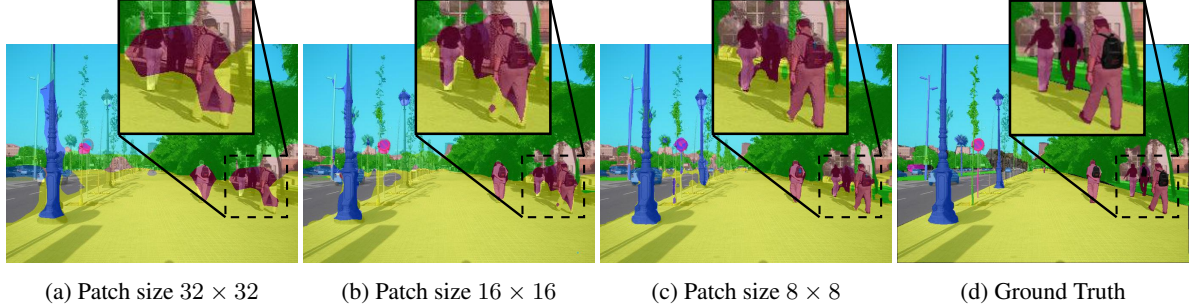


Figure 3: Impact of the model patch size on the segmentation maps.

We use publicly available models provided by the image classification library timm [52] and Google research [18]. Both models are pre-trained at an image resolution of 224 and fine-tuned on ImageNet-1k at a resolution of 384. We keep the patch size fixed and fine-tune the models for the semantic segmentation task at higher resolution depending on the dataset. As the patch size is fixed, increasing resolution results in longer token sequences. Following [17], we bilinearly interpolate the pre-trained position embeddings according to their original position in the image to match the fine-tuning sequence length. The decoders, described in Section B are initialized with random weights from a truncated normal distribution [23].

**Data augmentation.** During training, we follow the standard pipeline from the semantic segmentation library MM-Segmentation [12], which does mean subtraction, random resizing of the image to a ratio between 0.5 and 2.0 and random left-right flipping. We randomly crop large images and pad small images into a fixed size of  $512 \times 512$  for ADE20K,  $480 \times 480$  for Pascal-Context and  $768 \times 768$  for Cityscapes.

**Optimization.** To fine-tune the pre-trained models for the semantic segmentation task, we use the standard pixel-wise cross-entropy loss without weight rebalancing. We use stochastic gradient descent (SGD) [40] as the optimizer with a base learning rate  $\gamma_0$  and set weight decay to 0. Following the seminal work of DeepLab [32] we adopt the “poly” learning rate decay  $\gamma = \gamma_0 (1 - \frac{N_{iter}}{N_{total}})^{0.9}$  where  $N_{iter}$  and  $N_{total}$  represent the current iteration number and the total iteration number. For ADE20K, we set the base learning rate  $\gamma_0$  to  $10^{-3}$  and train for 160K iterations with a batch size of 8. For Pascal Context, we set  $\gamma_0$  to  $10^{-3}$  and train for 80K iterations with a batch size of 16. For Cityscapes, we set  $\gamma_0$  to  $10^{-2}$  and train for 81K iterations with a batch size of 8. The schedule is similar to DeepLab-V3+ [9] with learning rates divided by a factor 10 except for Cityscapes where we use a factor of 1.

**Inference.** To handle varying image sizes during inference, we use a sliding-window with a resolution matching the training size. For multi-scale inference, following standard practice [9] we use rescaled versions of the image with scaling factors of (0.5, 0.75, 1.0, 1.25, 1.5, 1.75) and left-right

Regularization		Model	
Dropout	Stochastic Depth	Seg-B	Seg-B <sup>†</sup>
✗	✗	45.21	45.90
✓	✗	45.06	42.82
✗	✓	<b>45.69</b>	<b>47.10</b>
✓	✓	45.16	43.61

Table 2: Comparison of different regularization schemes on ADE20K validation set with mean IoU.

Method	Backbone	Patch size	ImNet acc.	mIoU
DeepLab V3+	ResNet50-D	-	-	43.95
DeepLab V3+	ResNet101-D	-	-	<b>45.47</b>
Seg-S <sup>†</sup> /16	DeiT-S	16	81.20	42.40
Seg-B <sup>†</sup> /16	DeiT-B	16	<b>85.20</b>	<b>47.10</b>
Seg-B/32	ViT-B	32	81.65	40.92
Seg-B/16	ViT-B	16	83.95	45.69
Seg-B/8	ViT-B	8	<b>85.35</b>	<b>48.06</b>
Seg-L/32	ViT-L	32	81.68	42.64
Seg-L/16	ViT-L	16	<b>84.95</b>	<b>48.60</b>

Table 3: Performance comparison of different Segmenter models with varying backbones and input patch sizes on ADE20K validation set.

flipping and average the results.

### 4.3. Ablation study

In this section, we ablate different variants of our approach on the ADE20K validation set. We investigate model regularization, model size, patch size, model performance, training dataset size, compare Segmenter to convolutional approaches and evaluate different decoders. Unless stated otherwise we use the baseline linear decoder and report results using single-scale inference.

**Regularization.** We first compare two forms of regularization, dropout [44] and stochastic depth [29], and show that stochastic depth consistently improves transformer training for segmentation. CNN models rely on batch nor-

malization [31] which also acts as a regularizer. In contrast, transformers are usually composed of layer normalization [2] combined with dropout as a regularizer during training [15, 17]. Dropout randomly ignores tokens given as input of a block and stochastic depth randomly skips a learnable block of the model during the forward pass. We compare regularizations on Seg-B and Seg-B<sup>†</sup> based on ViT-B and DeiT-B backbones respectively. Table 2 shows that stochastic depth improves the performance independently of the backbone, with +0.48% for ViT-B and +1.20% for DeiT-B compared to the baseline without regularization. Dropout consistently hurts performances, either alone or when combined with stochastic depth. The biggest drop in performance is for the DeiT backbone, -3.08% compared to the baseline, which may come from the fact that DeiT has been trained on ImageNet without dropout. This is consistent with [46] which observed the negative impact of dropout for image classification. From now on, all the models will be trained with stochastic depth and without dropout.

**Transformer size.** We now study the impact of transformers size on performance by varying the number of layers and the tokens size for a fixed patch size of 16. Table 3 shows that performance scales nicely with the backbone capacity. When doubling the token dimension with the DeiT backbone, from Seg-S<sup>†</sup>/16 to Seg-B<sup>†</sup>/16, we get a 4.70% improvement. When doubling the number of layers with the ViT backbone, from Seg-B/16 to Seg-L/16, we get an improvement of 2.91%. We observe that while ViT-B and DeiT-B have the same capacity, DeiT-B performs better, the difference in mIoU is similar to the ImageNet top-1 accuracy which highlights the benefit of strong data-augmentation pre-training used in DeiT. Finally, our largest Segmenter model, Seg-L/16, achieves a *state-of-the-art* mIoU of 48.60 with a simple decoding scheme on the ADE20 validation dataset. The absence of tasks-specific layers vastly used in FCN models suggests that transformer based methods provide more expressive models, well suited for semantic segmentation.

**Patch size.** Representing an image with a patch sequence provides a simple way to trade-off between speed and accuracy by varying the patch size. While increasing the patch size results in a coarser representation of the image, it results in a smaller sequence that is faster to process. The third and fourth parts of Table 3 report the performance for ViT backbones and varying patch sizes. We observe that the patch size is a key factor for semantic segmentation performance. It is similarly important to the model size. Indeed, going from a patch size 32 to 16 we observe an improvement of 4.77% for Seg-B and 5.96% for Seg-L. For Seg-B, we also report results for a patch size of 8 and report an mIoU of 48.06%, reducing the gap from ViT-B to ViT-L to 0.54% while requiring substantially fewer parameters. This trend shows that reducing the patch size is a robust source

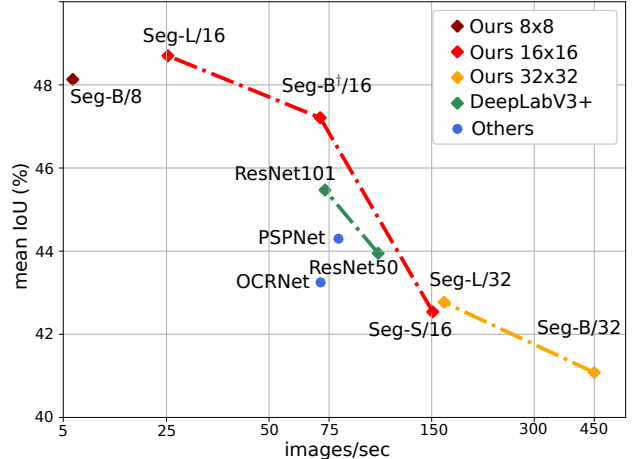


Figure 4: Images per second and mean IoU for our approach compared to other methods on ADE20K validation set. Segmenter models offer a competitive trade-off in terms of performance and precision.

Method	Small	Medium	Large
DeepLab V3+	36.88	50.33	49.19
Seg-B/32	27.15	44.98	50.63
Seg-B/16	35.12	49.78	52.46
Seg-B <sup>†</sup> /16	<b>38.31</b>	50.91	52.08
Seg-B/8	37.34	52.33	54.51
Seg-L/32	29.62	46.17	52.14
Seg-L/16	37.84	<b>52.66</b>	<b>55.25</b>
Seg-B-Mask/32	29.80	45.32	49.71
Seg-B <sup>†</sup> -Mask/16	<b>40.0</b>	50.28	52.75
Seg-L-Mask/16	38.71	<b>52.75</b>	<b>55.32</b>

Table 4: Evaluation with respect to the object size on ADE20k validation set (mean IoU). (Top) DeepLab v3+, (middle) different encoders with a linear decoder and (bottom) different encoders with a mask transformer as decoder.

Dataset Size	4k	8k	12k	16k	20k
mIoU	39.71	42.80	44.42	45.88	47.10

Table 5: Performance comparison of Seg-B<sup>†</sup>/16 models trained with increasing dataset size and evaluated on ADE20K validation set.

of improvement which does not introduce any parameters but requires to compute attention over longer sequences, increasing the compute time and memory footprint. If it was computationally feasible, ViT-L/8 would probably be the best performing model. Going towards more computation and memory efficient transformers handling larger sequence of smaller patches is a promising direction.

To further study the impact of patch size, we show segmentation maps generated by Segmenter models with de-

creasing patch size in Figure 3. We observe that for a patch size of 32, the model learns a globally meaningful segmentation but produces poor boundaries, for example the two persons on the left are predicted by a single blob. Reducing the patch size leads to considerably sharper boundaries as can be observed when looking at the contours of persons. Hard to segment instances as the thin streetlight pole in the background are only captured at a resolution of 8. In Table 4, we report mean IoU with respect to object size and compare Segmenter to DeepLab V3+. To reproduce DeepLabV3+ results, we used models from the MMSegmentation library [12]. We observe that overall DeepLab performs well on small instances. Segmenter outperforms DeepLab when using large transformer models or a patch size of 8.

**Transformer versus FCN.** Table 3 and Table 4 compares our approach to DeepLabV3+ [9] with ResNet-D backbone [25], one of the best fully-convolutional approaches. Our transformer approach provides a significant improvement over this state-of-the-art convolutional approach, highlighting the ability of transformers to capture global scene understanding. On large instances Segmenter consistently outperforms DeepLab with an improvement of 5% for Seg-L/16. However, DeepLab performs well on small instances and outperforms Seg-B/16 while having a similar number of parameters. Seg-B/8 and Seg-L/16 perform best on small instances though at higher computational cost.

**Performance.** In Figure 4, we compare our models to several *state-of-the-art* methods in terms of images per seconds and mIoU and show a clear advantage of Segmenter over FCN based models. To compute the images per second, we use a 32Go V100 GPU, fix the image resolution to 512 and for each model we maximize the batch size allowed by memory for a fair comparison. We observe that Seg/16 models perform best in terms of performance versus accuracy with Seg-B<sup>†</sup>/16 offering the best trade-off and outperforming FCN based approaches with similar inference speed. Seg/32 models learn coarser segmentation maps as discussed in the previous section and enable fast inference with 450 images per second for Seg-B/32.

**Dataset size.** Vision Transformers highlighted the importance of large datasets to attain good performance for the task of image classification. At the scale of a semantic segmentation dataset, we analyze Seg-B<sup>†</sup>/16 performance on ADE20k dataset in Table 5 when trained with a dataset of increasing size. We observe an important drop in performance when the training set size is below 8k images. This shows that even during fine-tuning transformers performs best with a sufficient amount of data.

**Decoder variants.** In this section, we compare different decoder variants. We evaluate the mask transformer introduced in Section B and compare it to the linear baseline. The mask transformer has 2 layers with the same token

Model	Decoder	mIoU
Seg-B/32	Linear	40.92
Seg-B-Mask/32	Mask Transformer	<b>42.14</b>
Seg-B <sup>†</sup> /16	Linear	47.10
Seg-B <sup>†</sup> -Mask/16	Mask Transformer	<b>47.67</b>
Seg-L/16	Linear	48.60
Seg-L-Mask/16	Mask Transformer	<b>48.82</b>

Table 6: Performance comparison of different decoders on ADE20K validation set.

Method	Backbone	mIoU	Pix. Acc.
PSANet [61]	ResNet-101	43.77	81.51
CFNet [58]	ResNet-101	44.89	-
PSPNet [60]	ResNet-269	44.94	81.69
CCNet [30]	ResNet-101	45.22	-
ANN [64]	ResNet-101	45.24	-
APCNet [24]	ResNet-101	45.38	-
DeepLabv3+ [9]	Xception-65	45.65	-
OCR [56]	HRNetV2-W48	45.66	-
ACNet [22]	ResNet-101	45.90	81.96
DNL [53]	ResNet-101	45.97	-
DRANet [20]	ResNet-101	46.18	81.91
CPNet [54]	ResNet-101	46.27	81.85
Seg-L-Mask/16 (SS)	ViT-L	48.82	82.77
Seg-L-Mask/16 (MS)	ViT-L	<b>50.77</b>	<b>83.72</b>

Table 7: State-of-the-art comparison on ADE20K validation set.

and hidden size as the encoder. Table 6 reports the performance. The mask transformer provides consistent improvements over the linear baseline. The most significant gain of +1.2% is obtained for Seg-B/32. This shows that non-linear decoding is more important for relatively large patch sizes. For Seg-B<sup>†</sup>-Mask/16 we obtain a +0.5% improvement and for Seg-L/16 a gain of +0.2%. In Table 4 we examine the gain of different models with respect to object size. In all cases the gain is more significant for small objects, i.e., +2.5% for Seg-B/32, +1.7% Seg-B<sup>†</sup>-Mask/16 and +1.4% for Seg-L/16. Again the difference is most pronounced for Seg-B/32. The class embedding learned by the mask transformer are semantically meaningful, i.e., similar classes are nearby, see Figure 8 for more details.

#### 4.4. Comparison with state of the art

In this section, we compare the performance of Segmenter with respect to *state-of-the-art* methods on ADE20K, Cityscapes and Pascal Context datasets. **ADE20K.** With single-scale inference only, Seg-L-Mask/16 outperforms the state of the art as shown in Table 7, achieving a 48.82% mIoU and 82.77% pixel accuracy. With multi-scale inference, Seg-L-Mask/16 attains 50.77%, outperforming by a margin of +4.50%

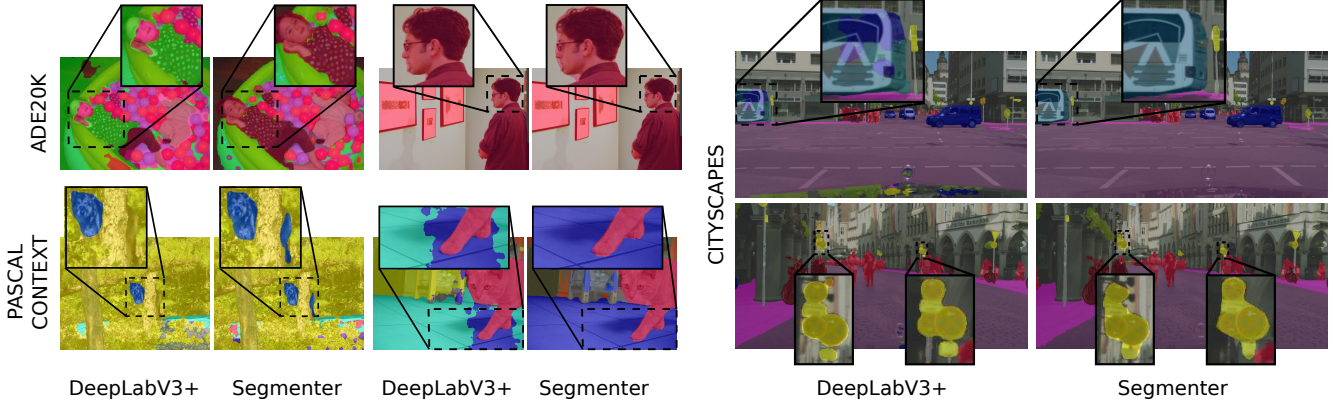


Figure 5: Qualitative comparison of Seg-L-Mask/16 performance with DeepLabV3+. See Section D for additional qualitative results.

mIoU the *state-of-the-art* method CPNet [54].

**Cityscapes.** Table 8 reports the performance of Segmenter on Cityscapes. For Cityscapes, we report results using Seg-B<sup>+</sup>-Mask/16 as larger variants did not fit into memory due to the large image resolution of 768×768. Seg-B<sup>+</sup>-Mask/16 with multi-scale is on-par with other *state-of-the-art* methods.

**Pascal Context** Table 9 reports the performance on Pascal Context. Seg-L-Mask/16 with multi-scale inference achieves 55.6% mIoU and is improving over many *state-of-the-art* methods. In particular, we report an improvement of 7% over DeepLabV3+.

**Discussion.** Segmenter outperforms the state of the art on ADE20K, the largest dataset in this paper (20k training images). It is on-par with the state of the art on Cityscapes and Pascal Context (3k and 5k training images). This is consistent with the results presented in Table 5 where we observe a drop in performance when the training dataset is less than 8k images. Such results suggest that even when fine-tuning on down-stream tasks vision transformers perform better for relatively large datasets.

**Qualitative results.** Figure 5 shows a qualitative comparison of Segmenter and DeepLabV3+, for which models were provided by the MMSegmentation [12] library. We can observe that DeeplabV3+ tends to generate sharper object boundaries while Segmenter provides more consistent labels on large instances and handles partial occlusions better.

## 5. Conclusion

This paper introduces a pure transformer approach for semantic segmentation. The encoding part builds up on the recent Vision Transformer (ViT), but differs in that we rely on the encoding of all images patches. We observe that the transformer captures the global context very well. Applying a simple point-wise linear decoder to the patch encodings already achieves excellent results. Decoding with a mask transformer further improves the performance. We believe

Method	Backbone	mIoU
PSANet [61]	ResNet-101	79.1
DeepLabV3+ [9]	Xception-71	79.6
ANN [64]	ResNet-101	79.9
MDEQ [4]	MDEQ	80.3
DNL [53]	ResNet-101	80.5
CCNet [30]	ResNet-101	81.3
Panoptic-Deeplab [11]	Xception-71	<b>81.5</b>
Seg-B <sup>+</sup> -Mask/16 (SS)	DeiT-B	79.0
Seg-B <sup>+</sup> -Mask/16 (MS)	DeiT-B	<b>80.7</b>

Table 8: State-of-the-art comparison on Cityscapes validation set.

Method	Backbone	mIoU
DeepLabv3+ [9]	ResNet-101	48.5
DANet [21]	ResNet-101	52.6
ANN [64]	ResNet101	52.8
CPNet [54]	ResNet-101	53.9
CFNet [58]	ResNet-101	54.0
ACNet [22]	ResNet-101	54.1
APCNet [24]	ResNet101	54.7
DNL [53]	HRNetV2-W48	55.3
DRANet [20]	ResNet-101	55.4
OCR [56]	HRNetV2-W48	<b>56.2</b>
Seg-L-Mask/16 (SS)	ViT-L	54.6
Seg-L-Mask/16 (MS)	ViT-L	<b>55.6</b>

Table 9: State-of-the-art comparison on Pascal Context validation set.

that our end-to-end encoder-decoder transformer is a first step towards a unified approach for semantic segmentation, instance segmentation and panoptic segmentation.



## 6. Acknowledgements

We thank Andreas Steiner for providing the ViT-Base model trained on  $8 \times 8$  patches. This work was partially supported by the HPC resources from GENCI-IDRIS (Grant 2020-AD011011163R1), Louis Vuitton ENS Chair on Artificial Intelligence, and the French government under management of Agence Nationale de la Recherche as part of the "Investissements d'avenir" program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute).

## References

- [1] Md Amirul Islam, Mrigank Rochan, Neil D. B. Bruce, and Yang Wang. Gated feedback refinement network for dense image labeling. In *CVPR*, 2017. 2
- [2] Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *arXiv preprint*, 2016. 6
- [3] V. Badrinarayanan, A. Kendall, and R. Cipolla. SegNet: A deep convolutional encoder-decoder architecture for image segmentation. *TPAMI*, 2017. 2
- [4] Shaojie Bai, Vladlen Koltun, and J. Zico Kolter. Multiscale deep equilibrium models. In *NeurIPS*, 2020. 8
- [5] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? *arXiv preprint*, 2021. 2
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with Transformers. In *ECCV*, 2020. 2
- [7] 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. *TPAMI*, 2018. 2
- [8] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint*, 2017. 2
- [9] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018. 1, 2, 5, 7, 8
- [10] Mark Chen, Alec Radford, Rewon Child, Jeff Wu, Heewoo Jun, Prafulla Dhariwal, David Luan, and Ilya Sutskever. Generative pretraining from pixels. *PLMR*, 2020. 2
- [11] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-Deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *CVPR*, 2020. 8
- [12] MMSegmentation Contributors. MMSegmentation: OpenMMLab semantic segmentation toolbox and benchmark. <https://github.com/open-mmlab/mms Segmentation>, 2020. 5, 7, 8
- [13] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 4
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 3
- [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. 6
- [16] Carl Doersch, Ankush Gupta, and Andrew Zisserman. CrossTransformers: spatially-aware few-shot transfer. In *NeurIPS*, 2020. 2
- [17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth  $16 \times 16$  words: Transformers for image recognition at scale. In *ICLR*, 2021. 2, 4, 5, 6
- [18] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. Vision transformer models. [https://console.cloud.google.com/storage/browser/vit\\_models](https://console.cloud.google.com/storage/browser/vit_models), 2021. 5
- [19] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *TPAMI*, 2013. 2
- [20] J. Fu, J. Liu, J. Jiang, Y. Li, Y. Bao, and H. Lu. Scene segmentation with Dual Relation-Aware attention Network. *TNNLS*, 2020. 1, 2, 7, 8
- [21] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual Attention Network for scene segmentation. In *CVPR*, 2019. 1, 2, 8
- [22] Jun Fu, Jing Liu, Yuhang Wang, Yong Li, Yongjun Bao, Jinhui Tang, and Hanqing Lu. Adaptive Context Network for scene parsing. *ICCV*, 2019. 7, 8
- [23] Boris Hanin and David Rolnick. How to start training: The effect of initialization and architecture. In *NeurIPS*, 2018. 5
- [24] Junjun He, Zhongying Deng, Lei Zhou, Yali Wang, and Yu Qiao. Adaptive Pyramid Context Network for semantic segmentation. In *CVPR*, 2019. 7, 8
- [25] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *CVPR*, 2019. 7
- [26] Mohammad Hesam Hesamian, Wenjing Jia, Xiangjian He, and Paul Kennedy. Deep learning techniques for medical image segmentation: Achievements and challenges. *JDI*, 2019. 1
- [27] Zhang-Wei Hong, Chen Yu-Ming, Shih-Yang Su, Tzu-Yun Shann, Yi-Hsiang Chang, Hsuan-Kung Yang, Brian Hsi-Lin Ho, Chih-Chieh Tu, Yueh-Chuan Chang, Tsu-Ching Hsiao, Hsin-Wei Hsiao, Sih-Pin Lai, and Chun-Yi Lee. Virtual-to-real: Learning to control in visual semantic segmentation. *IJCAI*, 2018. 1
- [28] Chia-Wei Hsieh, Chieh-Yun Chen, Chien-Lung Chou, Hong-Han Shuai, Jiaying Liu, and Wen-Huang Cheng. Fashionon: Semantic-guided image-based virtual try-on with detailed human and clothing information. In *ACM MM*, 2019. 1
- [29] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q. Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016. 5
- [30] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. CCNet: Criss-Cross attention for semantic segmentation. In *ICCV*, 2019. 1, 7, 8

- [31] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 6
- [32] Chen Liang-Chieh, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan Yuille. Semantic image segmentation with deep convolutional nets and fully connected CRFs. In *ICLR*, 2015. 2, 5
- [33] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. RefineNet: Multi-path refinement networks for high-resolution semantic segmentation. In *CVPR*, 2017. 2
- [34] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully Convolutional Networks for semantic segmentation. In *CVPR*, 2015. 2
- [35] Shervin Minaee, Yuri Boykov, F. Porikli, Antonio J. Plaza, N. Kehtarnavaz, and Demetri Terzopoulos. Image segmentation using deep learning: A survey. *TPAMI*, 2021. 1
- [36] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan L. Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, 2014. 4
- [37] F. Ning, D. Delhomme, Yann LeCun, F. Piano, Léon Bottou, and Paolo Emilio Barbano. Toward automatic phenotyping of developing embryos from videos. *TIP*, 2005. 2
- [38] Pedro Pinheiro and Ronan Collobert. Recurrent Convolutional Neural Networks for scene labeling. In *ICML*, 2014. 2
- [39] Tobias Pohlen, Alexander Hermans, Markus Mathias, and Bastian Leibe. Full-resolution residual networks for semantic segmentation in street scenes. In *CVPR*, 2017. 2
- [40] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 1951. 5
- [41] O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 2
- [42] Evan Shelhamer, Jonathan Long, and Trevor Darrell. Fully Convolutional Networks for semantic segmentation. *TPAMI*, 2017. 1
- [43] M. Siam, S. Elkerdawy, M. Jagersand, and S. Yogamani. Deep semantic segmentation for automated driving: Taxonomy, roadmap and challenges. In *ITSC*, 2017. 1
- [44] Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *JMLR*, 2014. 5
- [45] Farhana Sultana, Abu Sufian, and Paramartha Dutta. Evolution of image segmentation using deep convolutional neural network: A survey. *Knowledge-Based Systems*, 2020. 1
- [46] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training Data-Efficient image Transformers and distillation through attention. *arXiv preprint*, 2020. 2, 3, 4, 6
- [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2, 3
- [48] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan L. Yuille, and Liang-Chieh Chen. Max-deeplab: End-to-end panoptic segmentation with mask transformers. *arXiv preprint*, 2020. 2, 4
- [49] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-Deeplab: Stand-alone axial-attention for panoptic segmentation. In *ECCV*, 2020. 2
- [50] Xiaolong Wang, Ross B. Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018. 2
- [51] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. In *NeurIPS*, 2020. 4
- [52] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2020. 5
- [53] Minghao Yin, Zhuliang Yao, Yue Cao, Xiu Li, Zheng Zhang, Stephen Lin, and Han Hu. Disentangled Non-local Neural Networks. In *ECCV*, 2020. 1, 2, 7, 8
- [54] Changqian Yu, Jingbo Wang, Changxin Gao, Gang Yu, Chunhua Shen, and Nong Sang. Context Prior for scene segmentation. In *CVPR*, 2020. 1, 2, 7, 8
- [55] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016. 2
- [56] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-Contextual Representations for semantic segmentation. In *ECCV*, 2020. 1, 7, 8
- [57] Yuhui Yuan and Jingdong Wang. OCNet: Object Context Network for scene parsing. *arXiv preprint*, 2018. 2
- [58] H. Zhang, H. Zhang, C. Wang, and J. Xie. Co-occurrent features in semantic segmentation. In *CVPR*, 2019. 1, 7, 8
- [59] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip Torr, and Vladlen Koltun. Point Transformer. *arXiv preprint*, 2020. 2
- [60] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid Scene Parsing Network. In *CVPR*, 2017. 1, 2, 7
- [61] Hengshuang Zhao, Yi Zhang, Shu Liu, Jianping Shi, Chen Change Loy, Dahua Lin, and Jiaya Jia. PSANet: Point-wise Spatial Attention Network for scene parsing. In *ECCV*, September 2018. 1, 2, 7, 8
- [62] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip H.S. Torr, and Li Zhang. Rethinking semantic segmentation from a sequence-to-sequence perspective with Transformers. *arXiv preprint*, 2020. 2
- [63] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ADE20K dataset. *IJCV*, 2019. 2, 4
- [64] Zhen Zhu, Mengde Xu, Song Bai, Tengpeng Huang, and Xiang Bai. Asymmetric non-local neural networks for semantic segmentation. In *ICCV*, 2019. 7, 8

## Appendix

This appendix presents additional results. We study the impact of ImageNet pretraining on the performance and demonstrate its importance in Section A. We then ablate the parameters of our mask transformer and compare to alternative convolutional decoders in Section B. To gain more insight about our approach Segmter, we analyze its attention maps and the learned class embeddings in Section C. Finally, we give additional qualitative comparison of Segmter to DeepLabV3+ on ADE20K, Cityscapes and Pascal Context in Section D.

### A. ImageNet pre-training

To study the impact of ImageNet pre-training on Segmter, we compare our model pre-trained on ImageNet with equivalent models trained from scratch. To train from scratch, the weights of the model are initialized randomly with a truncated normal distribution. We use a base learning rate of  $10^{-3}$  and two training procedures. First, we follow the fine-tuning procedure and use SGD optimizer with "poly" scheduler. Second, we follow a more standard procedure when training a transformer from scratch where we use AdamW with a cosine scheduler and a linear warmup for 16K iterations corresponding to 10% of the total number of iterations. Table 10 reports results for Seg-B<sup>†</sup>/16. We observe that when pre-trained on ImageNet-21k, Seg-B<sup>†</sup>/16 reaches 47.21% yielding a 36% improvement over the best randomly initialized model.

Method	Pre-training	Optimizer	mIoU
Seg-B <sup>†</sup> /16	None	AdamW	6.64
Seg-B <sup>†</sup> /16	None	SGD	11.18
Seg-B <sup>†</sup> /16	ImageNet-21k	SGD	<b>47.21</b>

Table 10: Impact of pretraining on the performance on ADE20K validation set.

### B. Decoder variants

For the decoder ablation, we use Seg-B<sup>†</sup>/16. We compare our transformer-based decoder to convolution-based decoders. We also ablate the mask transformer capacity by varying the number of layers and the token dimension.

**CNN decoder.** The CNN baseline comprises a 1x1 convolution followed by a block of 3x3 convolutions, all of the convolutions have 256 units followed by BatchNorm

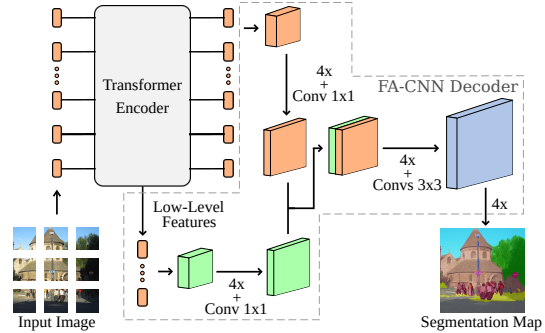


Figure 6: Overview of the FA-CNN decoder.

and ReLU. The encoder tokens are reshaped into a 2D image and passed to the decoder. The first convolution is applied, followed by 4x upsampling and the second convolution block. The segmentation map is then upsampled to the original image size.

**FA-CNN decoder.** The Features Aggregation CNN (FA-CNN) decoder is similar to the CNN decoder except that it includes features from a low transformer layer during decoding. We provide an overview of this decoder in Figure 6. After the first convolution, features from an intermediate layer of the encoder are concatenated to the decoder features. This decoder is identical to DeepLab V3+.

Decoder	Layers	Token dim	Upsampling	# Params	mIoU
Linear	0	-	-	0.1M	47.21
CNN	2	256	4	1.4M	<b>47.51</b>
FA-CNN	2	256	4	1.6M	47.07
Mask	1	768	2	9M	47.64
Mask	2	768	2	16M	<b>47.77</b>
Mask	3	768	2	23M	47.58
Mask	2	384	2	4.2M	<b>47.38</b>

Table 11: Ablation of the decoder on ADE20K validation set.

Results in Table 11 show that our mask transformer slightly outperforms the baseline as well as the two CNN-based decoders. The CNN decoder works better than the FA-CNN decoder and the linear baseline. Our mask transformer outperforms the best CNN decoder by +0.2%.

We also observe that our mask transformer obtains best results with 2 layers and 768 tokens. Reducing the token dimensions from 768 to 384 decreases performance. In summary, our new mask transformer decoder demonstrates competitive performance and is more general than its convolutional counterpart.

### C. Attention maps and class embeddings

To better understand how our approach Segmenter processes images, we display attention maps of Seg-B<sup>†</sup>-Mask/16 for 3 images in Figure 7. We resize attention maps to the original image size, apply a Gaussian filter matching the model input patch size and normalize the values in  $[0, 1]$ . For each image, we analyze attention maps of a patch on a small instance, for example lamp, cow or car. We also analyze attention maps of a patch on a large instance, for example bed, grass and road. We observe that the attention map field-of-view adapts to the input image and the instance size, gathering global information on large instances and focusing on local information on smaller instances. This adaptability is typically not possible with CNN which have a constant field-of-view, independently of the data. We also note there is progressive gathering of information from bottom to top layers, as for example on the cow instance, where the model first identifies the cow the patch belongs to, then identifies other cow instances. We observe that attention maps of lower layers depends strongly on the selected patch while they tend to be more similar for higher layers.

To gain some understanding of the class embeddings learned with the mask transformer, we project embeddings into 2D with a singular value decomposition. Figure 8 shows that these projections group instances such as means of transportation (bottom left), objects in a house (top) and outdoor categories (middle right). It displays an implicit clustering of semantically related categories.

### D. Qualitative Results

We present additional qualitative results including comparison with DeepLabV3+ and failure cases in Figures 9, 10 and 11. We can see in Figure 9 that Segmenter produces more coherent segmentation maps than DeepLabV3+. This is the case for the airplane signalmen’s helmet in Figure 9(b) or the car which is occluded by people in Figure 9(d). In Figure 11(a), we can see how DeepLabV3+ handles better the boundaries between different people entities. In Figure 10, we show how for some examples, different segments which look very similar are confused both in DeepLabV3+ and Segmenter. For example, the cushions and pillows in column (a), grass and ground in (b) and the gears of the pilot and his motorbike texture in (c). Finally, both Segmenter and DeepLabV3+ have problems segmenting small instances such as lamp, people or flowers in Figure 12 (a) or the cars and signals in Figure 12 (b).



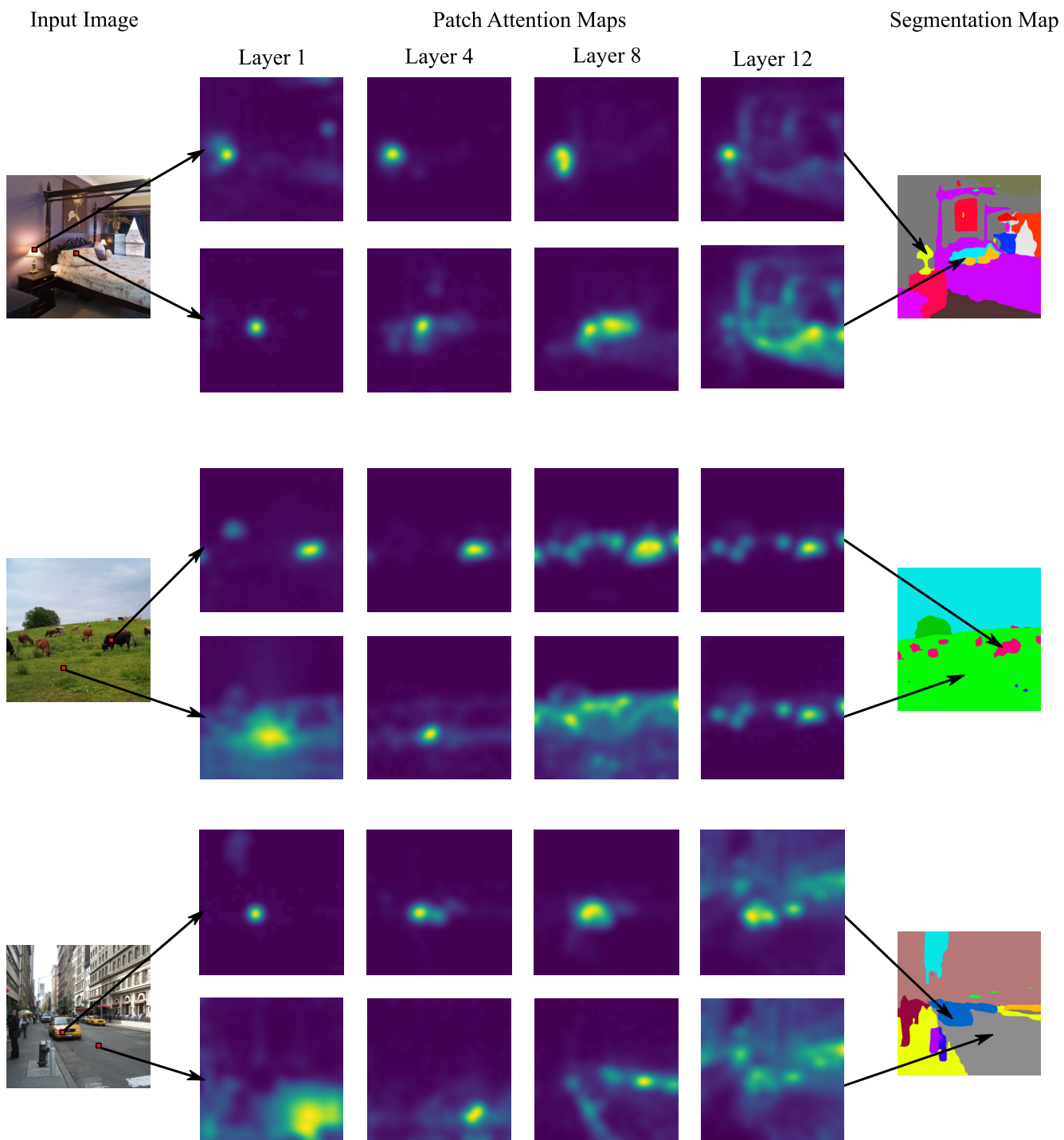


Figure 7: Seg-B<sup>†</sup>-Mask/16 patch attention maps for the layers 1, 4, 8 and 12.

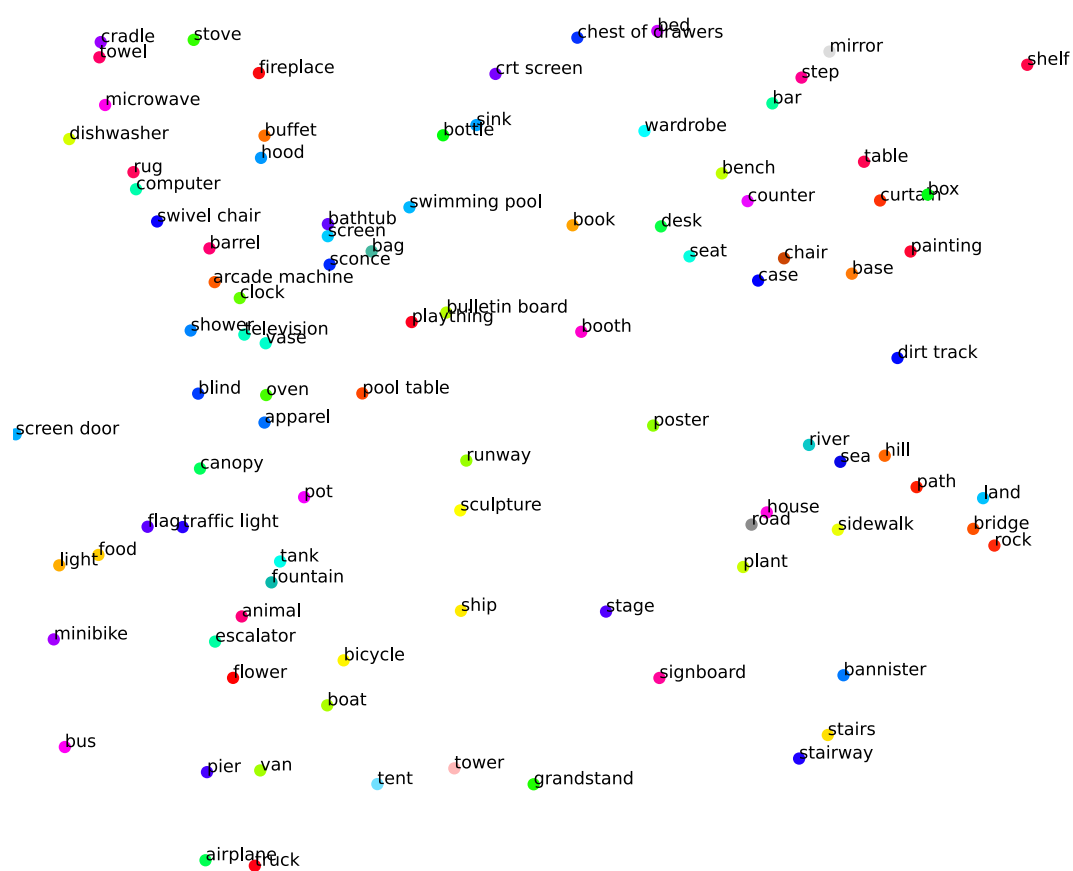


Figure 8: Singular value decomposition of the class embeddings learned with the mask transformer on ADE20K.

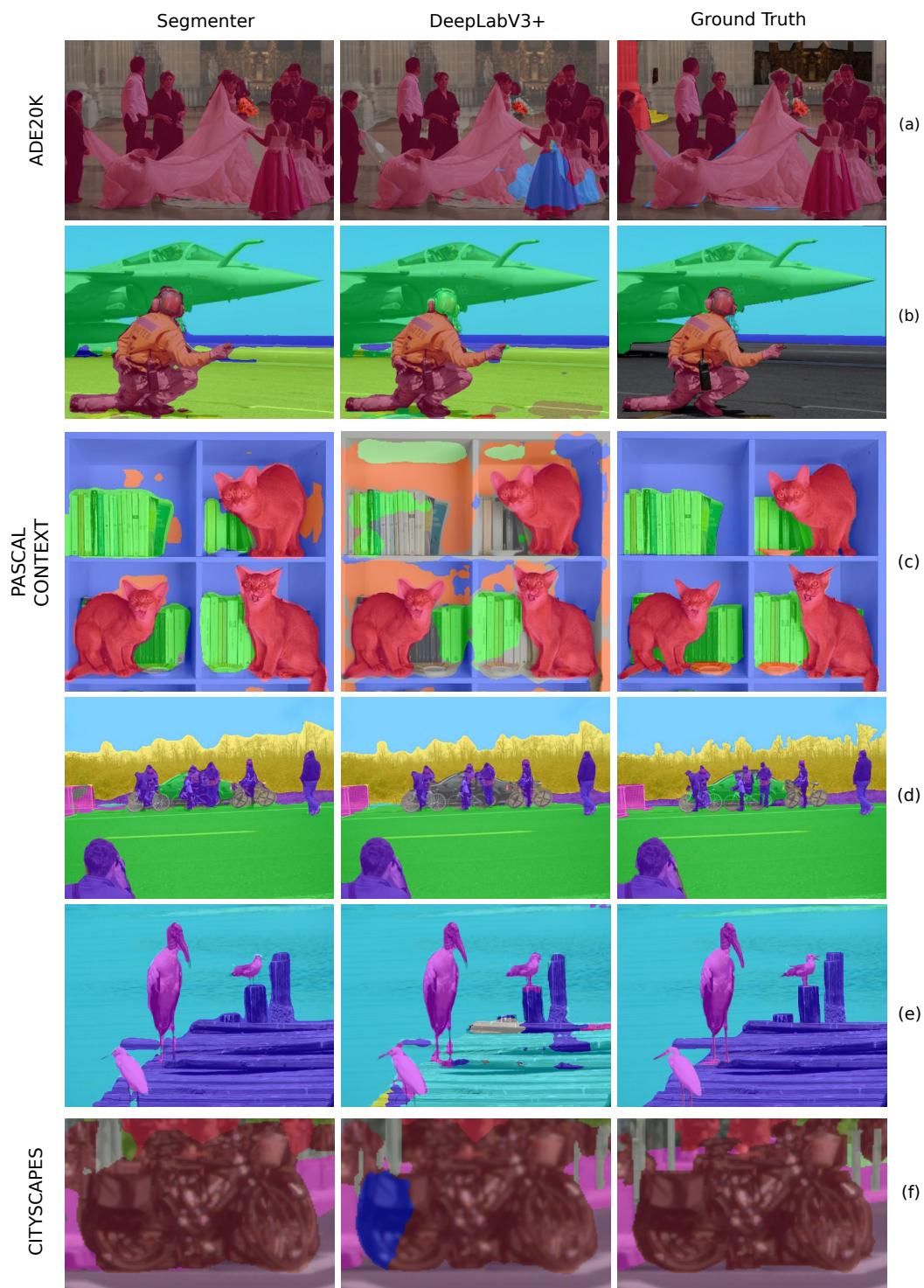


Figure 9: Segmentation maps where Seg-L-Mask/16 produces more coherent segmentation maps than DeepLabv3+.

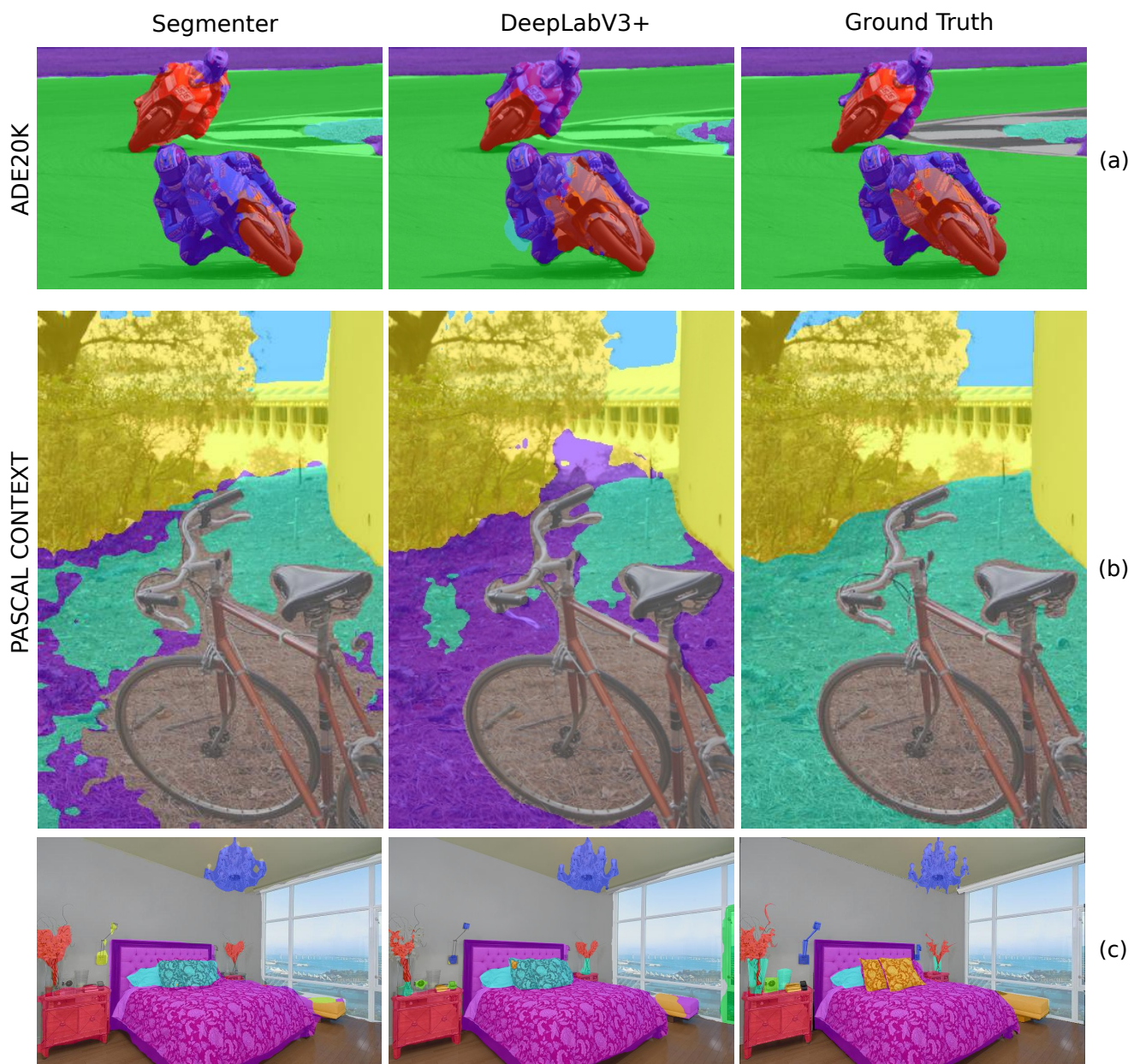


Figure 10: Example for Seg-L-Mask/16 and DeepLabV3+, where elements which look very similar are confused.



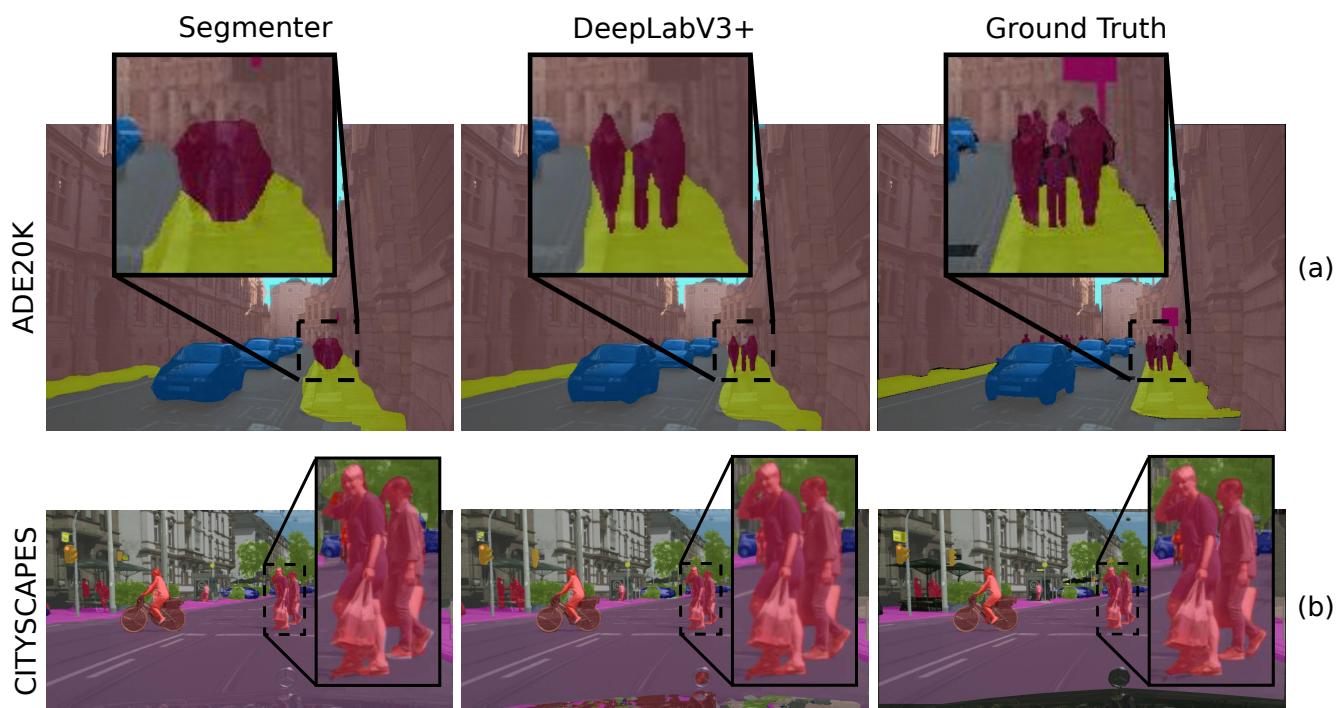


Figure 11: Comparison of Seg-L-Mask/16 and DeepLabV3+ for images with near-by persons. We can observe that DeepLabV3+ localizes boundaries better.

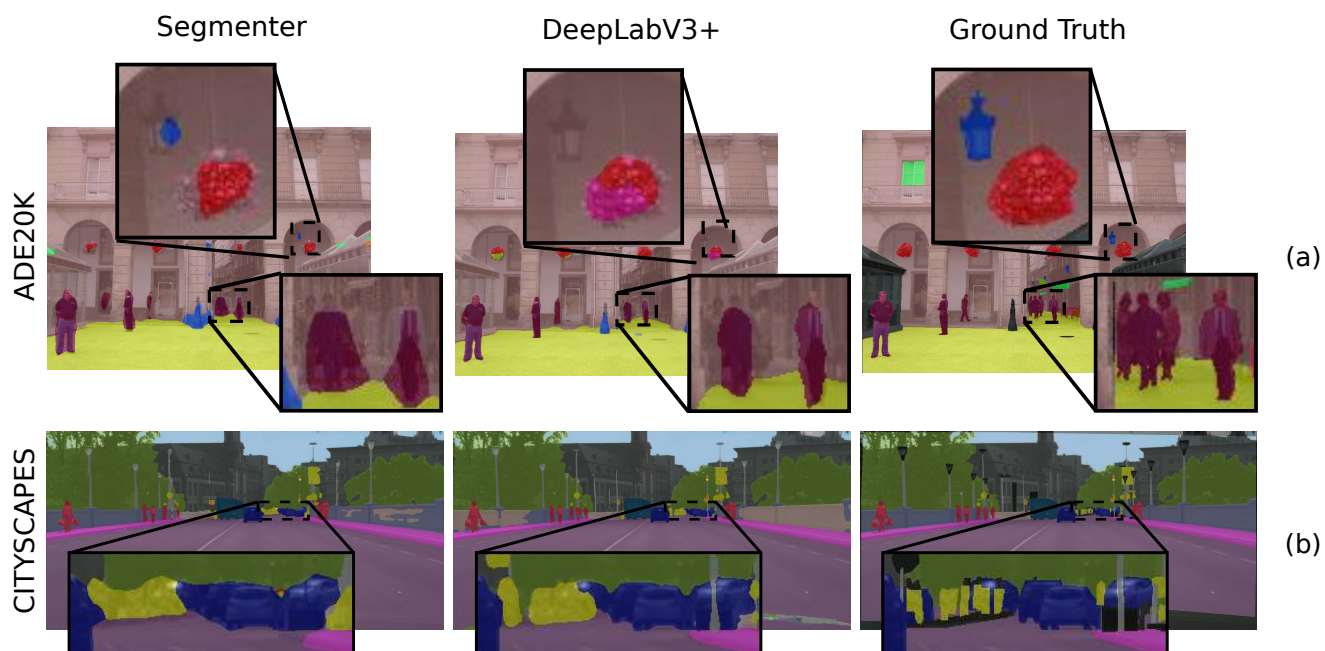


Figure 12: Failure cases of DeepLabV3+ and Seg-L-Mask/16, for small instances such as (a) lamp, people, flowers and (b) cars, signals.