# Progressive Semantic Segmentation

Chuong Huynh[1]    Anh Tuan Tran[1,2]    Khoa Luu[1,3]    Minh Hoai[1,4]

[1]VinAI Research, Hanoi, Vietnam, [2]VinUniversity, Hanoi, Vietnam
[3]University of Arkansas, Fayetteville, AR 72701, USA
[4]Stony Brook University, Stony Brook, NY 11790, USA

{v.chuonghm,v.anhtt152,v.khoal,v.hoainm}@vinai.io

## Abstract

*The objective of this work is to segment high-resolution images without overloading GPU memory usage or losing the fine details in the output segmentation map. The memory constraint means that we must either downsample the big image or divide the image into local patches for separate processing. However, the former approach would lose the fine details, while the latter can be ambiguous due to the lack of a global picture. In this work, we present Mag-Net, a multi-scale framework that resolves local ambiguity by looking at the image at multiple magnification levels. MagNet has multiple processing stages, where each stage corresponds to a magnification level, and the output of one stage is fed into the next stage for coarse-to-fine information propagation. Each stage analyzes the image at a higher resolution than the previous stage, recovering the previously lost details due to the lossy downsampling step, and the segmentation output is progressively refined through the processing stages. Experiments on three high-resolution datasets of urban views, aerial scenes, and medical images show that MagNet consistently outperforms the state-of-the-art methods by a significant margin. Code is available at https://github.com/VinAIResearch/MagNet.*

## 1. Introduction

The current state-of-the-art (SOTA) semantic image segmentation techniques [1, 4, 16, 19, 21, 23, 26] are based on deep learning, where a convolutional neural network (CNN) takes an input image and outputs a segmentation map. Most of the existing techniques, however, assume that the entire segmentation process can be performed with a single feed-forward pass of the input image and the entire process can be fitted into GPU memory. Unfortunately, most existing techniques cannot handle high-resolution input images due to memory and other computational constraints. One approach to handle a large input image is to downsample it, but this results in a low-resolution segmentation map, which is not adequate for applications that require high-resolution
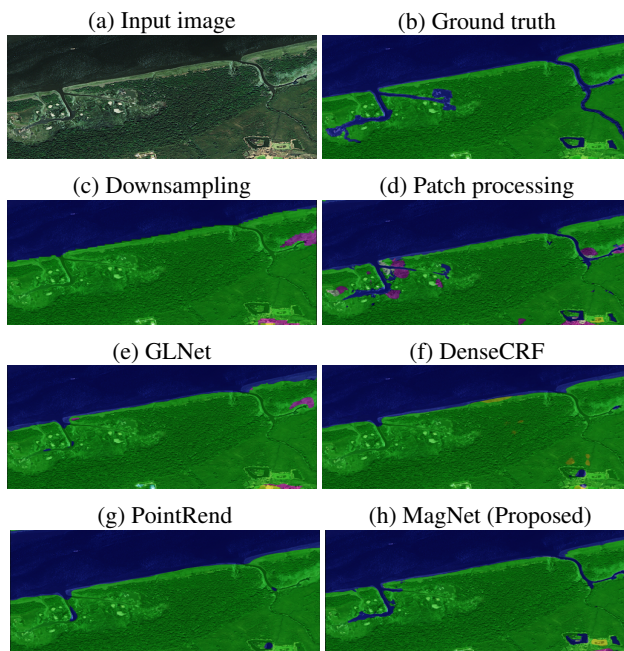


(a) Input image        (b) Ground truth
(c) Downsampling       (d) Patch processing
(e) GLNet              (f) DenseCRF
(g) PointRend          (h) MagNet (Proposed)

Figure 1: **Comparing several semantic segmentation and refinement approaches on a high-resolution input image**. Downsampling loses fine details, while Patch Processing wrongly classifies local patches due to the lack of the global context. The collaborative global-local network GLNet fails due to the large discrepancy between the global and local branches. Post-processing and refinement methods such as DenseCRF and PointRend can only correct small mistakes due to local inconsistency. MagNet outperforms other methods, thanks to a novel multi-scale segmentation and refinement framework. Best viewed in color.

output with fine details [12, 24], *e.g.*, for tracking the progression of malignant lesion [8]. Another approach to handle a large input image is to divide the image into small patches and process each patch independently. This approach, however, does not take into account the global information [22] that is needed to resolve ambiguity in local patches. The limitations of these two approaches are illustrated in Fig. 1(c) & (d).

One way to address the limitations of the two aforemen-

tioned approaches is to combine them, i.e., to fuse global and local segmentation processes. On the one hand, the global view of the entire image can be used to resolve the ambiguity in the appearance of local patches. On the other hand, by analyzing local patches, we can refine the segmentation boundaries and recover the lost details due to the downsampling procedure of the global segmentation process. This approach has been successfully demonstrated recently by the Global-Local Network (GLNet) [5]. However, given an input image with ultra-high resolution, there is a huge gap between the scale of the whole image and the scale of the local patches. This will lead to contrasting output segmentation maps, and it will be difficult to combine and reconcile differences with a single feed-forward processing stage (see Fig. 1e); the difficulty of this combination task is analogous to constructing a single-span bridge across a wide river.

To bridge the gap between the two extreme ends of the scale space, we propose to consider multiple scales in between. We introduce a novel multi-scale framework where the output segmentation map will be progressively refined as the image is analyzed from the coarsest to the finest scale. The core of our framework is a refinement module that can use one segmentation map to refine another. This refinement module is used at every stage of our multi-scale processing pipeline to refine the output segmentation map at its most uncertain locations. Our framework can integrate global contextual cues to produce more accurate segmentation, and it can output high-resolution detailed segmentation maps under a memory constraint. Fig. 1 shows the result of MagNet and compares it with other segmentation methods, including the recently proposed PointRend [14] method that seeks to refine only at the most uncertain pixels.

## 2. Related Work

**Multi-scale, multi-stage, context aggregation.** The combination of multiple scale levels helps the network aggregate different fields of view and provides more context to each pixel [3, 11]. ICNet [38] used a cascaded architecture for feature maps of different downsampled inputs, while RefineNet [18] fused upsampled outputs of the branches that handled different low-resolution inputs. Feature Pyramid Network (FPN) [13] upsampled feature maps in different scales and aggregated them with the output of low layers. The dilated convolution and Atrous Spatial Pyramid Pooling (ASPP) module in DeepLab [4] enlarged the receptive field, created a connection between far-apart pixels. The same effect was achieved by PSPNet [37], which combined different scaling feature maps to enlarge the receptive fields. High-resolution Net (HRNet) [30] proposed another scale fusion schema where a new branch with a larger receptive field was added after each stage. Attention technique was also used in many recent approaches [2, 34, 35] to add

more global information to every single point.

Another approach to handle high-resolution images is to use multi-stage networks, where images are segmented after several stages or sub-networks. Xia et al. [32] proposed Hierarchical Auto-Zoom Net, a strategy to scale the field of view when sliding the view through a big image. For ultra-high resolution images, Takahama et al. [27] solved the imbalance between background and foreground by predicting whether the whole patch contained foreground pixels or not before segmenting it.

Propagating global information to local patches is a promising approach to deal with high-resolution images. ParseNet [20] pooled global context to local field of view to have more information. BiSeNet [33] included one more branch for global pooling and the global context would be added to the feature map at the last stage. Although those methods were efficient, they consumed a huge amount of GPU memory for ultra-high resolution images. Tokunaga et al. [28] proposed a method for super-high-resolution images, using independent multi-scale networks and an adaptive weight generator. The outputs of network members were combined with corresponding trained weights to produce the final output, but there was no knowledge sharing between network branches. Unlike [28], [5] contained two sub-networks with shared information, where the global branch took the downsampled images to extract global context and the local branch took patches and corresponding global features to improve the details of high-resolution images. However, due to the ad-hoc combination of the global and local branches, it was difficult to extend to more than two scales. Moreover, in our experience, training the local network was difficult due to the domination of the strong global branch.

**Segmentation refinement.** There were several approaches to improve the segmentation outputs with post-processing. One approach was to use classical methods such as Conditional Random Fields (CRFs) [15] or Guided Filter (GF) [10] on the segmentation mask produced by deep learning networks. However, these methods were slow and the improvement was incremental. The inference speed could be improved with a deep learning version of Guided Filter (DGF) [31]. Another approach for post-processing was to use deep networks. Iterative Instance Segmentation (ISS) [17] refined the output by repeatedly passing the input image and the segmentation map through a refinement module several times. This method was based on self-reflection, the input image to each refinement stage was the same. Like ISS, CascadedPSP [6] used the same refinement scheme but the resolution of the input at each refinement stage was different. However, the wrong prediction at any middle stage could significantly affect the performance of later steps. Some methods aimed to refine parts of the output only, such as pixels at the boundaries [36, 39] or pixels

| Method | Dense CRF [15] | GF [10] | DGF [31] | ISS [17] | GLNet [5] | Cascade PSP [6] | PointRend [14] | SegFix [36] | DeepStrip [39] | Ours |
|---|---|---|---|---|---|---|---|---|---|---|
| Deep learning | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Using high-resolution input | ✓ | ✓ | ✓ | | ✓ | ✓ | | ✓ | ✓ | ✓ |
| Multi-scale processing | | | | | ✓ | ✓ | ✓ | | | ✓ |
| Can recognize new objects | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | ✓ |
| Partly refinement | | | | | | | ✓ | ✓ | ✓ | ✓ |

Table 1: Summary of key features of various semantic segmentation refinement approaches.

at uncertain locations (PointRend [14]). However, boundary refinement methods [36, 39] failed to recover tiny objects, while PointRend [14] only used the local context for refinement. Furthermore, because the input of the PointRend was the high-level features of a deep network, it must be trained specifically for each segmentation backbone. In this paper, we propose a modular framework for having any number of scale levels. It is simple but effective for refining a coarse segmentation output, being able to keep the overall structure of the coarse segmentation output while adding more details after each stage without suffering the domination problems. Table 1 compares the key features of different methods.

## 3. MagNet

We now describe MagNet, a multi-scale segmentation framework for ultra-high resolution images. It is a multi-stage network architecture, where each stage of the network corresponds to a particular scale. An input image will be inspected at multiple scales, from the coarsest to the finest.

The core of our framework is a segmentation module and a refinement module, which are used at every processing stage. At each stage, the inputs to the refinement module are two segmentation maps: (1) the cumulative result from the previous stages, and (2) the result obtained by running the segmentation module at and only at the current scale. The objective of the refinement module is to use the latter segmentation map to refine the former one, at selective locations determined based on the uncertainty of two estimated segmentation maps.

In our framework, the segmentation module can be any segmentation backbone, as long as it can output a segmentation map with uncertainty estimates. The refinement module is agnostic to the segmentation backbone; it can be trained with one backbone and used with another. In the following, we will describe the multi-stage processing pipeline and the refinement module in details.

### 3.1. Multi-stage processing pipeline

The architecture and processing pipeline of MagNet is depicted in Fig. 2. There is one segmentation module and one refinement module, which are used repeatedly in $m$ processing stages, where $m$ is a hyper-parameter for the num-

ber of scales that we want to analyze. We use $s$ to denote the processing stage, where $s = 1$ corresponds to the coarsest scale and $s = m$ corresponds to the finest scale. Let $X \in \mathbb{R}^{H \times W \times 3}$ be an input image, where $H, W$ are the height and width of the image. We consider the case when $H$ and $W$ are too big for image $X$ to be processed without downsampling, and let $h$ and $w$ be the largest (or desirable) sizes that can be handled by the segmentation module. We use $h^s$ and $w^s$ to denote the height and width for the scale level $s$. We determine the scale levels so that they span the entire scale space: $H = h^1 > \cdots > h^m = h$ and $W = w^1 > \cdots > w^m = w$.

For a particular scale level $s$, we divide the input image $X$ into patches of size $h^s \times w^s$ and perform semantic segmentation on these patches. The locations of these patches are defined by a set of rectangular windows, and let $\mathcal{P}^s$ denote the set of these windows: $\mathcal{P}^s = \{\mathbf{p} | \mathbf{p} = (x, y, w^s, h^s)\}$, where each window is specified by the top-left corner, width, and height. As the scale level $s$ increases, the width and height of the rectangular windows decrease, but the cardinality of $\mathcal{P}^s$ increases. For a particular window $\mathbf{p}$, we will use $X_{\mathbf{p}}$ to denote the image patch extracted at the window $\mathbf{p}$.

Our network will take an image $X \in \mathbb{R}^{H \times W \times 3}$ and produce a sequence of segmentation maps $Y^1, \cdots, Y^m \in \mathbb{R}^{H \times W \times C}$, where $C$ is the number of semantic categories in consideration. At stage $s$, we first determine the set of rectangular windows $\mathcal{P}^s$ for patch division and refine the segmentation map of each patch. Specifically, for each window $\mathbf{p} \in \mathcal{P}^s$, do:

1. Extract the image patch $X_{\mathbf{p}}^s$ and previous segmentation output $Y_{\mathbf{p}}^{s-1}$ defined by the window $\mathbf{p}$. The height and width of these tensors are $h^s$ and $w^s$.

2. Downsample $X_{\mathbf{p}}^s$ and $Y_{\mathbf{p}}^{s-1}$ so that the new height and width are $h$ and $w$, which are the size that can be fitted into GPU memory and be processed by the segmentation and refinement modules. Let $\bar{X}_{\mathbf{p}}^s$ and $\bar{Y}_{\mathbf{p}}^{s-1}$ denote the downsampled tensors.

3. With $\bar{X}_{\mathbf{p}}^s$ as the input, use the segmentation module to obtain the scale-specific segmentation map $\bar{O}_{\mathbf{p}}^s$.
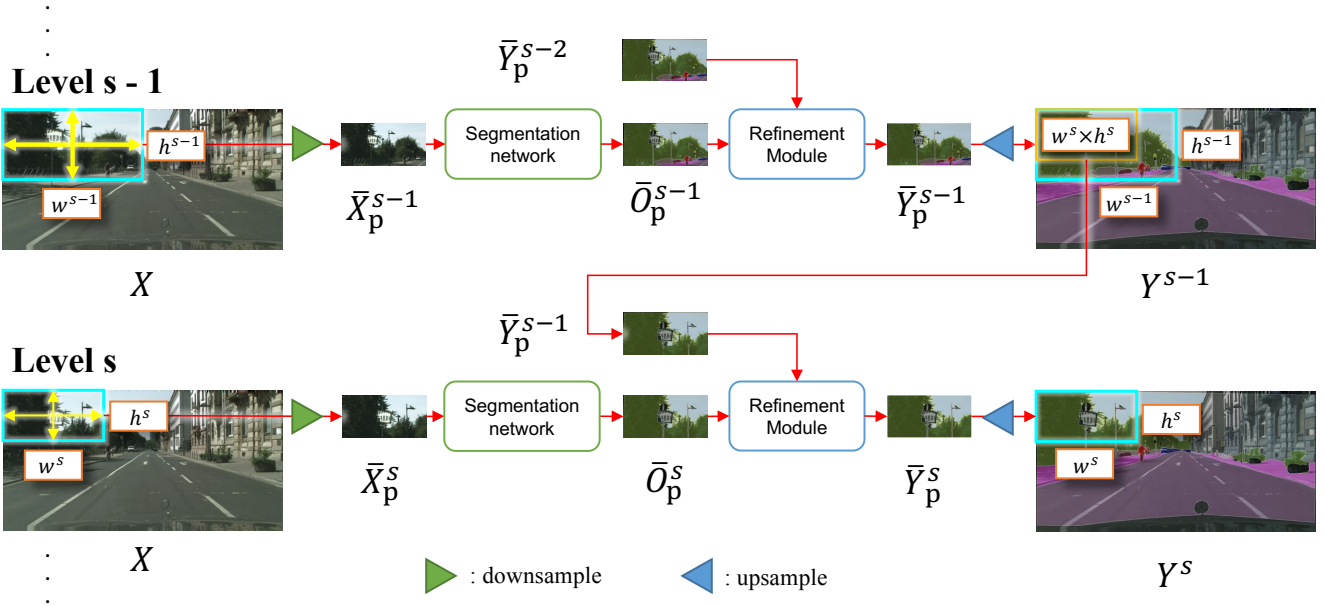
Figure 2: **Overview of our proposed MagNet**. The segmentation network produces the scale-specific prediction while the refinement module selectively refines the coarse prediction from previous stages based on that local prediction.
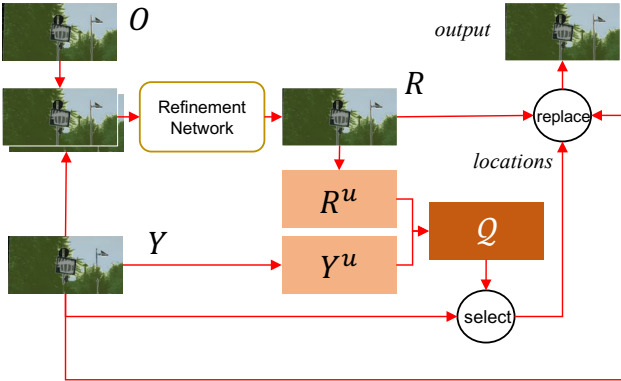


Figure 3: **The overview of the refinement module**. The cumulative segmentation $Y$ is partly replaced with the scale-specific segmentation map $O$ based on the score $\mathcal{Q}$.

4. With $\bar{Y}_{\mathbf{p}}^{s-1}$ and $\bar{O}_{\mathbf{p}}^{s}$ as the input, use the refinement module to refine $\bar{Y}_{\mathbf{p}}^{s-1}$ to obtain $\bar{Y}_{\mathbf{p}}^{s}$ (See Sec. 3.2).

5. Upsample $\bar{Y}_{\mathbf{p}}^{s}$ to get $Y_{\mathbf{p}}^{s}$ of size $h^s \times w^s \times C$.

These processing steps are illustrated in Fig. 2.

### 3.2. Refinement module

The refinement module is a core component of our framework, which is used to refine the individual patches of a segmentation map at every processing stage of our pipeline. The input to this module is two segmentation maps of size $h \times w \times C$: (1) the scale-cumulative segmentation map $Y$, from all previous scales, and (2) the scale-specific segmentation map $O$, from the current scale. The output of the module is the updated scale-cumulative segmentation map.

Fig. 3 depicts the refinement process, which contains the following steps. First, using a small network with $Y$ and

$O$ as the input, we obtain an initial combined segmentation map $R$. We then calculate the prediction uncertainty maps for both $Y$ and $R$. Specifically, for each pixel of $Y$, the prediction confidence at this location is defined as the absolute difference between the highest probability value and the second-highest value (among the $C$ probability values for $C$ classes). The uncertainty score is then computed based on the confidence score such that the two scores must add up to one. Similarly, we can compute the prediction uncertainty map for $R$. Let $Y^u$ and $R^u$ denote the prediction uncertainty maps for $Y$ and $R$ respectively.

Next, we will describe how to use two uncertainty maps to select $k$ locations of $Y$ for refinement. While only one uncertainty prediction is used in the previous work PointRend [14], we extend this approach to have a better selection strategy. These are the locations where $Y$ is uncertain about its prediction, while $R$ is certain about its prediction. The score map for ranking the pixels is calculated as $\mathcal{Q} = \mathcal{F}(Y^u \odot (1 - R^u))$, where $\odot$ denotes point-wise multiplication and $\mathcal{F}$ denotes median blurring to smooth the score map. Empirical comparison between the effect of each element in the formula can be found in Sec. 4.2.3.

### 3.3. MagNet-Fast

There is trade-off between the accuracy and the run-time efficiency of our framework. One way to reduce the running time is to decrease the number of scales to process. Another approach is to perform segmentation and refinement on a subset of image patches at each scale level. MagNet-Fast combines these two approaches when it runs on a smaller number of scales and only selects the patches with the highest prediction uncertainty $Y^u$ for refinement. In MagNet-Fast, the total number of image patches that need to be fed
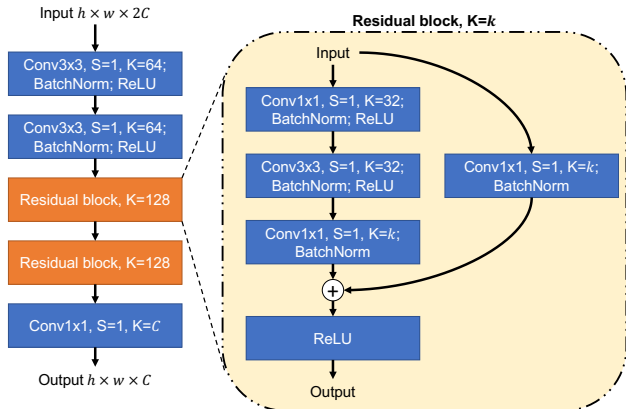
Figure 4: **The two residual blocks are trained to refine the segmentation at each scale**. This module outputs the same size $h \times w$ as the input.

into the segmentation module might be even smaller than the number of image patches used in the patch processing approach. Moreover, MagNet-Fast can leverage both global context and detailed information for segmentation, leading to superior results as will be seen in our experiments.

## 4. Experiments

We evaluated the performance of MagNet on three high resolution datasets: Cityscapes [7], DeepGlobe [9], and Gleason [29]. Some information about these datasets is listed in Table 2. The number of pixels of each image is from 2 to 25 million. In this section, we present experiments comparing MagNet with other state-of-the-art frameworks in semantic segmentation and also describe some ablation studies on Cityscapes.

| Dataset | Content | Resolution | No.classes |
|---|---|---|---|
| Cityscapes [7] | urban scene | 2048×1024 | 19 |
| DeepGlobe [9] | aerial scene | 2448×2448 | 6 |
| Gleason [29] | histopathology | 5000×5000 | 4 |

Table 2: **Details of high-resolution datasets used to evaluate our framework.** All images have from 2 to 25 million pixels with a lot of details.

### 4.1. Implementation details

**Architecture of the refinement module.** Fig. 4 depicts the architecture of the refinement module used in all experiments. The main components are the two residual blocks. With the input of size $h \times w \times 2C$, the refinement module produces the output of size $h \times w \times C$.

**Training.** For each dataset, we trained a state-of-the-art segmentation model on the downsampled images and a refinement module to refine the coarse output on sliced images. While training the refinement module, we randomly extracted image patches and also applied the following data

| Refinement steps | mIoU(%) | Time(s) |
|---|---|---|
| 256 | 63.23 | 0.03 |
| 256→512 | 65.73 | 0.19 |
| 256 → 1024 | 65.23 | 0.61 |
| 256 → 2048 | 65.21 | 2.24 |
| 256→512 → 2048 | 67.13 | 2.38 |
| 256 → 1024→2048 | 66.95 | 2.79 |
| 256→512→1024→2048 | **67.57** | 2.93 |

Table 3: **Performance of MagNet on Cityscapes with and without intermediate scale levels**. It is essential to have the intermediate scales.

augmentation: rotation, and horizontal and vertical flipping. We used SGD optimizer with momentum $0.9$, decayed weight $5 \times 10^{-4}$, and initial learning rate $10^{-3}$. We trained the refinement module for 50 epochs, and the learning rate was decreased tenfold at epoch 20, 30, 40, and 45. We used cross-entropy as the loss function for training segmentation and refinement modules. We implemented Mag-Net using PyTorch [25] starting from the public implementation of HRNet-OCR [35]. We use a batch size of 16 for training on a DGX-1 workstation with Tesla V100 GPUs.

**Testing.** During inference, at each scale, we extracted non-overlapping patches for processing. The evaluation metric is mean Intersection over Union (mIoU). For memory and speed comparison, we ran benchmarking experiments on a machine with an Intel i7 CPU and an RTX2080Ti GPU.

### 4.2. Experiments on the Cityscapes dataset

Cityscapes is a dataset of high-resolution urban scenes, containing images of size $2048 \times 1024$ pixels. The task is to segment objects in videos captured by auto cameras. There are two kinds of data, with coarse and fine labels. In our experiments, we used the fine-labeled dataset, with 2,975 training images and 500 testing images.

We considered four possible scale levels ($m$=4), corresponding to patch sizes of $2048 \times 1024$, $1024 \times 512$, $512 \times 256$, and $256 \times 128$. The size of the input to the segmentation module was always $256 \times 128$ to satisfy the memory constraint, so any larger patch would be downsampled.

#### 4.2.1 Benefits of multiple scale levels

Table 3 shows the results of MagNet for a different number of scales. While the direct refinement from the lowest to highest scale improves about 2% mIoU, from 63.23% to 65.23%, adding the two intermediate scales between the smallest and largest scales improve the performance by 4.34% mIoU. Qualitative improvements through each processing stage can be observed in Fig. 5. After each stage, the errors decrease and the segmentation masks become finer.

| Model | mIoU(%) | Mem.(MB) | Time(s) |
|---|---|---|---|
| Patching | 52.19 | 1575 | 1.77 |
| Downsampling | 63.23 | 1575 | 0.02 |
| DenseCRF [15] | 62.95 (↓0.28) | 1575 | 26.02 |
| DGF [31] | 63.33 (↑0.10) | **1727** (↑152) | 0.32 |
| PointRend [14] | 64.39 (↑1.16) | 2033 (↑458) | **0.14** |
| SegFix [36] | 65.83 (↑2.60) | 2961 (↑1386) | 0.38 |
| MagNet-Fast | 66.91 (↑3.68) | 2007 (↑432) | 0.32 |
| MagNet | **67.57** (↑4.34) | 2007 (↑432) | 2.93 |

Table 4: **Performance of MagNet and other segmentation refinement methods on the Cityscapes dataset**. The backbone HRNetV2-W18+OCR [35] was used as the segmentation module for all refinement methods.

### 4.2.2 Comparing segmentation approaches

Table 4 and Table 5 compare the performance of MagNet with several state-of-the-art semantic segmentation refinement methods. All methods were trained and applied to the output of the pretrained HRNet+OCR [35] model, which is among the leading methods on this dataset. Although there are various HRNet models [30], we used HRNetV2-W18



| | 256 x 128 | 512 x 256 | 1024 x 512 | 2048 x 1024 |
|---|---|---|---|---|

Figure 5: **The visualization of segmentation output through each processing stage on Cityscapes dataset**. The first one is the image. Others in the first row are the selective points to be refined (red color). The second row is the segmentation output. The third row is the errors comparing with the ground-truth. Best viewed in color.

| Model | fence | pole | traffic sign | person | motor-cycle |
|---|---|---|---|---|---|
| Patching | 33.32 | **42.87** | 59.39 | 61.23 | 22.12 |
| Downsample | 45.38 | 35.58 | 54.97 | 64.64 | 36.16 |
| DenseCRF [15] | 45.30 | 32.27 | 54.69 | 64.32 | 36.20 |
| DGF [31] | 45.42 | 35.63 | 55.27 | 64.85 | 36.20 |
| PointRend [14] | 45.01 | 42.71 | **60.18** | 67.10 | **39.17** |
| SegFix [36] | **46.17** | 38.77 | 59.23 | **67.86** | 37.12 |
| MagNet-Fast | 48.57 | 46.38 | **64.84** | 69.65 | 41.98 |
| MagNet | **50.59** | **49.39** | 64.15 | **72.16** | **45.19** |
| Improvement | 4.42 | 6.52 | 4.66 | 4.30 | 6.02 |

Table 5: **IoU for some specific categories on Cityscapes.** The best and previous best method is highlighted in red and blue color respectively, and the difference between them is shown in the last row.

given its manageable complexity that was required for our experiments, especially in terms of memory constraint.

MagNet-Fast is an efficient version of MagNet, in which only the most uncertain patches are refined at each scale. In this experiment, we selected the number of patches so that MagNet-Fast had a similar processing speed as SegFix [36] and DGF [31]. This model ran on only three scales $256 \rightarrow 512 \rightarrow 2048$ and the three highest uncertain patches for each scale. In total, MagNet-Fast needed to run inference on $1+3+3 = 7$ patches, comparing to $1+4+16+64 = 85$ patches of MagNet, and 64 patches of the patch processing approach.

In this experiment, except for DenseCRF [15], we fine-tuned other frameworks to achieve the best result with the segmentation backbone. For SegFix [36], with the offset prediction published by the authors, the best result was achieved with the offset width of 10. Both DGF [31] and PointRend [14] were trained with the output of the segmentation backbone. Besides, DGF ran on patches to be fairly compared with our method in speed and memory.

As can be observed, DenseCRF [15] is the only method that cannot improve the coarse segmentation. PointRend [14] is the fastest method, but the improvement is small. The inference times of MagNet-Fast, SegFix [36], and DGF [31] are similar, but MagNet-Fast outperforms the others significantly. MagNet was slower, but it had the highest mIoU.

The cumulative IoU distributions of these methods in our experiment are shown on Fig. 6a. There is a big gap between MagNet and the other methods, especially when looking at the zoom-in window.

Fig. 6b shows the results of several methods on two Cityscapes images. Both MagNet and MagNet-Fast yield the best refinement. SegFix [36] cannot recover small objects, such as sign poles, that are wrongly merged with a bigger region, while PointRend [14] performs poorly due to
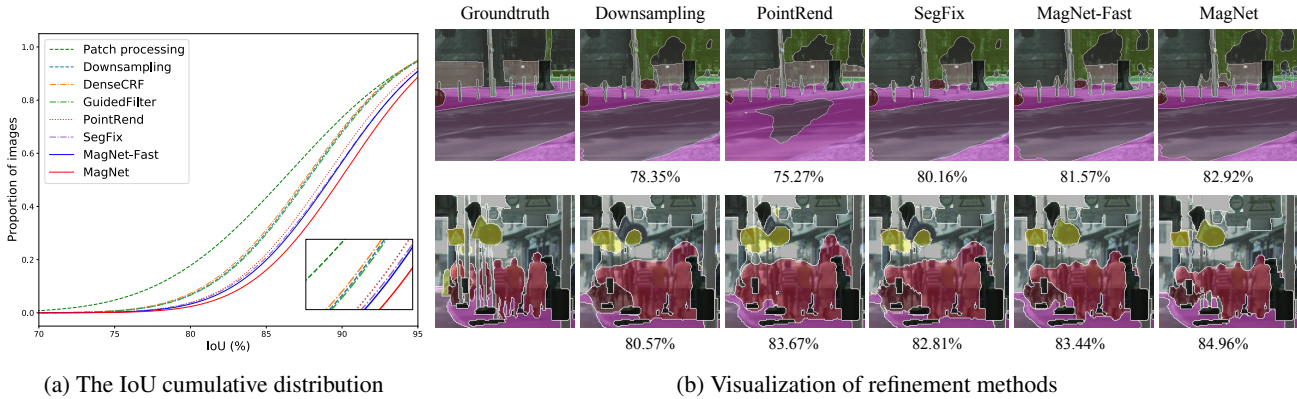
| Groundtruth | Downsampling | PointRend | SegFix | MagNet-Fast | MagNet |
|---|---|---|---|---|---|
| | 78.35% | 75.27% | 80.16% | 81.57% | 82.92% |
| | 80.57% | 83.67% | 82.81% | 83.44% | 84.96% |

(a) The IoU cumulative distribution      (b) Visualization of refinement methods

Figure 6: **Our methods outperform other refinement frameworks on the Cityscapes dataset**. (a) The cumulative distribution of mIoU of each image on the dataset (lower is better). The MagNet and MagNet-Fast achieve the best result among others. (b) Some segmentation results of refinement frameworks and our MagNet. The mIoU numbers are below the images. More tiny objects are recognized and boundaries are refined better with MagNet and MagNet-Fast. Best viewed in a digital device with magnification.

| $Y^u$ | $(1 - R^u)$ | $\mathcal{F}$ | mIoU (%) |
|:---:|:---:|:---:|---:|
| ✓ | | | 63.22 |
| ✓ | | ✓ | 63.25 |
| | ✓ | | 66.36 |
| | ✓ | ✓ | 66.46 |
| ✓ | ✓ | | 67.37 |
| ✓ | ✓ | ✓ | **67.57** |

Table 6: **Performance of MagNet on Cityscapes with different ranking scores**. The initial segmentation has mIoU of 63.23%. With $k=2^{16}$, the framework achieves the best performance when using both $Y^u$ and $(1 - R^u)$ with smoothing operation.

the lack of global context.
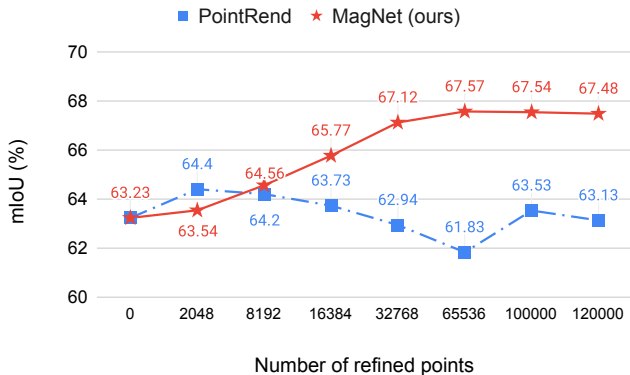
### 4.2.3 Ablation study: point selection



Figure 7: **Correlation between the number of selective points of each scale and mean IoU on the Cityscapes dataset**. When the quantity of points increases, the performance of MagNet continuously grows while the mIoU of PointRend decreases.

Table 6 shows our ablation study on the importance of using the prediction uncertainty maps $Y^u$, $R^u$, and the median filtering function $\mathcal{F}$ in selecting points for refinement.

The best performance was achieved when both uncertainty maps were used. Also, smoothing with median filtering improved the result in every case.

We also studied how the number of refinement points correlates with accuracy. The results of MagNet and PointRend [14] are shown in Fig. 7. As can be seen, the performance of MagNet improved when the number of points increased. The performance stopped increasing after $2^{16}$ points, and it dropped to 66.86% if all points were selected for refinement. Meanwhile, the performance of PointRend decreased significantly when the number of selected points increased beyond 2048; it even dipped below the initial value where no refinement was applied.

### 4.2.4 Ablation study: segmentation backbones

We also tested the MagNet framework with two different segmentation backbone networks, and the results are shown in Table 7. In both cases, MagNet improved the segmentation results of the original networks significantly, from 2% to 5%. In this experiment, we used four scale levels: $256 \rightarrow 512 \rightarrow 1024 \rightarrow 2048$ and the number of refinement lo-

| Model | mIoU(%) |
|---|---:|
| ***Backbone:*** DeepLabV3+ [4] | |
| Patch processing | 59.64 |
| Downsampling | 52.01 |
| MagNet | **61.99** |
| ***Backbone:*** HRNetV2-W48 + OCR [35] | |
| Patch processing | 54.30 |
| Downsampling | 63.92 |
| MagNet | **68.90** |

Table 7: **Results of using MagNet with two backbone networks**. MagNet can be used with different segmentation backbones, and improve their segmentation results. See Table 4 for the results for using MagNet with HRNetV2-W18 backbone.

| Model | mIoU(%) | Mem.(MB) | Time(s) |
|---|---|---|---|
| *Downsampling* | | | |
| U-net[26] | 50.11 | 1813 | |
| FCN-8s[21] | 52.86 | 10569 | |
| SegNet[1] | 60.93 | 2645 | |
| DeepLabv3+[4] | 63.50 | 1541 | |
| FPN[13] | 67.86 | 1247 | 0.01 |
| *Patch processing* | | | |
| U-net[26] | 46.53 | 1813 | |
| FCN-8s[21] | 62.43 | 10569 | |
| SegNet[1] | 68.40 | 2645 | |
| DeepLabv3+[4] | 69.69 | 1541 | |
| FPN[13] | 70.98 | 1247 | 0.31 |
| DenseCRF[15] | 70.36 (↓0.62) | 1247 | 39.68 |
| DGF[31] | 70.38 (↓0.6) | 1435 (↑188) | 0.25 |
| GLNet[5] | 71.60 (↑0.62) | 1865 (↑618) | 0.37 |
| PointRend[14] | 71.78 (↑0.8) | 1593 (↑346) | 0.16 |
| MagNet-Fast | 71.85 (↑0.87) | 1559 (↑312) | 0.29 |
| MagNet | **72.96** (↑1.98) | 1559 (↑312) | 1.19 |

Table 8: **Segmentation results on the DeepGlobe dataset.** We used the same segmentation backbone (FPN) for all refinement methods in the last part.

cations for each patch was $k = 2^{16}$.

## 4.3. DeepGlobe

DeepGlobe is a dataset of high-resolution satellite images. The dataset contains 803 images, annotated with seven landscape classes, including the *unknown* class. Following the evaluation protocol of [5], the *unknown* class is ignored in mIoU calculation, so there are only six classes to consider. The size of the images is $2448\times2448$ pixels. We used the same train/validation/test split as reported in [5], with 455, 207, and 142 images for training, validation, and testing, respectively.

The Feature Pyramid Network (FPN) [13] with Resnet-50 backbone was used as the segmentation network as in the previous work GLNet [5]. We also used the same input size $508\times508$ as GLNet. We used three refinement stages with three scales $612\rightarrow1224\rightarrow2448$ and selected $2^{16}$ points for refinement at each scale. The results are shown in Table 8. For MagNet-Fast, at each of the three scale levels, we selected the top three patches with the highest level of prediction uncertainty for refinement. PointRend [14] was also trained with the same segmentation backbone and it achieved higher accuracy than GLNet. Both MagNet and MagNet-Fast outperformed other methods.

## 4.4. Gleason

Gleason [29] is a medical dataset with histopathological images of prostate cancer. The task is to segment and

| Model | mIoU(%) | Mem.(MB) | Time(s) |
|---|---|---|---|
| Experts | 65.48 | - | - |
| Patching | 46.56 | 1903 | 2.42 |
| Downsampling | 68.90 | 1903 | 0.02 |
| DenseCRF[15] | 69.46 (↑0.56) | 1903 | 141.79 |
| DGF[31] | 68.91 (↑0.01) | 2223 (↑320) | 0.29 |
| PointRend[14] | 68.97 (↑0.07) | 2655 (↑752) | 0.21 |
| MagNet-Fast | 69.75 (↑0.85) | 2621 (↑718) | 0.33 |
| MagNet | **70.60** (↑1.7) | 2621 (↑718) | 2.74 |

Table 9: **Performance of MagNet and other frameworks on Gleason dataset with PSPNet [37] as the backbone.**

grade lesions on ultra-high-resolution images. There are four classes in the dataset that need to be segmented: benign, Grade 3, Grade 4, and Grade 5. The dataset contains 244 images with a size of $5000\times5000$ pixels with segmentation labels provided by six clinical experts. The combined final label is based on majority voting. We randomly split the dataset into 195 training and 49 testing images.

PSPNet [37], the highest-ranked method on the leaderboard for Gleason, was used as the segmentation network with the backbone Resnet-101. We used the input size of $512\times512$, and four refinement stages with four scales: $625\rightarrow1250\rightarrow2500\rightarrow5000$. MagNet-Fast was also run on four scales, but only on three patches with the highest level of prediction uncertainty at each scale. There are $2^{16}$ refinement points for each scale. The results of MagNet and MagNet-Fast, together with the result of the winning solution PSPNet and the mIoU agreement between medical experts, are shown in Table 9. MagNet was run with the PSPNet segmentation backbone, and it improved the performance of PSPNet by 1.7%.

## 5. Conclusions

We have proposed MagNet, a multi-scale segmentation framework for high-resolution images. MagNet can generate high-resolution segmentation output without exploding the GPU memory usage by dividing input images into patches. To avoid the problem of being too global or local, patches of multiple scales are considered, from the coarsest to the finest levels. MagNet has multiple segmentation stages, where the output of one stage will be used as the input for the next stage, and the segmentation output will be progressively refined. We have demonstrated the benefits of MagNet on three challenging high-resolution image datasets, where MagNet outperforms the previous state-of-the-art methods by a margin of 1% to 2% in terms of mean Intersection over Union (mIoU).

## References

[1] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture

for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(12):2481–2495, 2017. 1, 8

[2] Hao Chen, Kunyang Sun, Zhi Tian, Chunhua Shen, Yongming Huang, and Youliang Yan. Blendmask: Top-down meets bottom-up for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2

[3] Liang-Chieh Chen, Yi Yang, Jiang Wang, Wei Xu, and Alan L Yuille. Attention to scale: Scale-aware semantic image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2

[4] 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*, 40(4):834–848, 2017. 1, 2, 7, 8

[5] Wuyang Chen, Ziyu Jiang, Zhangyang Wang, Kexin Cui, and Xiaoning Qian. Collaborative global-local networks for memory-efficient segmentation of ultra-high resolution images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2, 3, 8

[6] Ho Kei Cheng, Jihoon Chung, Yu-Wing Tai, and Chi-Keung Tang. Cascadepsp: Toward class-agnostic and very high-resolution segmentation via global and local refinement. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2, 3

[7] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 5

[8] Jorge Cuadros and George Bresnick. EyePACS: an adaptable telemedicine system for diabetic retinopathy screening. *Journal of Diabetes Science and Technology*, 3(3):509–516, 2009. 1

[9] Ilke Demir, Krzysztof Koperski, David Lindenbaum, Guan Pang, Jing Huang, Saikat Basu, Forest Hughes, Devis Tuia, and Ramesh Raska. Deepglobe 2018: A challenge to parse the earth through satellite images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshop*, 2018. 5

[10] Kaiming He, Jian Sun, and Xiaoou Tang. Guided image filtering. In *Proceedings of the European Conference on Computer Vision*, 2010. 2, 3

[11] Le Hou, Tomas F. Yago Vicente, Minh Hoai, and Dimitris Samaras. Large scale shadow annotation and detection using lazy annotation and stacked cnns. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(4):1337–1351, 2021. 2

[12] Minh-Chuong Huynh, Trung-Hieu Nguyen, and Minh-Triet Tran. Context learning for bone shadow exclusion in chexnet accuracy improvement. In *International Conference on Knowledge and Systems Engineering (KSE)*, 2018. 1

[13] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2, 8

[14] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. Pointrend: Image segmentation as rendering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2, 3, 4, 6, 7, 8

[15] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *Advances in Neural Information Processing Systems*, 2011. 2, 3, 6, 8

[16] Hieu Le, Tomas F. Yago Vicente, Vu Nguyen, Minh Hoai, and Dimitris Samaras. A+D Net: Training a shadow detector with adversarial shadow attenuation. In *Proceedings of the European Conference on Computer Vision*, 2018. 1

[17] Ke Li, Bharath Hariharan, and Jitendra Malik. Iterative instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2, 3

[18] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2

[19] Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan L Yuille, and Li Fei-Fei. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1

[20] Wei Liu, Andrew Rabinovich, and Alexander C Berg. Parsenet: Looking wider to see better. *arXiv preprint arXiv:1506.04579*, 2015. 2

[21] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 1, 8

[22] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014. 1

[23] Vu Nguyen, Tomas F. Yago Vicente, Maozheng Zhao, Minh Hoai, and Dimitris Samaras. Shadow detection with conditional generative adversarial networks. In *Proceedings of the International Conference on Computer Vision*, 2017. 1

[24] Yuval Nirkin, Iacopo Masi, Anh Tran Tuan, Tal Hassner, and Gerard Medioni. On face segmentation, face swapping, and face perception. In *Proceedings of the International Conference on Automatic Face and Gesture Recognition*, 2018. 1

[25] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, 2019. 5

[26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Proceedings of the International Conference on*

*Medical Image Computing and Computer Assisted Intervention*, 2015. 1, 8

[27] Shusuke Takahama, Yusuke Kurose, Yusuke Mukuta, Hiroyuki Abe, Masashi Fukayama, Akihiko Yoshizawa, Masanobu Kitagawa, and Tatsuya Harada. Multi-stage pathological image classification using semantic segmentation. In *Proceedings of the International Conference on Computer Vision*, 2019. 2

[28] Hiroki Tokunaga, Yuki Teramoto, Akihiko Yoshizawa, and Ryoma Bise. Adaptive weighting multi-field-of-view cnn for semantic segmentation in pathology. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2

[29] Danielle Walker. Miccai automatic prostate gleason grading challenge, 2019. 5, 8

[30] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2, 6

[31] Huikai Wu, Shuai Zheng, Junge Zhang, and Kaiqi Huang. Fast end-to-end trainable guided filter. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2, 3, 6, 8

[32] Fangting Xia, Peng Wang, Liang-Chieh Chen, and Alan L Yuille. Zoom better to see clearer: Human and object parsing with hierarchical auto-zoom net. In *Proceedings of the European Conference on Computer Vision*, 2016. 2

[33] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *Proceedings of the European Conference on Computer Vision*, 2018. 2

[34] Changqian Yu, Jingbo Wang, Changxin Gao, Gang Yu, Chunhua Shen, and Nong Sang. Context prior for scene segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2

[35] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation. In *Proceedings of the European Conference on Computer Vision*, 2020. 2, 5, 6, 7

[36] Yuhui Yuan, Jingyi Xie, Xilin Chen, and Jingdong Wang. Segfix: Model-agnostic boundary refinement for segmentation. In *Proceedings of the European Conference on Computer Vision*, 2020. 2, 3, 6

[37] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2, 8

[38] Hengshuang Zhao, Xiaojuan Qi, Xiaoyong Shen, Jianping Shi, and Jiaya Jia. Icnet for real-time semantic segmentation on high-resolution images. In *Proceedings of the European Conference on Computer Vision*, 2018. 2

[39] Peng Zhou, Brian Price, Scott Cohen, Gregg Wilensky, and Larry S Davis. Deepstrip: High-resolution boundary refinement. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2, 3