

Superpixel Transformers for Efficient Semantic Segmentation

Alex Zihao Zhu^{1*}, Jieru Mei^{2*}, Siyuan Qiao³, Hang Yan¹,
Yukun Zhu³, Liang-Chieh Chen⁴, Henrik Kretzschmar⁵

Abstract— Semantic segmentation, which aims to classify every pixel in an image, is a key task in machine perception, with many applications across robotics and autonomous driving. Due to the high dimensionality of this task, most existing approaches use local operations, such as convolutions, to generate per-pixel features. However, these methods are typically unable to effectively leverage global context information due to the high computational costs of operating on a dense image. In this work, we propose a solution to this issue by leveraging the idea of superpixels, an over-segmentation of the image, and applying them with a modern transformer framework. In particular, our model learns to decompose the pixel space into a spatially low dimensional superpixel space via a series of local cross-attentions. We then apply multi-head self-attention to the superpixels to enrich the superpixel features with global context and then directly produce a class prediction for each superpixel. Finally, we directly project the superpixel class predictions back into the pixel space using the associations between the superpixels and the image pixel features. Reasoning in the superpixel space allows our method to be substantially more computationally efficient compared to convolution-based decoder methods. Yet, our method achieves state-of-the-art performance in semantic segmentation due to the rich superpixel features generated by the global self-attention mechanism. Our experiments on Cityscapes and ADE20K demonstrate that our method matches the state of the art in terms of accuracy, while outperforming in terms of model parameters and latency.

I. INTRODUCTION

The problem of semantic segmentation, or classifying every pixel in the image, is increasingly common in many robotics applications. A dense, fine-grained, understanding of the world is necessary for navigation in cluttered environments, particularly for applications such as autonomous driving, where scene understanding is deeply safety-critical. On the other hand, many robotics systems are combinations of highly complex and specialized systems, and latency is an ever-present issue for real time operation.

The balance between safety and performance and latency is critical for modern robotic systems. While the state of the art in semantic segmentation is able to achieve strong performance in terms of metrics such as mean Intersection over Union (mIoU), many methods still rely on dense decoders which produce predictions for every pixel in the scene. As a result, these methods tend to be relatively expensive, and arguably produce a lot of redundant computation for nearby pixels that are often very similar.

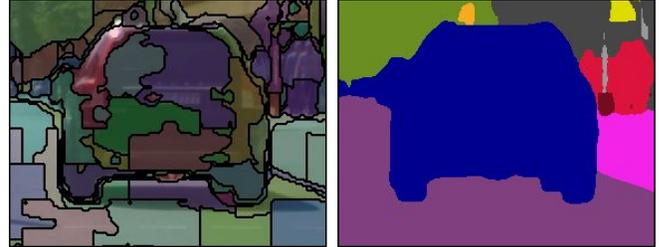


Fig. 1. Our method enables efficient segmentation of high-resolution camera images by learning to decompose the images into a set of superpixels. Specifically, we oversegment the image pixels into a small set of soft superpixels (left) via a series of local cross attentions. The superpixels are then refined via a set of multi-head self attentions, and directly classified. Finally, we fuse the class predictions with the superpixel-pixel associations to produce a dense semantic segmentation (right).

To address this issue, we aim to bring the classical ideas surrounding superpixels into modern deep learning. The premise of using superpixels is to decompose and over-segment the image into a series of irregular patches. By grouping similar pixels into superpixels and then operating on the superpixel level, one can significantly reduce the computational cost of dense prediction tasks, such as semantic segmentation. Classical superpixel algorithms, such as SLIC [1], however, rely on hard associations between each image pixel and superpixel. This makes it hard to embed the superpixel representation into neural network architectures [29] as this association is not differentiable for back propagation. Recent works, such as superpixel sampling networks (SSN) [24], resolve this issue by turning the hard association into a soft one. While this is a step towards incorporating superpixels into neural networks, their segmentation quality still lags behind other models that adopt the per-pixel or per-mask representation.

In this work, we propose a novel architecture that aims to revive the differentiable superpixel generation pipeline in a modern transformer framework [40]. In place of the iterative clustering algorithms used in SLIC [1] and SSN [24], we propose to learn the superpixel representation by developing a series of *local* cross-attentions between a set of learned superpixel queries and pixel features. The outputs of cross-attention modules act as the superpixel features, directly used for semantic segmentation prediction. As a result, the proposed transformer decoder effectively converts object queries to superpixel features, enabling the model to learn the superpixel representation end-to-end.

Operating on the superpixel level provides a number of

arXiv:2309.16889v2 [cs.CV] 2 Oct 2023

¹Waymo LLC

²Johns Hopkins University (Work done as an intern at Waymo)

³Google Research

⁴ByteDance Research (Work done while at Google Research)

⁵Work done while at Waymo

*Equal contributions

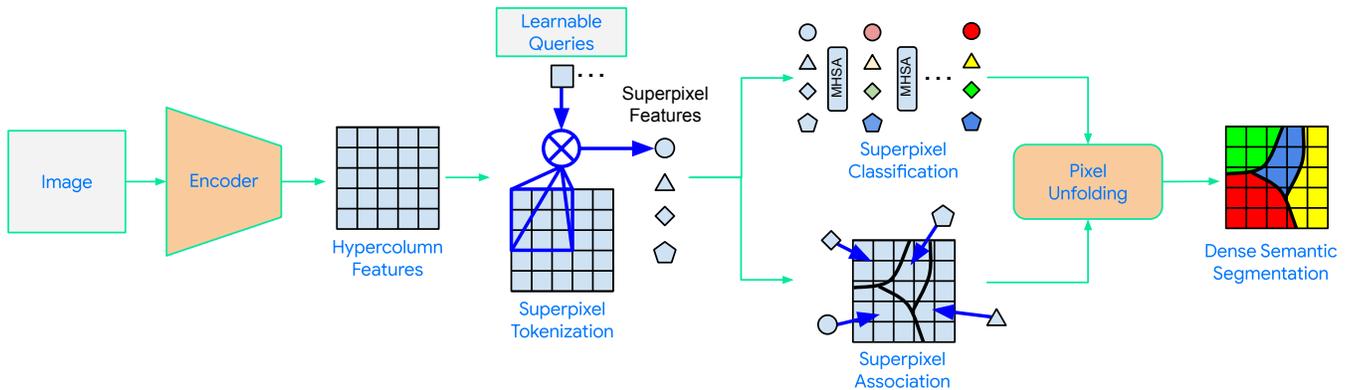


Fig. 2. Our proposed Superpixel Transformer architecture. Given an image, we first generate hypercolumn features with an off-the-shelf encoder backbone. Our superpixel tokenization module uses a series of local dual-path cross-attentions to generate features for each superpixel. The superpixel features are then enriched by several multi-head self-attention (MHSA) layers to produce a class prediction for each superpixel, while the associations between each superpixel and pixel feature are computed from their respective features. Finally, the superpixel class predictions are unfolded into the dense pixel space using the associations. Note that the figure illustrates a hard assignment between pixels and superpixels for simplicity, while in practice we apply a differentiable soft assignment.

notable benefits. Conventional pixel-based approaches are limited by the high dimensionality of the pixel space, making global self-attention computationally intractable. Numerous approaches such as axial attention [42] or window attention [31] have been developed to work around these issues by relaxing the global attention to a local one. By over-segmenting the image into a small set of superpixels, we are able to efficiently apply *global* self-attention on the superpixels, providing full global context to the superpixel features, even when reasoning about high-resolution images. Despite applying global self-attention (vs conventional convolutional neural networks) in our model, our method is more efficient than existing methods due to the low dimensionality of the superpixel space. Finally, we directly produce semantic classes for the superpixel features, and then back-project the predicted classes onto the image space using the superpixel-pixel associations.

We perform extensive evaluations on the Cityscapes [14] and ADE20K [56] datasets, where our method matches state-of-the-art performance, but at significantly lower computational cost.

In summary, the main contributions of this work are as follows:

- The first work that revives the superpixel representation in the modern transformer framework, where the object queries are used to learn superpixel features.
- A novel network architecture that uses local cross-attention to significantly reduce the spatial dimensionality of pixel features to a small set of superpixel features, enabling learning the global context between them and the direct classification of each superpixel.
- A superpixel association and unfolding scheme that projects each superpixel class prediction back to a dense pixel segmentation, discarding the CNN pixel decoder.
- Experiments on the Cityscapes and ADE20k datasets, where our method outperforms the state of the art at substantially lower computational cost.

II. RELATED WORK

Superpixels for Segmentation Before the deep learning era, the superpixel representation, paired with graphical models, was the main paradigm for image segmentation. Superpixel methods [33], [37], [13], [1] are usually used in the pre-processing step to reduce the computation cost. A shallow classifier, *e.g.*, SVM [15], predicts the semantic labels of each superpixel [18], which aggregates hand-crafted features. The graphical models, particularly conditional random fields [28], are then employed to refine the segmentation results [22], [27], [26].

ConvNets for Segmentation Convolutional neural networks (ConvNets) [29] deployed in a fully convolutional manner [34] performs semantic segmentation by pixel-wise classification. Typical ConvNet-based approaches include the DeepLab series [4], [6], [8], PSPNet [54], UPerNet [46], and OCRNet [51]. Alternatively, there are some works [19], [24] that employ superpixels to aggregate features extracted by ConvNets and show promising results.

Transformers for Segmentation Transformers [40] and their vision variants [17] have been adopted as the backbone encoders for image segmentation [7], [55], [3]. Transformer encoders can be instantiated as augmenting ConvNets with self-attention modules [43], [42]. When used as stand-alone backbones [36], [17], [31], [47], they also demonstrate strong performance compared to the previous ConvNet baselines.

Transformers are also used as the decoders [2] for image segmentation. A popular design is to generate masks embedding vectors from object queries and then multiply them with the pixel features to generate masks [39], [44]. For example, MaX-DeepLab [41] proposes an end-to-end mask transformer framework that directly predicts class-labeled object masks. Segmenter [38] and MaskFormer [12] tackle semantic segmentation from the view of mask classification. K-Net [52] generates segmentation masks by a group of learnable kernels. Inspired by the similarity between mask

transformers and clustering algorithms [33], clustering-based mask transformers are proposed to segment images [49], [48], [50]. Deformable transformer [57] is also used for improving the image segmentation as in Panoptic SegFormer [30] and Mask2Former [11].

Similar to this work, Region Proxy [53] (RegProxy) also incorporates the idea of superpixels into a deep segmentation network by using a CNN decoder to learn the association between each pixel and superpixel. However, RegProxy uses features on the regular pixel grid to represent each superpixel, and on which to apply self-attention. In comparison, we apply a set of learned weights, which correspond to each superpixel, and use cross-attention with the pixel features to compute the pixel-superpixel associations. Our experiments demonstrate that our methodology provides significant performance improvements.

In summary, all of the prior works that apply transformers to segmentation have, in some way, relied on a dense CNN decoder to generate the final dense features, and then combined these features with an attention mechanism to improve performance. Our method, in comparison, uses cross attention to reduce the image into a small set of superpixels, and only applies self-attention in this superpixel space in the decoder. This allows our method to operate on a significantly lower dimensional space (often $32^2 \times$ smaller than the image resolution), while utilizing the benefits that come with global self-attention to achieve state-of-the-art performance.

III. METHOD

Our proposed Superpixel Transformer architecture, summarized in Figure 2, consists of four main components:

- 1) *Pixel Feature Extraction*: A convolutional encoder backbone to generate hypercolumn features.
- 2) *Superpixel Tokenization*: A series of *local* dual-path cross-attentions, between a set of learned queries and pixel features, to generate a set of superpixel features.
- 3) *Superpixel Classification*: A series of multi-head self-attention layers to refine the superpixel features and produce a semantic class for each superpixel.
- 4) *Superpixel Association*: Associating the predicted superpixel classes and pixel features to obtain the final dense semantic segmentation.

We detail each component in the following subsections.

A. Pixel Feature Extraction

Typical convolutional neural networks, such as ResNet [21] and ConvNeXt [32], are employed as the encoder backbone. On top of the encoder output, we apply a multi-layer perceptron (MLP), and bilinear resize to the features after stage-1 (stride 2), stage-3 (stride 8), and stage-5 (stride 32). The multi-scale features are combined with addition to form hypercolumn features [20]. Each pixel feature is represented by their corresponding hypercolumn features, which are fed to the following Superpixel Tokenization module.

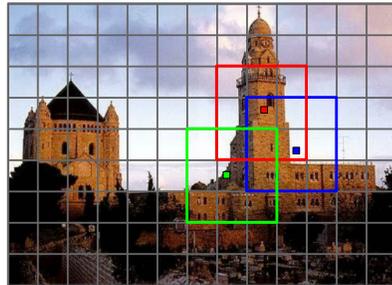


Fig. 3. Visualization of the superpixel-pixel association. Each superpixel is assigned to a gray grid cell in the image. For each pixel (small dots, size exaggerated) inside a given superpixel, we compute its cross attention with its neighboring 3×3 superpixels, highlighted with the same color. The essence of our method is that these neighborhoods overlap in a sliding window fashion.

B. Superpixel Tokenization

Before introducing our proposed Superpixel Tokenization module, we briefly review the previous works on Differentiable SLIC [1], [24], which we modernize with the transformer framework.

Preliminary: Differentiable SLIC Simple Linear Iterative Clustering (SLIC) [1] adopts the classical iterative k -means algorithm [33] to generate superpixels by clustering pixels based on their features (*e.g.*, color similarity and location proximity). Given a set of pixel features I_p and initialized superpixel features S_i^0 at iteration 0, the algorithm iterates between two steps at iteration t :

- 1) (*Hard*) *Assignment*: Compute the similarity Q_{pi}^t between each pixel feature I_p and superpixel feature S_i^t . Assign each pixel to a *single* superpixel based on its maximum similarity.
- 2) *Update*: Update the superpixel features S_i^t based on the pixels features assigned to it.

The Superpixel Sampling Networks (SSN) [24] make the whole process differentiable by replacing the hard assignment between each pixel and superpixel with a soft weight:

$$Q_{pi}^t = e^{-\|I_p - S_i^{t-1}\|^2} \quad (1)$$

$$S_i^t = \frac{1}{Z_i^t} \sum_{p=1}^n Q_{pi}^t I_p, \quad (2)$$

where $Z_i^t = \sum_p Q_{pi}^t$ is the normalization constant. In practice, at $t = 0$, the superpixel features are initialized as the mean feature within a set of rectangular patches that are evenly distributed in the image. In order to reduce the computational complexity and to apply a spatial locality constraint, the distance computation $\|I_p - S_i^{t-1}\|^2$ is restricted to a local 3×3 superpixel neighborhood around each pixel, although larger window sizes are possible.

Superpixel Tokenization We propose to unroll the SSN iterations and replace the k -means clustering steps with a set of *local* cross-attentions. We initialize the superpixel features, which are distributed on a regular grid in the image, (Figure 3) with a set of randomly-initialized, learnable queries, S_i^0 , and perform the superpixel update step using cross-attention

between superpixel features and pixel features by adapting the dual-path cross-attention [41], giving

$$S_i^t = S_i^{t-1} + \sum_{p \in \mathcal{N}(i)} \text{softmax}_p(q_{S_i^{t-1}} \cdot k_{I_p^{t-1}}) v_{I_p^{t-1}} \quad (3)$$

$$I_p^t = I_p^{t-1} + \sum_{i \in \mathcal{N}(p)} \text{softmax}_i(q_{I_p^{t-1}} \cdot k_{S_i^{t-1}}) v_{S_i^{t-1}}, \quad (4)$$

where $\mathcal{N}(x)$ denotes the neighborhood of x and q , k and v are the query, key and value, generated applying a MLP to each respective feature plus an additive learned position embedding. For each superpixel neighborhood corresponding to a superpixel, we share the same set of position embeddings. For each superpixel S_i , there are $9 \cdot h \cdot w$ pixel neighbors, where $[h, w]$ is the size of the patch covered by one superpixel, while each pixel has 9 superpixel neighbors. We illustrate this neighborhood in Figure 3. The local dual-path cross-attention repeats n times to generate the output superpixel, $S_i^{t_n}$ and pixel, $I_p^{t_n}$, features.

This *local* dual-path cross-attention serves three purposes:

- Reduce complexity compared to a full cross-attention.
- Stabilize training, as the final softmax is only between 9 superpixel features or $9 \cdot h \cdot w$ pixel features.
- Encourage spatial locality of the superpixels, forcing them to focus on a coherent, local over-segmentation.

C. Superpixel Classification

Given the updated superpixel features from the Superpixel Tokenization module, we directly predict a class for each superpixel using a series of self-attentions. In particular, we apply k multi-head self-attention (MHSA) layers [40] to learn global context information between superpixels, producing outputs F_i . Performing MHSA on the superpixel features is significantly more efficient than on the pixel features, since the number of superpixels is much smaller. In our experiments, we typically use a superpixel resolution that is $32^2 \times$ smaller than the input resolution. Finally, we apply a linear layer as a classifier, *producing a semantic class prediction for each refined superpixel feature*, C_i . As opposed to the CNN pixel decoders used in other approaches [10], [12], our superpixel class predictions C_i can be directly projected back to the final pixel-level semantic segmentation output *without any additional layers*, as described in Section III-D.

D. Superpixel Association

To project the superpixel class predictions back into the pixel space, we use the outputs of the Superpixel Tokenization module, $I_p^{t_n}$ and $S_i^{t_n}$, to compute the association between each pixel and its 9 neighboring superpixels:

$$Q_{pi} = \text{softmax}_{i \in \mathcal{N}(p)}(I_p^{t_n} \cdot S_i^{t_n}). \quad (5)$$

The final dense semantic segmentation, Y , is then computed at each pixel, p , as the combination of each predicted superpixel class from the Superpixel Classification module, C_i , weighted by the above associations:

$$Y_p = \sum_{i \in \mathcal{N}(p)} Q_{pi} \cdot C_i. \quad (6)$$

During training, the dense semantic segmentation Y is supervised by the semantic segmentation ground truth.

IV. EXPERIMENTS

In this section, we evaluate the size, latency and accuracy of a small (ResNet-50 backbone) and large (ConvNeXt-L backbone) variant of our model against prior works, and provide ablations for the superpixel tokenization module and a fine grained latency analysis.

We evaluate our work on the Cityscapes [14] and ADE20K [56]. Cityscapes is a driving dataset, consisting of 5,000 high resolution street-view images, and 19 semantic classes. ADE20K is a general scene parsing dataset, consisting of 20,210 images with 150 semantic classes.

A. Implementation Details

Pixel Feature Extraction The hypercolumn features are generated by applying a MLP to project each encoder feature to 256 channels, and bilinear resizing to stride 8.

Superpixel Tokenization For both datasets, we apply 2 sequential local dual-path cross-attention to generate the superpixel embeddings, each with 256 channels and 2 heads. We use a single set of learned position embeddings for each superpixel and pixel feature, initialized at $\frac{1}{4}$ resolution of each feature, bilinear upsampled to the feature resolution and added to the feature.

Superpixel Classification 4 multi-head attention layers are applied in the superpixel classification stage, with 4 heads each, outputting 256 channels in each layer.

Superpixel Association The associations between the superpixel features and pixel features are computed at stride 8. The association is then bilinear upsampled to the input resolution before applying the softmax in (5).

Training Our training hyperparameters closely follow prior work such as [5], [50]. Specifically, we employ the polynomial learning rate scheduler, and the backbone learning rate multiplier is set to 0.1. The AdamW optimizer [25], [35] is used with weight decay 0.05. For regularization and augmentations, we use random scaling, color jittering [16], and drop path [23] with a keep probability of 0.8 (for ConvNeXt-L). We use a softmax cross-entropy loss, applied to the top 20% of pixels of the dense segmentation output.

Cityscapes Our models are trained with global batch size 32 over 32 TPU cores for 60k iterations. The initial learning rate is 10^{-3} with 5,000 steps of linear warmup. The model accepts the full resolution 1024×2048 images as input, and produces 128×256 hypercolumn features (*i.e.*, stride 8). Taking as input the hypercolumn features, the superpixel tokenization module uses 32×64 superpixels.

ADE20K Our models are trained with global batch size 64 over 32 TPU cores and crop size 640×640 . Our Resnet-50 and ConvNeXt-L models are trained for 100k and 150k iterations, respectively. The initial learning rate is 10^{-3} with 5,000 steps of linear warmup. 160×160 hypercolumn features are generated, and 40×40 superpixels are used.

Method	Backbone	Params ↓	FLOPs ↓	FPS ↑	mIoU ↑
MaskFormer [12]	ResNet-50 [21]	-	-	-	78.5
Mask2Former[11]	ResNet-50 [21]	-	-	-	79.4
Panoptic-DeepLab [10]	ResNet-50 [21]	43M	517G	-	78.7
RegProxy* [53]	ViT-S [17]	23M	270G	-	79.8
kMaX-DeepLab [†] [50]	ResNet-50 [21]	56M	434G	9.0	79.7
SP-Transformer	ResNet-50 [21]	29M	253G	15.3	80.4
Mask2Former [‡] [11]	Swin-L [31]	-	-	-	83.3
RegProxy* [53]	ViT-L/16 [17]	307M	-	-	81.4
SegFormer [47]	MiT-B5 [47]	85M	1,448G	2.5	82.4
kMaX-DeepLab [†] [50]	ConvNeXt-L [32]	232M	1,673G	3.1	83.5
SP-Transformer	ConvNeXt-L [32]	202M	1,557G	3.6	83.1

TABLE I

CITYSCAPES *val* SET RESULTS. WE EVALUATE FLOPS AND FPS WITH INPUT 1024×2048 FOR OUR SP-TRANSFORMER ON A TESLA V100-SXM2 GPU. SP-TRANSFORMER WITH RESNET-50 OUTPERFORMS PRIOR ARTS IN TERMS OF PARAMETERS, LATENCY, AND PERFORMANCE. FOR THE LARGE MODELS, SP-TRANSFORMER WITH CONVNEXT-L BACKBONE ACHIEVES SIMILAR REDUCTIONS IN PARAMETERS, WHILE ACHIEVING THE LOWEST LATENCY, AND COMPETITIVE MIOU PERFORMANCE. * REGPROXY EVALUATES USING A 768^2 SLIDING WINDOW. [†]kMAX-DEEPLAB IS TRAINED FOR PANOPTIC SEGMENTATION.

Method	Backbone	Crop	Params ↓	FLOPs ↓	FPS ↑	mIoU ↑
RegProxy [53]	ViT-Ti/16 [17]	512	6M	3.9G	38.9	42.1
MaskFormer [12]	ResNet-50 [21]	512	41M	53G	24.5	44.5
kMaX-DeepLab [†] [50]	ResNet-50 [21]	641	57M	75G	38.7	45.0
SP-Transformer	ResNet-50 [21]	640	29M	78G	40.8	43.7

TABLE II

ADE20K *val* SET RESULTS. WE EVALUATE FLOPS AND FPS WITH INPUT 640×640 FOR SP-TRANSFORMER ON A TESLA V100-SXM2 GPU. OUR METHOD OUTPERFORMS THE PRIOR REGPROXY WORK, WHILE REMAINING COMPETITIVE WITH OTHER PRIOR WORKS, AND AT THE HIGHEST FPS.

[†]kMAX-DEEPLAB IS TRAINED FOR PANOPTIC SEGMENTATION.

B. Results

1) *Cityscapes Dataset*: Table I compares the results of our Superpixel Transformer model to other transformer-based state-of-the-art models for semantic segmentation on the Cityscapes *val* set. In these experiments, we compare models with backbones roughly the same size as ResNet-50 and ConvNeXt-L. In addition to other semantic segmentation models, we also compare against the state of the art panoptic segmentation method, kMaX-DeepLab [50]. While this method is trained on a slightly different task, we find that, for computational cost, it is the most fair comparison, as our training schedule and pipeline is most similar to theirs.

With the smaller ResNet-50 backbone, our model is roughly half the size and latency of kMaX-DeepLab, while improving upon mIoU by 0.6 against the previous SOTA, RegProxy [53]. As the ResNet-50 backbone is relatively small, most existing models are dominated by the size of their decoders, which allows our model’s reduction in decoder size to have significant impacts on the overall size and performance, with a 70% improvement in FPS compared to kMaX-DeepLab. In addition, we expect to see this effect even more for even smaller backbones.

For comparison, we also provided results with a larger ConvNeXt-L backbone. Here, our model has a similar absolute reduction in params and FLOPs, as compared to the equivalent kMaX-DeepLab model. However, as the model is largely dominated by the size of the backbone, the overall improvements are more modest. Nonetheless, our

model is able to achieve near state-of-the-art performance, especially compared to the prior semantic segmentation methods, where we outperform most of the prior works, except for Mask2Former, where we are within 0.2mIoU. We note that, for this comparison, the equivalent Mask2Former [9] model is pre-trained with a significantly larger ImageNet-22k dataset, whereas our model is pre-trained on ImageNet-1k.

Figure 4 provides qualitative examples of our ConvNeXt-L model. The semantic segmentation predictions suggest that the model is able recover thin structures, such as poles, despite largely operating in a 32×64 superpixel space.

We also provide a visualization of the learned superpixels. We convert the soft association in Section III-D to a hard assignment by selecting the argmax over superpixels, i :

$$\bar{Q}_p = \operatorname{argmax}_{i \in \mathcal{N}(p)} Q_{pi} \quad (7)$$

We visualize the boundaries of these assignments overlaid on top of the input image in the left-most column of Figure 4.

From these visualizations, we can see that, despite the model being trained with a soft association, the superpixels generated by the hard assignment tightly follow the boundaries in the image. We note that this is particularly interesting as we do not provide any direct supervision to the superpixel associations, and instead these are learned implicitly by the network. In addition, we find that these boundaries tend to be more faithful to the edges of an object than the labels.

2) *ADE20K Dataset*: We provide quantitative results on the ADE20K dataset in Table II. We choose one of the most

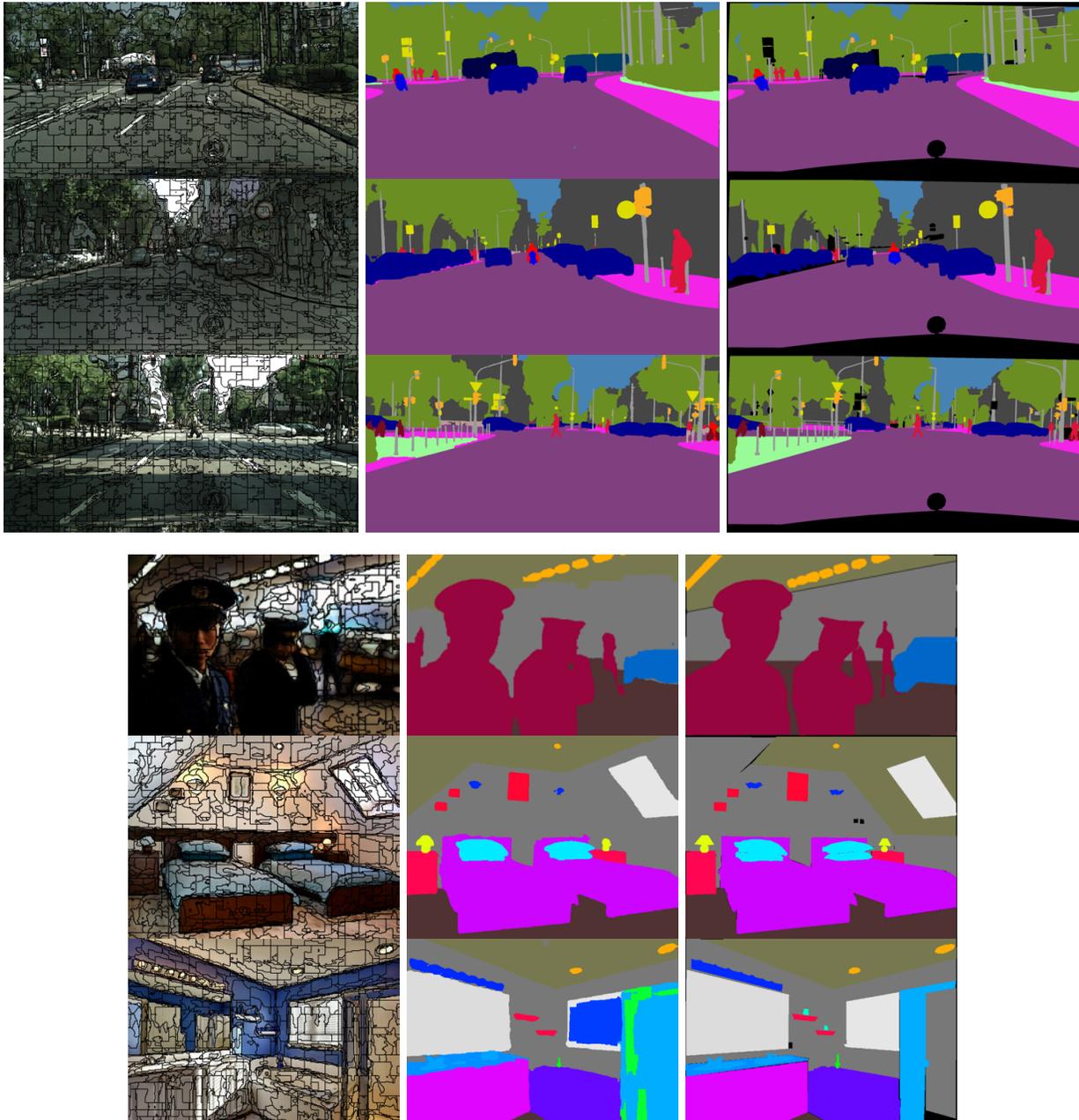


Fig. 4. Qualitative examples of our ConvNeXt-L backbone model on the CityScapes (top) and ADE20k (bottom) datasets. Left: Input image with superpixel boundaries overlaid. Middle: Semantic prediction. Right: Semantic ground truth. The superpixels are visualized by generating a hard assignment by taking the argmax of the soft assignment. Despite only being trained with a semantic segmentation loss and with soft assignments, the hard assignment superpixels faithfully follow boundaries in the image. The superpixel overlays are best viewed zoomed in.

commonly used crop sizes (640×640) and provide results for the ResNet-50 backbone.

For ADE20K, we achieve the highest FPS and the second lowest # params (behind the surprisingly small RegProxy model), while outperforming RegProxy.

However, we do note that the gap in performance is larger for ADE20K than Cityscapes. Our hypothesis is that the large number of classes in ADE20K (152) results in ambiguities when a pixel could belong to multiple classes. This results in inconsistencies for object classes, and in particular where

boundaries are drawn (see Figure 5 for examples). As our superpixel tokenization module operates before any semantic prediction, and each superpixel query only operates on a local neighborhood in the pixel space, the model must learn a consistent way to divide the image into a set of superpixels. When the label boundaries are inconsistent, our model is less able to effectively learn this over-segmentation, leading to a small decrease in performance.

3) *Superpixel Tokenization Ablation:* In Table III, we provide an ablation of the number of cross attention layers



Fig. 5. Inconsistent label boundaries in the ADE20K dataset make it difficult for our model to effectively learn superpixels. See the spaces in between fence railings, the books on the shelves, the unlabeled people on the bleachers and the trees/vegetation.

# CA	sp. res.	params	FLOPs	Runtime(ms)	mIoU
1	32×64	28M	234G	50.0	78.0
2	16×32	29M	240G	53.5	74.9
2	32×64	29M	253G	63.4	80.4
2	64×128	31M	404G	120.5	79.7

TABLE III

EXPERIMENTS ABLATING THE NUMBER OF LOCAL CROSS-ATTENTION LAYERS (# CA) AND THE RESOLUTION OF THE SUPERPIXELS. ABLATIONS ARE PERFORMED WITH A RESNET-50 BACKBONE ON CITYSCAPES.

and the superpixel resolution, evaluated using our ResNet-50 backbone model on Cityscapes.

From these experiments, we find that increasing the number of cross attention layers from 1 to 2 improves mIoU significantly by 2.4. Further increasing the number of these layers may have further improvements in performance, although we are currently bottlenecked by accelerator memory constraints. This is largely due to our naive TensorFlow implementation [45], which we discuss in Section IV-C.2.

Reducing the number of superpixels to 16×32 also has a significant regression in mIoU of 5.5. On the other hand, increasing the superpixels to 64×128 also results in a mild regression of 0.7. We posit that this is because our pixel features used in the association are at stride 8 (128×256). As a result, having 64×128 superpixels provides each local-cross attention with a neighborhood of $128/64 \times 3 = 6 \times 6$ pixels, which makes the receptive field for each superpixel too small to learn the oversegmentation as effectively. This can be resolved by increasing size of the neighborhood around each superpixel, but at significantly higher computational cost.

C. Discussion

1) *Superpixel Quality*: Despite our network not being trained with any explicit superpixel-based loss, we find that the associations learned by the network closely resemble classical superpixels. That is, the superpixels are aligned such that they follow the dominant edges in the image. We posit that this is due to the limited receptive field for each superpixel’s cross attention. As any single superpixel may not have visibility to all of the pixels for a given mask, the model must use local edges and boundaries to separate the

Method	Time (ms)
Backbone: ResNet-50	28.0
Hypercolumn	7.1
Superpixel Tokenization	25.4
Superpixel Self-Attention	4.0
Superpixel Association	1.0
Total	65.5

TABLE IV

LATENCY INFORMATION FOR EACH SUB-COMPONENT. THE BULK OF THE NON-BACKBONE LATENCY IS CONSUMED BY THE LOCAL CROSS ATTENTION IN THE SUPERPIXEL TOKENIZATION MODULE.

superpixels.

2) *Latency*: As seen in Tables I and II, our method provides a significant improvement in FPS. As the main processing of our model operates on the small superpixel space, this allows for a large reduction in model complexity, while achieving state of the art performance.

However, we believe that there is still a large further reduction in latency available. In particular, the local cross attention operation is not efficient for standard accelerators and native TensorFlow or PyTorch implementations. This is because it requires a sliding window with overlapping patches, but with different operands in each patch (as opposed to convolutions where the weights are the same for all patches). However, as we operate on the superpixel grid level (32×64 for Cityscapes and 40×40 for ADE20K), this impact of this inefficiency is low enough to make our method overall faster compared to methods which operate on the dense pixel space.

Nonetheless, the Superpixel Tokenization module takes up the majority of the decoder runtime, as can be seen in Table IV, where we provide timing information for our method with a ResNet-50 backbone on a 1024×2048 input. Our experiments are currently designed with a relatively naive, pure TensorFlow implementation, and involves the duplication of each superpixel or pixel 9 times (depending on the cross-attention operation). We believe that a CUDA implementation could remove this redundant copy, and provide even further speedups for our method.

V. CONCLUSIONS

We presented a novel network architecture for semantic segmentation that leverages superpixels to project the dense image segmentation problem into a low dimensional superpixel space. Operating on this space enables us to significantly reduce the size and inference latency of our network compared to prior works, while achieving state-of-the-art performance.

REFERENCES

- [1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(11):2274–2282, 2012.
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.
- [3] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. *arXiv:2102.04306*, 2021.

- [4] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015.
- [5] 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 Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [6] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv:1706.05587*, 2017.
- [7] Liang-Chieh Chen, Yi Yang, Jiang Wang, Wei Xu, and Alan L Yuille. Attention to scale: Scale-aware semantic image segmentation. In *CVPR*, 2016.
- [8] 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.
- [9] Bowen Cheng, Anwesa Choudhuri, Ishan Misra, Alexander Kirillov, Rohit Girdhar, and Alexander G Schwing. Mask2former for video instance segmentation. In *CVPR*, 2022.
- [10] 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.
- [11] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022.
- [12] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In *NeurIPS*, 2021.
- [13] Dorin Comaniciu and Peter Meer. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(5):603–619, 2002.
- [14] 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.
- [15] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- [16] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. In *CVPR*, 2019.
- [17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020.
- [18] Brian Fulkerson, Andrea Vedaldi, and Stefano Soatto. Class segmentation and object localization with superpixel neighborhoods. In *ICCV*, 2009.
- [19] Raghudeep Gadde, Varun Jampani, Martin Kiefel, Daniel Kappler, and Peter V Gehler. Superpixel convolutional networks using bilateral inceptions. In *ECCV*, 2016.
- [20] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Hypercolumns for object segmentation and fine-grained localization. In *CVPR*, 2015.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [22] Xuming He, Richard S Zemel, and Miguel A Carreira-Perpinán. Multiscale conditional random fields for image labeling. In *CVPR*, 2004.
- [23] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016.
- [24] Varun Jampani, Deqing Sun, Ming-Yu Liu, Ming-Hsuan Yang, and Jan Kautz. Superpixel sampling networks. In *ECCV*, pages 352–368, 2018.
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [26] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *NeurIPS*, 2011.
- [27] L’ubor Ladický, Chris Russell, Pushmeet Kohli, and Philip HS Torr. Associative hierarchical crfs for object class image segmentation. In *ICCV*, 2009.
- [28] John Lafferty, Andrew McCallum, and Fernando CN Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *ICML*, 2001.
- [29] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [30] Zhiqi Li, Wenhai Wang, Enze Xie, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, Ping Luo, and Tong Lu. Panoptic segformer: Delving deeper into panoptic segmentation with transformers. In *CVPR*, 2022.
- [31] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- [32] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, 2022.
- [33] Stuart Lloyd. Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137, 1982.
- [34] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- [35] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
- [36] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. In *NeurIPS*, 2019.
- [37] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- [38] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segformer: Transformer for semantic segmentation. In *ICCV*, 2021.
- [39] Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. In *ECCV*, 2020.
- [40] 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*, volume 30, 2017.
- [41] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Max-deeplab: End-to-end panoptic segmentation with mask transformers. In *CVPR*, 2021.
- [42] 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.
- [43] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018.
- [44] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. In *NeurIPS*, 2020.
- [45] Mark Weber, Huiyu Wang, Siyuan Qiao, Jun Xie, Maxwell D Collins, Yukun Zhu, Liangzhe Yuan, Dahun Kim, Qihang Yu, Daniel Cremers, et al. Deeplab2: A tensorflow library for deep labeling. *arXiv:2106.09748*, 2021.
- [46] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *ECCV*, 2018.
- [47] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *NeurIPS*, 2021.
- [48] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *CVPR*, 2022.
- [49] Qihang Yu, Huiyu Wang, Dahun Kim, Siyuan Qiao, Maxwell Collins, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Cmt-deeplab: Clustering mask transformers for panoptic segmentation. In *CVPR*, 2022.
- [50] Qihang Yu, Huiyu Wang, Siyuan Qiao, Maxwell Collins, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. k-means mask transformer. In *ECCV*, 2022.
- [51] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation. In *ECCV*, 2020.
- [52] Wenwei Zhang, Jiangmiao Pang, Kai Chen, and Chen Change Loy. K-net: Towards unified image segmentation. In *NeurIPS*, 2021.
- [53] Yifan Zhang, Bo Pang, and Cewu Lu. Semantic segmentation by early region proxy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1258–1268, 2022.
- [54] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *CVPR*, 2017.
- [55] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *CVPR*, 2021.
- [56] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *CVPR*, 2017.
- [57] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *ICLR*, 2020.