# An End-to-End Network for Panoptic Segmentation

Huanyu Liu[1][†], Chao Peng[2], Changqian Yu[3][†], Jingbo Wang[4][†], Xu Liu[5][†], Gang Yu[2], Wei Jiang[1]

[1]Zhejiang University, [2]Megvii Inc. (Face++), [3]Huazhong University of Science and Technology
[4]Peking University, [5]The University of Tokyo

{liuhy, jiangwei_zju}@zju.edu.cn, mikejay0520@163.com, changqian_yu@hust.edu.cn
wangjingbo1219@pku.edu.cn, liuxu@kmj.iis.u-tokyo.ac.jp, yugang@megvii.com

## Abstract

*Panoptic segmentation, which needs to assign a category label to each pixel and segment each object instance simultaneously, is a challenging topic. Traditionally, the existing approaches utilize two independent models without sharing features, which makes the pipeline inefficient to implement. In addition, a heuristic method is usually employed to merge the results. However, the overlapping relationship between object instances is difficult to determine without sufficient context information during the merging process. To address the problems, we propose a novel end-to-end Occlusion Aware Network (OANet) for panoptic segmentation, which can efficiently and effectively predict both the instance and stuff segmentation in a single network. Moreover, we introduce a novel spatial ranking module to deal with the occlusion problem between the predicted instances. Extensive experiments have been done to validate the performance of our proposed method and promising results have been achieved on the COCO Panoptic benchmark.*

## 1. Introduction

Panoptic segmentation [18] is a new challenging topic for scene understanding. The goal is to assign each pixel with a category label and segment each object instance in the image. In this task, the stuff segmentation is employed to predict the amorphous regions (noted *Stuff*) while the instance segmentation [14] solves the countable objects (noted *Thing*). Therefore, this task can provide more comprehensive scene information and can be widely used in autonomous driving and scene parsing.

Previous panoptic segmentation algorithms [18] usually contain three separated components: the instance segmentation block, the stuff segmentation block and the merging block, as shown in Figure 1 (a). Usually in these algorithms, the instance and stuff segmentation blocks are independent

---
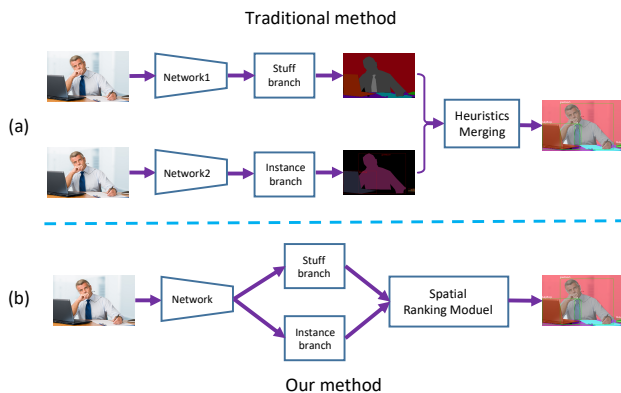[†]This work was done during an internship at Megvii Inc.



Figure 1. An illustration of our end-to-end network in contrast to the traditional method. Traditional methods [18] train two sub-network and do heuristics merge. Our method can train a single network for two subtasks, and realize a learnable fusion approach.

without any feature sharing. This results in apparent computational overhead. Furthermore, because of the separated models, these algorithms have to merge the corresponding separated predictions with post-processing. However, without the context information between the stuff and thing, the merge process will face the challenge of overlapping relationships between instances and stuff. Exactly as mentioned above, with the three separate parts, it is hard to apply this complex pipeline to the industrial application.

In this paper, we propose a novel end-to-end algorithm shown in Figure 1 (b). As far as we know, this is the first algorithm which can deal with the issues above in an end-to-end pipeline. More specifically, we incorporate the instance segmentation and stuff segmentation into one network, which shares the backbone features but applies different head branches for the two tasks. During the training phase, the backbone features will be optimized by the accumulated losses from both the stuff and thing supervision while the head branches will only fine-tune on the specific task.

To solve the problem of overlapping relationship between object instances, we also introduce a new algorithm called *Spatial Ranking Module*. This module learns the ranking score and offers an ordering accordance for instances.

In general, we summarize the contributions of our algorithm as follows:

- We are the first to propose an end-to-end occlusion aware pipeline for the problem of panoptic segmentation.

- We introduce a novel spatial ranking module to address the ambiguities of the overlapping relationship, which commonly exists in the problem of panoptic segmentation.

- We obtain state-of-the-art performance on the COCO panoptic segmentation dataset.

## 2. Related Work

### 2.1. Instance Segmentation

There are currently two main frameworks for instance segmentation, including the proposal-based methods and segmentation-based methods. The proposal-based approaches [8, 14, 24, 25, 28, 29, 33] first generate the object detection bounding boxes and then perform mask prediction on each box for instance segmentation. These methods are closely related to object detection algorithms such as Fast/Faster R-CNN and SPPNet [12, 15, 36]. Under this framework, the overlapping problem raises due to the independence prediction of distinct instances. That is, pixels may be allocated to wrong categories when covered by multiple masks. The segmentation-based methods use the semantic segmentation network to predict the pixel class, and obtain each instance mask by decoding the object boundary [19] or the custom field [2, 9, 27]. Finally, they use the bottom-up grouping mechanism to generate the object instances. RNN method was leveraged to predict a mask for each instance at a time in [35, 37, 46] .

### 2.2. Semantic Segmentation

Semantic segmentation has been extensively studied, and many new methods have been proposed in recent years. Driven by powerful deep neural networks [16, 21, 39, 40], FCN [30] successfully applied deep learning to pixel-by-pixel image semantic segmentation by replacing the fully connected layer of the image classification network with the convolutional layer.

The encoder-decoder structure such as UNet [1, 38, 42, 43, 47] can gradually recover resolution and capture more object details. Global Convolutional Network [34] proposes large kernel method to relieve the contradiction between

classification and localization. DFN [43] designs a channel attention block to select feature maps. DeepLab [4, 6] and PSPNet [48] use atrous spatial pyramid pooling or spatial pyramid pooling to get multi-scale context. Method of [44] uses dilated convolution to enlarge field of view. Multi-scale features were using to obtain sufficient receptive field in [5, 13, 41] .

Related datasets are also constantly being enriched and expanded. Currently, there are public datasets such as VOC [11], Cityscapes [7], ADE20K [49], Mapillary Vistas [32], and COCO Stuff [3].

### 2.3. Panoptic Segmentation

Panoptic Segmentation task was first proposed in [18] and the research work for this task is not too much. A weakly supervised model that jointly performs semantic and instance segmentation was proposed by [22]. It uses the weak bounding box annotations for "thing" classes, and image level tags for "stuff" classes. JSIS-Net [10] proposes a single network with the instance segmentation head [14] and the pyramid stuff segmentation head [48], following heuristics to merge two kinds of outputs. Li et al. [23] propose AUNet that can leverage proposal and mask level attention and get better background results.

### 2.4. Multi-task learning

Panoptic segmentation can also be treated as multi-task learning problem. Two different task can be trained together through strategies [17, 31]. UberNet [20] jointly handles low-, mid-, and high-level vision tasks in a single network, including boundary detection, semantic segmentation and normal estimation. Zamir et al. [45] build a directed graph named taskonomy, which can effectively measure and leverage the correlation between different visual tasks. It can avoid repetitive learning and enable learning with less data.

## 3. Proposed End-to-end Framework

The overview of our algorithm is illustrated in Figure 2. There are three major components in our algorithm: 1) The stuff branch predicts the stuff segmentation for the whole input. 2) The instance branch provides the instance segmentation predictions. 3) The spatial ranking module generates a ranking score for the each instance.

### 3.1. End-to-end Network Architecture

We employ FPN [26] as the backbone architecture for the end-to-end network. For instance segmentation, we adopt the original Mask R-CNN [14] as our network framework. We apply top-down pathway and lateral connections to get feature maps. And then, a $3 \times 3$ convolution layer is appended to get RPN feature maps. After that, we apply the
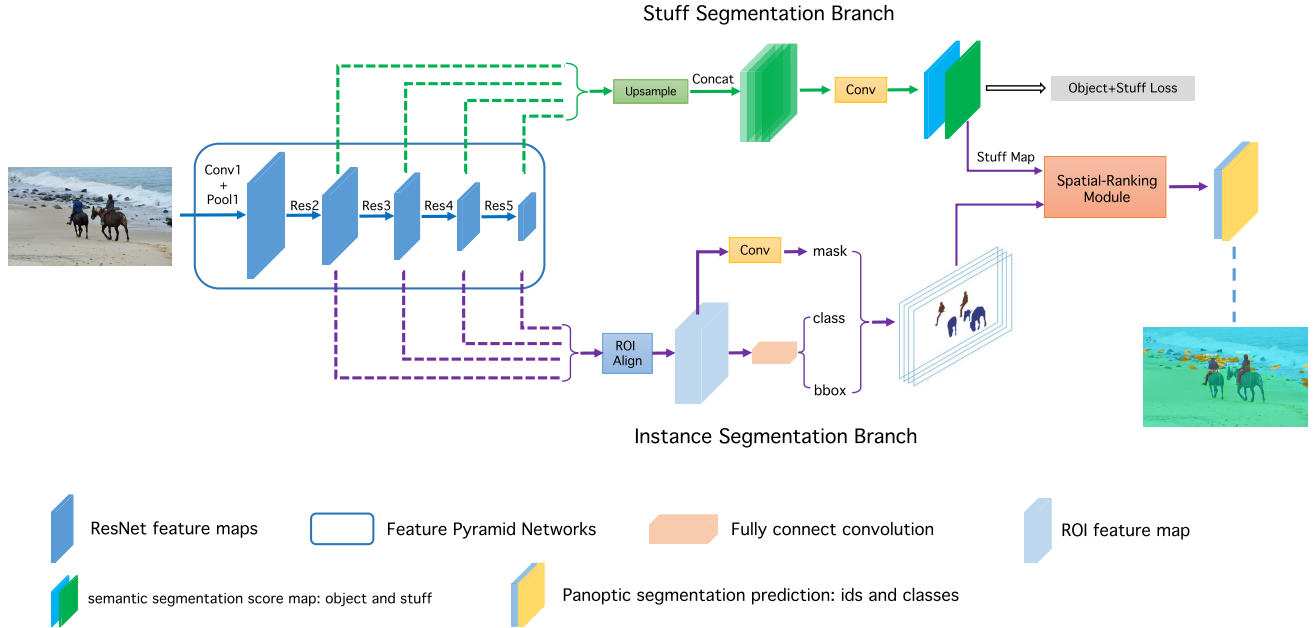
Figure 2. The illustration of overall framework. Given an input image, we use the FPN network to provide feature maps for stuff branch and instance branch. The two branches generate intermediate results, while later are passed to our spatial ranking module. The spatial ranking module learns a ranking score for each instance as the final merging evidence.

ROIAlign [14] layer to extract object proposal features and get three predictions: proposal classification score, proposal bounding box coordinates, and proposal instance mask.

For stuff segmentation, two $3 \times 3$ convolution layers are stacked on RPN feature maps. For the sake of multi-scale feature extraction, we then concatenate these layers with succeeding one $3 \times 3$ convolution layer and $1 \times 1$ convolution layer. Figure 3 presents the details of stuff branch. During training, we supervise the stuff segmentation and thing segmentation simultaneously, as the auxiliary objection information could provide object context for stuff prediction. In inference, we only extract the stuff predictions and normalized them to probability.

To break out the information flow barrier during training and to make the whole pipeline more efficient, we share the features from the backbone network of two branches. The issue raised here could be divided into two parts: 1) the sharing granularity on feature maps and 2) the balance between instance loss and stuff loss. In practice, we find that as more feature maps are shared, better performance we can obtain. Thus, we share the feature maps until skip-connection layers, that is the $3 \times 3$ convolution layer before RPN head shown in Figure 3.

$$
\begin{aligned}
L_{\text{total}} = &\underbrace{L_{\text{rpn\_cls}} + L_{\text{rpn\_bbox}} + L_{\text{cls}} + L_{\text{bbox}} + L_{\text{mask}}}_{\text{instance branch}} \\
&+ \underbrace{\lambda \cdot L_{\text{seg\_(stuff+object)}}}_{\text{stuff branch}} + L_{\text{srm}}
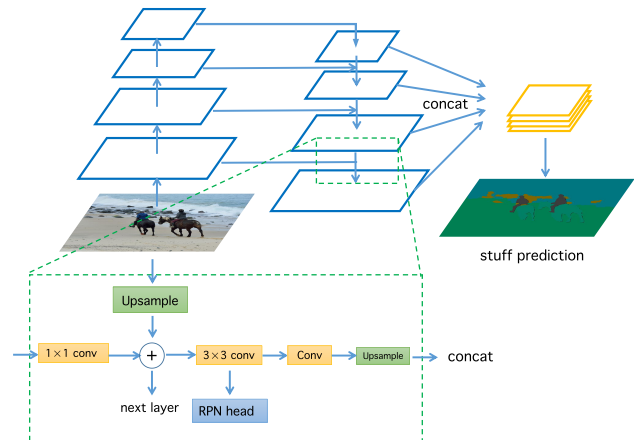\end{aligned}
\tag{1}
$$



Figure 3. A building block illustrating the stuff segmentation sub-network. Here we share both the backbone and skip-connection feature maps in stuff branch and instance branch. Besides, we predict both object and stuff category for the stuff branch.

As for the balance of two supervisions, we first present the multiple losses in Equation 1. The instance branch contains 5 losses: $L_{rpn\_cls}$ is the RPN objectness loss, $L_{rpn\_bbox}$ is the RPN bounding-box loss, $L_{cls}$ is the classification loss, $L_{bbox}$ is the object bounding-box regression loss, and $L_{mask}$ is the average binary cross-entropy loss for mask prediction. As the stuff branch, there is only one se-
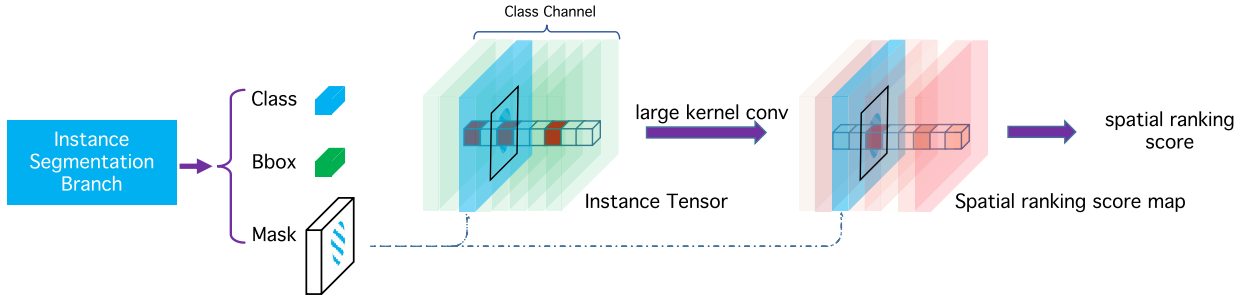
3

Figure 4. An illustration of spatial ranking score map prediction. The pixel vector in instance feature map represents instance prediction result in this pixel. The red color means that the corresponding category object includes this pixel and the multiple red channels indicate the occlusion problem between instances. We use the panoptic segmentation category label to supervise spatial ranking score map.

mantic segmentation loss named $L_{seg\_(stuff+object)}$. The hyperparameter $\lambda$ is employed for loss balance and will be discussed later. $L_{\mathrm{srm}}$ represents the loss function of the spatial ranking module, which is described in the next section.

## 3.2. Spatial Ranking Module

The modern instance segmentation framework is often based on object detection network with an additional mask prediction branch, such as the Mask RCNN [14] which is usually based on FPN [26]. Generally speaking, the current object detection framework does not consider the overlapping problems between different classes, since the popular metrics are not affected by this issue, e.g., the AP and AR. However, in the task of panoptic segmentation, since the number of pixels in one image is fixed, the overlapping problem, or specifically the multiple assignments for one pixel, must be resolved.

Commonly, the detection score was used to sort the instances in descending order, and then assign them to the stuff canvas by the rule of larger score objects on top of lower ones. However, this heuristic algorithm could easily fail in practice. For example, let's consider a person wearing a tie, shown in Figure 7. As the person class is more frequent than the tie in the COCO dataset, its detection score is tend to be higher than the tie bounding box. Thus through the above simple rule, the tie instance is covered by the person instance, and leading to the performance drops.

Could we alleviate this phenomenon through the panoptic annotation? That is if we force the network learns the person annotation with a hole in the place of the tie, could we avoid the above situation? As shown in Table 3, we conduct the experiment with the above mentioned annotations, but only find the decayed performance. Therefore, this approach is not applicable currently.

To counter this problem, we resort to a semantic-like approach and propose a simple but very effective algorithm for dealing with occlusion problems, called *spatial ranking module*. As shown in the Figure 4, we first map the results

of the instance segmentation to the tensor of input size. The dimension of the feature map is the number of object categories, and the instances of different categories are mapped to the corresponding channels. The instance tensor is initialized to zero, and the mapping value is set to one. We then append the large kernel convolution [34] after the tensor to obtain the ranking score map. In the end , we use the pixel-wise cross entropy loss to optimize the ranking score map, as the Equation 2 shows. $S_{map}$ represents the output ranking score map and $S_{label}$ represents the corresponding non-overlap semantic label.

$$L_{\mathrm{srm}} = CE(S_{map}, S_{label}) \tag{2}$$

After getting the ranking score map, we calculate the ranking score of each instance object as Equation 3. Here, $S_{i,j,cls}$ represents the ranking score value in $(i, j)$ of class $cls$, note that $S_{i,j,cls}$ has been normalized to a probability distribution. $m_{i,j}$ is the mask indicator, representing if pixel $(i, j)$ belongs to the instance. The ranking score of the whole instance $P_{objs}$ is computed by the average of pixel ranking scores in a mask.

$$P_{objs} = \frac{\sum_{(i,j) \in objs} S_{i,j,cls} \cdot m_{i,j}}{\sum_{(i,j) \in objs} m_{i,j}} \tag{3}$$

$$m_{i,j} = \begin{cases} 0 & \text{(i,j)} \in \text{instance} \\ 1 & \text{(i,j)} \notin \text{instance} \end{cases} \tag{4}$$

Let's reconsider the example mentioned above through our proposed spatial ranking module. When we forward the person mask and tie mask into this module, we obtain the spatial ranking scores using Equation 3 for these two objects. Within the ranking score, the sorting rule in the previous method could be more reliable, and the performance is improved, as the experiments present in the next section.

4

# 4. Experiments

## 4.1. Dataset and Evaluation Metrics

**Dataset:** We conduct all experiments on COCO panoptic segmentation dataset [18]. This dataset contains 118K images for training, 5k images for validation, with annotations on 80 categories for the thing and 53 classes for stuff. We only employ the training images for model training and test on the validation set. Finally, we submit test-dev result to the COCO 2018 panoptic segmentation leaderboard.

**Evaluation metrics:** We use the standard evaluation metric defined in [18], called Panoptic Quality (PQ). It contains two factors: 1) the Segmentation Quality (SQ) measures the quality of all categories and 2) the Detection Quality (DQ) measures only the instance classes. The mathematical formations of PQ, SQ and DQ are presented in Equation 5, where p and g are predictions and ground truth, and TP, FP, FN represent true positives, false positives and false negatives. It is easy to find that SQ is the common mean IOU metric normalized for matching instances, and DQ could be regarded as a form of detection accuracy. The matching threshold is set to 0.5, that is if the pixel IOU of prediction and ground truth is larger than 0.5, the prediction is regarded matched, otherwise unmatched. For stuff classes, each stuff class in an image is regarded as one instance, no matter the shape of it.

$$PQ = \underbrace{\frac{\sum_{(p,g)\in TP} IOU(p,g)}{|TP|}}_{\text{Segmentation Quality (SQ)}} \times \underbrace{\frac{|TP|}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}}_{\text{Detection Quality (DQ)}}$$

(5)

## 4.2. Implementation Details

We choose the ResNet-50 [16] pretrained on ImageNet for ablation studies. We use the SGD as the optimization algorithm with momentum 0.9 and weight decay 0.0001. The multi-stage learning rate policy with warm up strategy [33] is adopted. That is, in the first $2,000$ iterations, we use the linear gradual warmup policy by increasing the learning rate from 0.002 to 0.02. After $60,000$ iterations, we decrease the learning rate to 0.002 for the next $20,000$ iterations and further set it to 0.0002 for the rest $20,000$ iterations. The batch size of input is set to 16, which means each GPU consumes two images in one iteration. For other details, we employ the experience from Mask-RCNN [14].

Besides the training for the two branches of our network, a little bit more attention should be paid to the spatial ranking module. During the training process, the supervision label is the corresponding non-overlap semantic label and training it as a semantic segmentation network. We set the non-conflicting pixels ignored to force the network focus on the conflicting regions.

During the inference, we set the max number of boxes for each image to 100 and the min area of a connected stuff area to 4,900. As for the spatial ranking module, since we do not have the ground truth now, the outputs of instance branch will be passed through this module to resolve the overlapping issue.

## 4.3. Ablation Study on Network Structure

In this subsection, we focus on the properties of our end-to-end network design. There are three points should be discussed as follows: the loss balance parameter, the object context for stuff branch and the sharing mode of two branches. To avoid the Cartesian product of experiments, we only modify the specific parameters and control the other optimally.
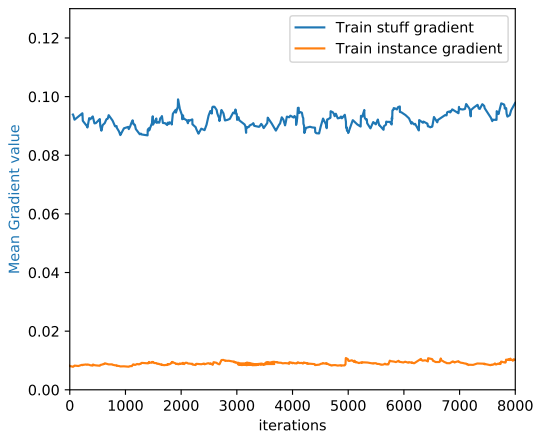


Figure 5. The plot of the specified layer mean gradient value in two branches. We chose one epoch iterations in the training process. The learning rate of the two branches is the same. The horizontal axis is the number of iterations, and the vertical axis is the mean gradient value of the backbone last layer.

| $\lambda$ | PQ | PQ$^{\text{Th}}$ | PQ$^{\text{St}}$ |
|---|---|---|---|
| 0.2 | 36.9 | 45.0 | 24.6 |
| 0.25 | **37.2** | 45.4 | 24.9 |
| 0.33 | 36.9 | 44.4 | 25.4 |
| 0.50 | 36.5 | 43.5 | 25.9 |
| 0.75 | 35.3 | 41.9 | 25.4 |
| 1.0 | - | - | - |

Table 1. Loss balance between instance segmentation and stuff segmentation

**The loss balance** issue comes from the reality that the gradients from stuff branch and instance branch are not close. We make a statistic on the mean gradients of two branches with respect to the last feature map of the backbone, where the hyperparameter $\lambda$ is set to 1 for fairness.

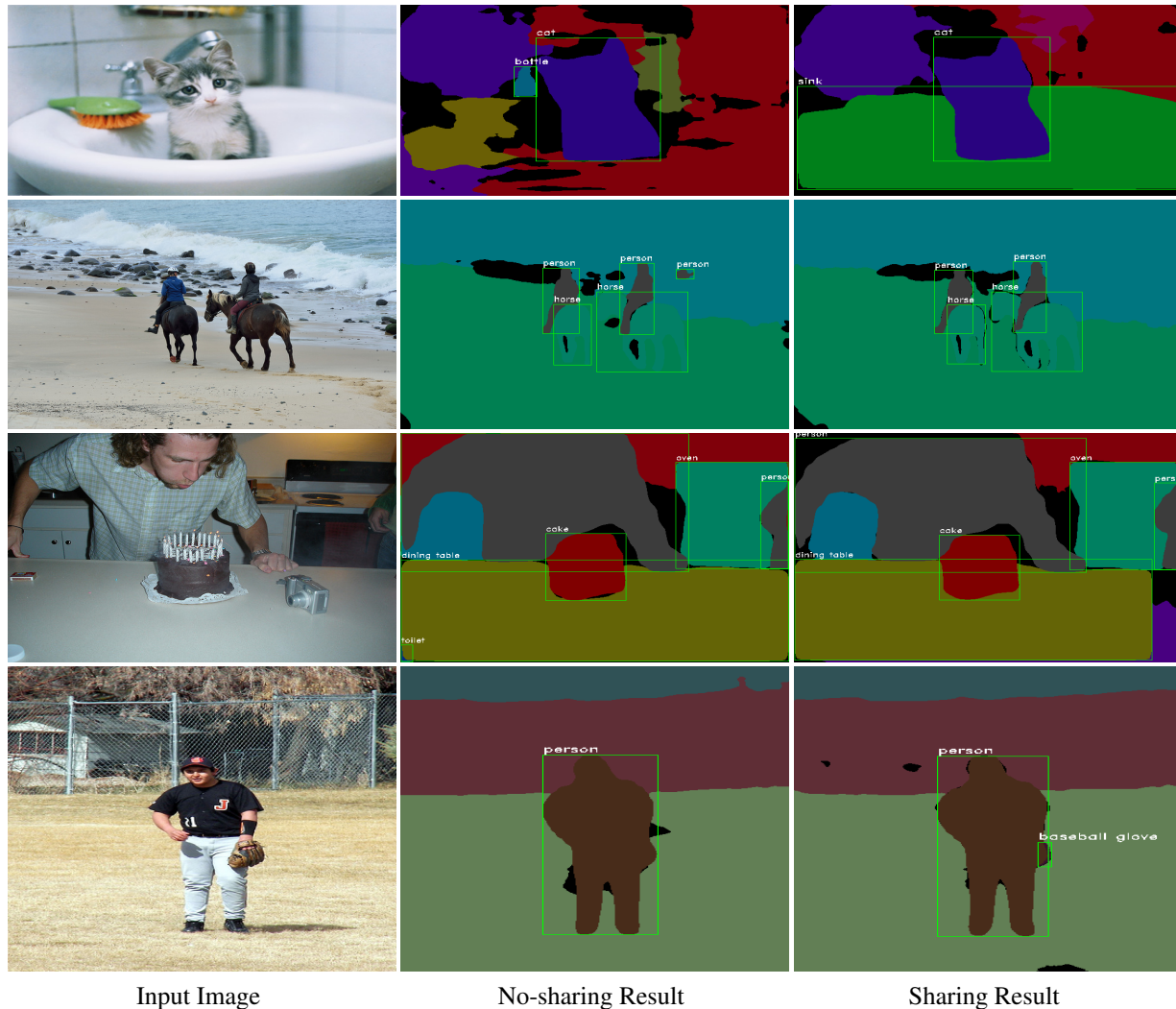|  | Input Image | No-sharing Result | Sharing Result |

Figure 6. Feature Sharing Mode Visualization. The first column is the original image, and the second column is sharing backbone feature result, and the last column is the no-sharing result.

As shown in the Figure 5, it is brief and clear that the gradient from stuff branch dominates the penalty signals. Therefore, we obtain a hyperparameter $\lambda$ in Equation 1 to balance the gradients. We conduct the experiments with $\lambda \in [0.2, 0.25, 0.33, 0.5, 0.75, 1.0]$. The interval is not uniform for the sake of searching efficiency. As the Table 1 summarizes, $\lambda = 0.25$ is the optimal choice. Note that if we set $\lambda = 1$, which means the instance and stuff branches trained like separate models, the network could not converge through our default learning rate policy.

**Object context** is a natural choice in stuff segmentation. Although we only want the stuff predictions from this branch, the lack of object supervision will introduce holes on ground truth, resulting in discontinuous context around objects. Therefore, we conduct a pair of comparative experiments where all 133 categories is supervised , and the other

| Stuff-SC | Object-SC | PQ | PQ$^{\text{Th}}$ | PQ$^{\text{St}}$ |
|:---:|:---:|:---:|:---:|:---:|
| ✓ | - | 36.7 | 43.8 | 25.9 |
| ✓ | ✓ | **37.2** | 45.4 | 24.9 |

Table 2. Ablation study results on stuff segmentation network design. Stuff-SC represents stuff supervision classes. It refers to predict stuff classes. While both Stuff-SC and Object-SC mean predicting all classes.

is trained on 53 stuff classes. The results in Table 2 shows the 0.5 improvement on overall PQ with object context.

**Sharing features** is a key point of our network design. The benefits of sharing have two parts: 1) two branches could absorb useful information from other supervision signals, and 2) the computation resources could be saved if the shared network is computed only once. To investigate

| Methods | backbone | PQ | PQ$^{Th}$ | PQ$^{St}$ | SQ | SQ$^{Th}$ | SQ$^{St}$ | DQ | DQ$^{Th}$ | DQ$^{St}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| baseline | ResNet-50 | 37.2 | 45.4 | 24.9 | 77.1 | 81.5 | 70.6 | 45.7 | 54.4 | 32.5 |
| w/pano-instance GT | ResNet-50 | 36.1 | 43.5 | 24.9 | 76.1 | 80.0 | 70.3 | 44.5 | 52.4 | 32.7 |
| w/spatial ranking module | ResNet-50 | **39.0** | 48.3 | 24.9 | 77.1 | 81.4 | 70.6 | 47.8 | 58.0 | 32.5 |
| baseline | ResNet-101 | 38.8 | 46.9 | 26.6 | 78.2 | 82.0 | 72.5 | 47.4 | 55.9 | 34.5 |
| w/spatial ranking module | ResNet-101 | **40.7** | 50.0 | 26.6 | 78.2 | 82.0 | 72.5 | 49.6 | 59.7 | 34.5 |

Table 3. Results on MS-COCO panoptic segmentation validation dataset which use our spatial ranking module method. W/pano-instance GT represents using panoptic segmentation ground truth to generate instance segmentation ground truth. It is trained in two separate networks. All results in this table are based on backbone ResNet-50.

| Methods | backbone | PQ | PQ$^{Th}$ | PQ$^{St}$ |
|---|---|---|---|---|
| no | ResNet-50 | 36.5 | 44.4 | 24.6 |
| res1-res5 | ResNet-50 | 37.0 | 44.8 | 25.2 |
| + skip-conection | ResNet-50 | **37.2** | 45.4 | 24.9 |
| no | ResNet-101 | 38.2 | 46.3 | 26.0 |
| + skip-conection | ResNet-101 | **38.8** | 46.9 | 26.6 |

Table 4. Results on whether share stuff segmentation and instance segmentation features. On ResNet-50 backbone, sharing features method gets a gain of 0.7 in PQ, and ResNet-101 gets a gain of 0.7. Ablation study results on different sharing feature way. The res1-res5 means just share the backbone ResNet features. The +skip-connection means share both the backbone features and FPN skip-connection branch.

| Conv Settings | PQ | PQ$^{Th}$ | PQ$^{St}$ |
|---|---|---|---|
| $1 \times 1$ | 38.4 | 47.4 | 24.9 |
| $3 \times 3$ | 38.7 | 47.8 | 24.9 |
| $1 \times 7 + 7 \times 1$ | **39.0** | 48.3 | 24.9 |

Table 5. Results on the convolution settings of spatial ranking module. $1 \times 1$ represents the convolution kernel size is 1. Results shows that the large receptive field can help the spatial ranking module get more context features and better results.

the granularity on sharing features, we conduct two experiments, where the *shallow share model* only shares the backbone features and the *deep share model* further shares the feature maps before RPN head, as Figure 3 presents. Table 4 shows the comparisons between different settings, and *deep share model* outperform the separate training baseline by 0.7% on PQ. Figure 6 presents the visualization of sharing features.

## 4.4. Ablation Study on Spatial Ranking Module

**Supervised by no-overlapping annotations** are a straightforward idea to resolve the object ranking issue. Here, we process the panoptic ground truth and extract the non-overlapping annotations for instance segmentation. The 3rd row of Table 3 gives the results of this idea. Unfortunately, merely replacing the instance ground truth do not help improve the performance, and may reversely reduce

| Methods | PQ | PQ$^{Th}$ | PQ$^{St}$ |
|---|---|---|---|
| Artemis | 16.9 | 16.8 | 17.0 |
| JSIS-Net [10] | 27.2 | 29.6 | 23.4 |
| MMAP-seg | 32.1 | 38.9 | 22.0 |
| LeChen | 37.0 | 44.8 | 25.2 |
| Ours(OANet) | **41.3** | 50.4 | 27.7 |

Table 6. Results on the COCO 2018 panoptic segmentation challenge test-dev. Results verifies the effectiveness of our feature sharing mode and the spatial ranking module. We use the ResNet-101 as our basemodel.

the accuracy and recall for objects. This phenomenon may come from the fact that most of the objects in COCO do not meet the overlapping issue, and forcing the network to learn non-overlapping hurts the overall performance.
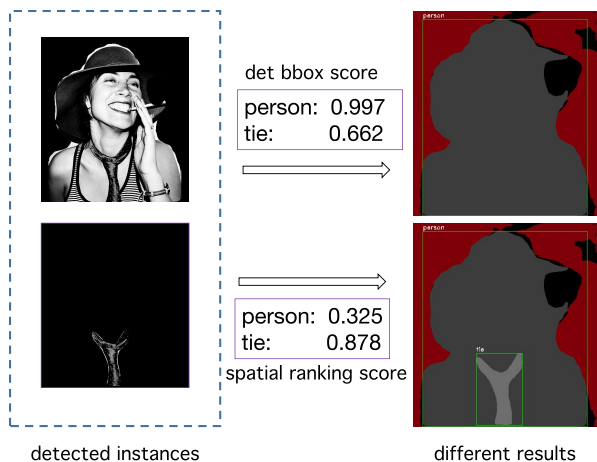


Figure 7. The visualization result of the spatial ranking module. The left two segments denote the detected instances through our model, the det bbox score represents the object detection score predicted in detection. The spatial ranking score represents the values from our approach.

**Spatial ranking module** proposed in this paper is aimed to solve the overlapping issue in panoptic segmentation. Table 5 shows that large receptive field can help the spatial

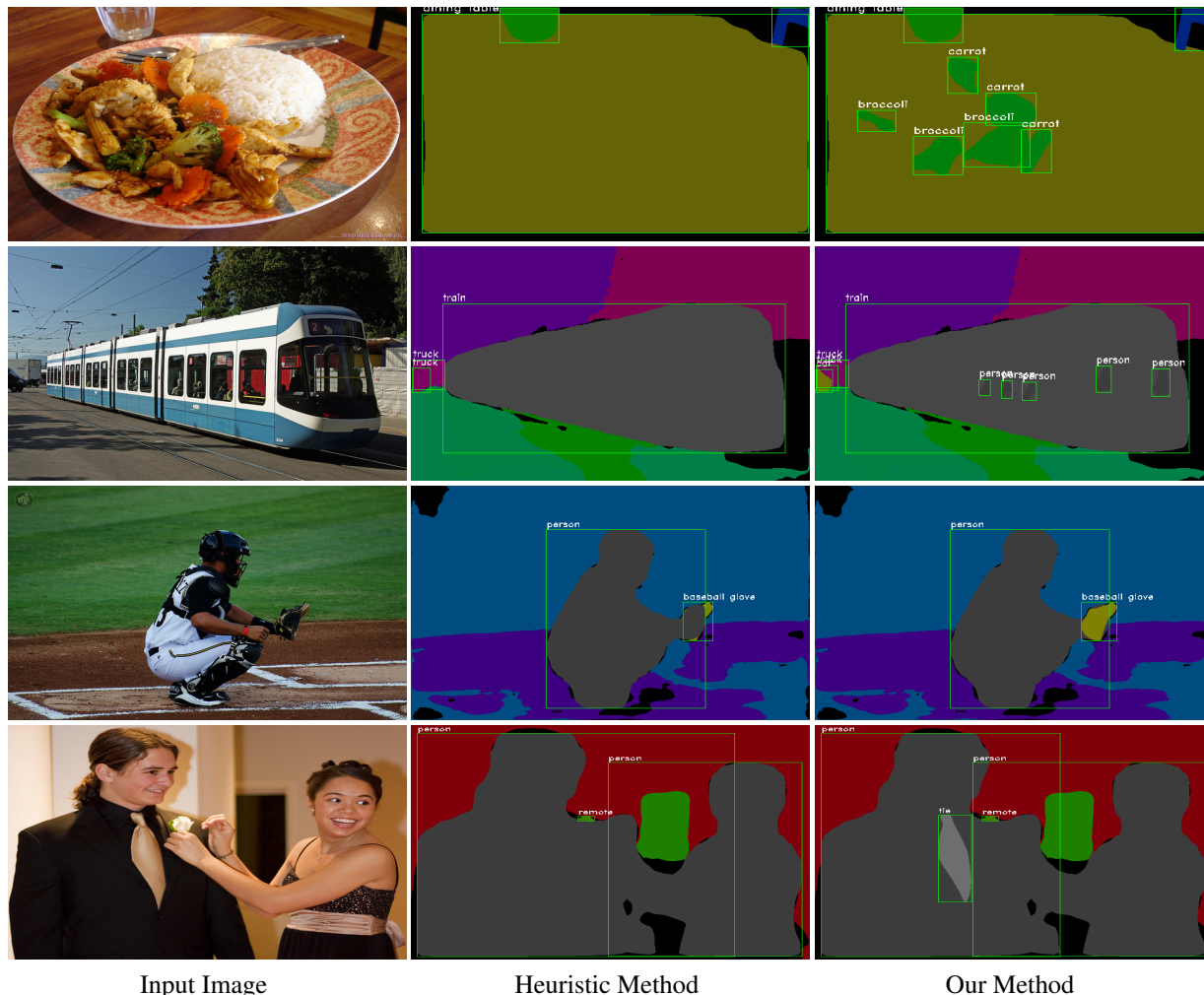|  |  |  |
|---|---|---|
| Input Image | Heuristic Method | Our Method |

Figure 8. The visualization results using our spatial ranking module. The first column is the input image, and second column is the panoptic segmentation result with heuristic method, and last column is the result using our approach.

ranking module get more context features and better results. As we can see in 3rd row or 2nd row for the Resnet101 of Table 3, compared with the above end-to-end baseline, our spatial ranking module improve the PQ by 1.8%. Specifically, the PQ$^{Th}$ is increased by 2.9%, while the metrics for stuff remains the same. These facts prove the purpose of our spatial ranking module is justified. We test our OANet on COCO test-dev, as shown in Table 6. Compared with the results of other methods, our method achieves the state-of-the-art result. For detailed results, please refer to the leardboard [‡].

Figure 7 explaines the principle of our spatial ranking module. For the example input image, the network predicts a person plus a tie, and their bounding box scores are 0.997 and 0.662 respectively. If we use the scores to decide the results, the tie will definitely be covered by the person. How-

ever, in our method, we can get a spatial ranking score for each instance, 0.325 and 0.878 respectively. With the help of the new scores, we can get the right predictions. Figure 8 summarizes more examples.

## 5. Conclusion

In this paper, we propose a novel end-to-end occlusion aware algorithm, which incorporates the common semantic segmentation and instance segmentation into a single model. In order to better employ the different supervisions and reduce the consumption of computation resources, we investigate the feature sharing between different branches and find that we should share as many features as possible. Besides, we have also observed the particular ranking problem raised in the panoptic segmentation, and design the simple but effective spatial ranking module to deal with this issue. The experiment results show that our approach outperforms the previous state-of-the-art models.

# References

[1] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE TPAMI*, 2017.

[2] M. Bai and R. Urtasun. Deep watershed transform for instance segmentation. In *CVPR*, 2017.

[3] H. Caesar, J. Uijlings, and V. Ferrari. Coco-stuff: Thing and stuff classes in context. In *CVPR*, 2018.

[4] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE TPAMI*, 2018.

[5] L.-C. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille. Attention to scale: Scale-aware semantic image segmentation. In *CVPR*, 2016.

[6] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018.

[7] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.

[8] J. Dai, K. He, and J. Sun. Instance-aware semantic segmentation via multi-task network cascades. In *CVPR*, 2016.

[9] B. De Brabandere, D. Neven, and L. Van Gool. Semantic instance segmentation with a discriminative loss function. *arXiv:1708.02551*, 2017.

[10] D. de Geus, P. Meletis, and G. Dubbelman. Panoptic segmentation with a joint semantic and instance segmentation network. *arXiv:1809.02110*, 2018.

[11] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. In *IJCV*, 2015.

[12] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.

[13] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Hypercolumns for object segmentation and fine-grained localization. In *CVPR*, 2015.

[14] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *ICCV*, 2017.

[15] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *ECCV*, 2014.

[16] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

[17] A. Kendall, Y. Gal, and R. Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *CVPR*, 2018.

[18] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár. Panoptic segmentation. *arXiv:1801.00868*, 2018.

[19] A. Kirillov, E. Levinkov, B. Andres, B. Savchynskyy, and C. Rother. Instancecut: from edges to instances with multicut. In *CVPR*, 2017.

[20] I. Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In *CVPR*, 2017.

[21] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.

[22] Q. Li, A. Arnab, and P. H. Torr. Weakly-and semi-supervised panoptic segmentation. In *ECCV*, 2018.

[23] Y. Li, X. Chen, Z. Zhu, L. Xie, G. Huang, D. Du, and X. Wang. Attention-guided unified network for panoptic segmentation. *arXiv:1812.03904*, 2018.

[24] Y. Li, H. Qi, J. Dai, X. Ji, and Y. Wei. Fully convolutional instance-aware semantic segmentation. In *CVPR*, 2017.

[25] Z. Li, C. Peng, G. Yu, X. Zhang, Y. Deng, and J. Sun. Detnet: Design backbone for object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 334–350, 2018.

[26] T.-Y. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017.

[27] S. Liu, J. Jia, S. Fidler, and R. Urtasun. Sgn: Sequential grouping networks for instance segmentation. In *ICCV*, 2017.

[28] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia. Path aggregation network for instance segmentation. In *CVPR*, 2018.

[29] S. Liu, X. Qi, J. Shi, H. Zhang, and J. Jia. Multi-scale patch aggregation (mpa) for simultaneous detection and segmentation. In *CVPR*, 2016.

[30] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.

[31] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert. Crossstitch networks for multi-task learning. In *CVPR*, 2016.

[32] G. Neuhold, T. Ollmann, S. R. Bulò, and P. Kontschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *ICCV*, 2017.

[33] C. Peng, T. Xiao, Z. Li, Y. Jiang, X. Zhang, K. Jia, G. Yu, and J. Sun. Megdet: A large mini-batch object detector. In *CVPR*, 2018.

[34] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun. Large kernel mattersimprove semantic segmentation by global convolutional network. In *CVPR*, 2017.

[35] M. Ren and R. S. Zemel. End-to-end instance segmentation with recurrent attention. In *CVPR*, 2017.

[36] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.

[37] B. Romera-Paredes and P. H. S. Torr. Recurrent instance segmentation. In *ECCV*, 2016.

[38] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.

[39] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*, 2014.

[40] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.

[41] F. Xia, P. Wang, L.-C. Chen, and A. L. Yuille. Zoom better to see clearer: Human and object parsing with hierarchical auto-zoom net. In *ECCV*, 2016.

[42] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *ECCV*, 2018.

[43] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang. Learning a discriminative feature network for semantic segmentation. In *CVPR*, 2018.

[44] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv:1511.07122*, 2015.

[45] A. R. Zamir, A. Sax, W. Shen, L. J. Guibas, J. Malik, and S. Savarese. Taskonomy: Disentangling task transfer learning. In *CVPR*, 2018.

[46] Z. Zhang, A. G. Schwing, S. Fidler, and R. Urtasun. Monocular object instance segmentation and depth ordering with cnns. In *CVPR*, 2015.

[47] Z. Zhang, X. Zhang, C. Peng, X. Xue, and J. Sun. Exfuse: Enhancing feature fusion for semantic segmentation. In *ECCV*, 2018.

[48] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *CVPR*, 2017.

[49] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Scene parsing through ade20k dataset. In *CVPR*, 2017.