

# You Only Segment Once: Towards Real-Time Panoptic Segmentation

Jie Hu, Linyan Huang, Tianhe Ren, Shengchuan Zhang, Rongrong Ji, and Liujuan Cao\*  
 Key Laboratory of Multimedia Trusted Perception and Efficient Computing,  
 Ministry of Education of China, Xiamen University

## Abstract

*In this paper, we propose YOSO, a real-time panoptic segmentation framework. YOSO predicts masks via dynamic convolutions between panoptic kernels and image feature maps, in which you only need to segment once for both instance and semantic segmentation tasks. To reduce the computational overhead, we design a feature pyramid aggregator for the feature map extraction, and a separable dynamic decoder for the panoptic kernel generation. The aggregator re-parameterizes interpolation-first modules in a convolution-first way, which significantly speeds up the pipeline without any additional costs. The decoder performs multi-head cross-attention via separable dynamic convolution for better efficiency and accuracy. To the best of our knowledge, YOSO is the first real-time panoptic segmentation framework that delivers competitive performance compared to state-of-the-art models. Specifically, YOSO achieves 46.4 PQ, 45.6 FPS on COCO; 52.5 PQ, 22.6 FPS on Cityscapes; 38.0 PQ, 35.4 FPS on ADE20K; and 34.1 PQ, 7.1 FPS on Mapillary Vistas. Code is available at <https://github.com/hujiectpp/YOSO>.*

## 1. Introduction

Panoptic segmentation is a task that involves assigning a semantic label and an instance identity to each pixel of an input image. The semantic labels are typically classified into two types, *i.e.*, *stuff* including amorphous and uncountable concepts (such as sky and road), and *things* consisting of countable categories (such as persons and cars). This division of label types naturally separates panoptic segmentation into two sub-tasks: semantic segmentation for *stuff* and instance segmentation for *things*. Thus, one of the major challenges for achieving real-time panoptic segmentation is the requirement for separate and computationally intensive branches to perform semantic and instance segmentation respectively. Typically, instance segmentation employs boxes or points to distinguish between different *things*, while se-

manic segmentation predicts distribution maps over semantic categories for *stuff*. As shown in Fig. 1, numerous efforts [9, 16, 23, 24, 31, 45] have been made to unify panoptic segmentation pipelines for improved speed and accuracy. However, achieving real-time panoptic segmentation still remains an open problem. On the one hand, heavy necks, *e.g.*, the multi-scale feature pyramid network (FPN) used in [27, 54], and heads, *e.g.*, the Transformer decoder used in [10, 58], are required to ensure accuracy, making real-time processing unfeasible. On the other hand, reducing the model size [16, 23, 24] leads to a decrease in model generalization. Therefore, developing a real-time panoptic segmentation framework that delivers competitive accuracy is challenging yet highly desirable.

In this paper, we present YOSO, a real-time panoptic segmentation framework. YOSO predicts panoptic kernels to convolute image feature maps, with which you only need to segment once for the masks of background *stuff* and foreground *things*. To make the process lightweight, we design a feature pyramid aggregator for extracting image feature maps, and a separable dynamic decoder for generating panoptic kernels. In the aggregator, we propose convolution-first aggregation (CFA) to re-parameterize the interpolation-first aggregation (IFA), resulting in an approximately 2.6 $\times$  speedup in GPU latency without compromising performance. Specifically, we demonstrate that the order, *i.e.*, interpolation-first or convolution-first, of applying bilinear interpolation and  $1\times 1$  convolution (w/o bias) does not affect results, but the convolution-first way provides a considerable speedup to the pipeline. In the decoder, we propose separable dynamic convolution attention (SDCA) to perform multi-head cross-attention in a weight-sharing way. SDCA achieves better accuracy (+1.0 PQ) and higher efficiency (approximately 1.2 $\times$  faster GPU latency) than traditional multi-head cross-attention.

In general, YOSO has three notable advantages. First, CFA reduces computational burden without re-training the model or compromising performance. CFA can be adapted to any task that uses the combination of bilinear interpolation and  $1\times 1$  convolution operations. Second, SDCA performs multi-head cross-attention with better accuracy and

\*Corresponding author

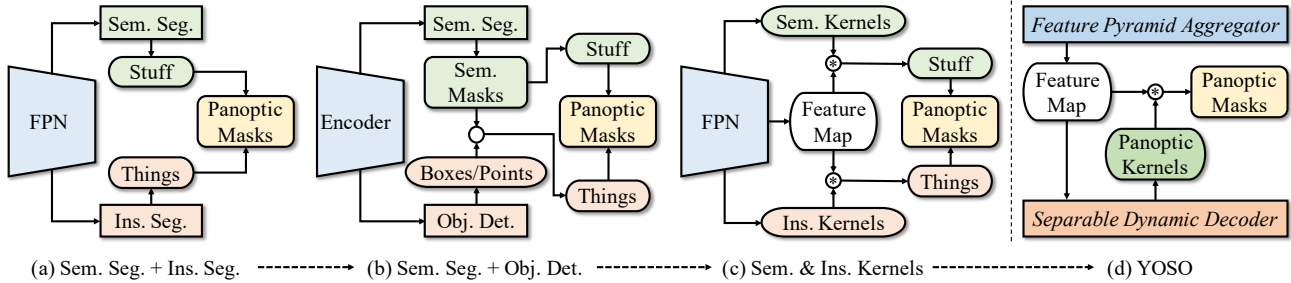


Figure 1. **Towards real-time panoptic segmentation.** (a) Semantic and instance segmentation are performed using shared FPN but separated task branches (e.g., in PanopticFPN [27] and UPSNet [54]). (b) Semantic segmentation generates masks for all categories, and instance recognition is achieved by object detection using boxes or points (e.g., in RealTimePan [24] and PanopticDeepLab [9]). (c) Kernels for *stuff* and *things* are generated to convolute image feature maps via heavy modules (e.g., in PanopticFCN [31], K-Net [58], and MaskFormer [10, 11]). (d) YOSO employs an efficient feature pyramid aggregator and a lightweight separable dynamic decoder to produce image feature maps and panoptic kernels. The figures do not include input images and backbone for concision.

efficiency. Third, YOSO runs faster and has competitive accuracy compared to state-of-the-art panoptic segmentation models, and its generalization is validated on four popular datasets: COCO (46.4 PQ, 45.6 FPS), Cityscapes (52.5 PQ, 22.6 FPS), ADE20K (38.0 PQ, 35.4 FPS), and Mapillary Vistas (34.1 PQ, 7.1 FPS).

## 2. Related Work

**Real-Time Panoptic Segmentation.** Panoptic segmentation aims to jointly perform semantic and instance segmentation, where each pixel in an input image is assigned both a semantic label and a unique instance identity. Many studies have been conducted for fast panoptic segmentation [8, 9, 16, 18, 23, 24, 31, 37, 41, 42, 45, 54, 55]. For instance, UPSNet [54] utilizes a deformable convolution based semantic segmentation head and a Mask R-CNN [21] style instance segmentation head. FPSNet [16] proposes a fast architecture for panoptic segmentation, avoiding instance mask prediction and merging outputs via soft attention masks. Recently, PanopticDeepLab [9], LPSNet [23], and RealTimePan [24] generate semantic masks for all categories first, then locate masks of instances via boxes or points, enabling efficient object segmentation. Meanwhile, PanopticFCN [31], K-Net [58], and MaskFormer [10, 11] attempt to predict masks for both *things* and *stuff* simultaneously through dynamic convolutions. Despite significant advances in this field, achieving real-time panoptic segmentation remains an open problem. In this paper, YOSO enables real-time panoptic segmentation with competitive accuracy by utilizing the proposed feature pyramid aggregator and separable dynamic decoder.

**Real-Time Instance Segmentation.** Instance segmentation aims to predict masks and classes for each instance in an image. To achieve real-time instance segmentation, various approaches have been proposed in recent literature. YOLACT [2, 3] proposes to multiply the predicted mask

coefficients with prototype masks, and SipMask [5] utilizes spatial mask coefficients for more accurate segmentation. CenterMask [30] employs an efficient anchor-free framework, and DeepSnake [40] explores the use of object contours for fast segmentation of instances. OrientMask [19] designs discriminative orientation maps that recover the masks without additional foreground segmentation, and SOLO [49, 50] segments objects by locations, with a decoupled branch to speed up the framework. Recently, SparseInst [12] introduced a sparse set of instance activation maps that highlight informative regions for each object in an image, constructing a real-time instance segmentation framework. As complex operations are required to distinguish different *things*, solving instance segmentation efficiently is the key to real-time panoptic segmentation. In this paper, YOSO predicts unified panoptic kernels for *stuff* and *things*. Implementing bipartite matching loss [6] for fast discrimination of different *things*, YOSO avoids time-consuming object localization operations such as RoIAlign [21] and post-processes such as non-maximum suppression. The output masks naturally represent independent instances for the categories of *things*. Moreover, the experimental results in our supplementary material show that YOSO can also achieve competitive performance on real-time instance segmentation.

**Real-Time Semantic Segmentation.** Semantic segmentation aims to predict pixel-wise categories for input images. In recent years, many approaches have been developed to enable real-time semantic segmentation. For example, E-Net [39] proposes a lightweight architecture for high-speed segmentation, and SegNet [1] combines a small network architecture with skip connections to achieve fast segmentation. ICNet [59] uses an image cascade algorithm to speed up the pipeline, and ESPNet [35, 36] introduces an efficient spatial pyramid dilated convolution. Additionally, BiSeNet [56, 56] separates spatial details and categorical semantics to enable both high accuracy and high efficiency

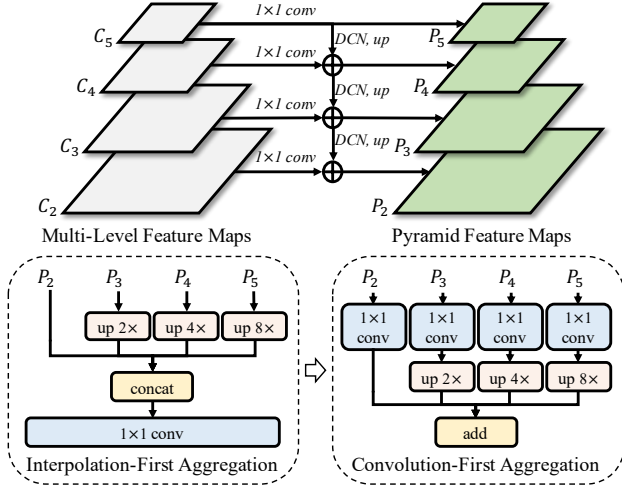


Figure 2. **Feature pyramid aggregator.** Multi-level feature maps are fused to generate pyramid features via  $1 \times 1$  convolutional layers and DCNs. Then, the pyramid feature maps are aggregated via CFA that can re-parameterize IFA for efficiency.

in semantic segmentation. More recently, SegFormer [53] employs Transformers with a lightweight multi-layer perceptron decoder for fast semantic segmentation. In contrast to traditional methods that predict distribution maps over classes for semantic masks, YOSO predicts kernels with their corresponding categories for segmentation. This enables an efficient way to jointly solve semantic and instance segmentation for panoptic segmentation.

### 3. Method

#### 3.1. Task Formulation

**Unified Panoptic Segmentation.** Panoptic segmentation maps each pixel of an image to a semantic class and an instance identity. We propose a unified approach that considers the set of foreground and background classes as a single entity. To achieve this, we aim to predict  $n$  binary masks  $\mathbf{M} \in \mathbb{B}^{n \times h \times w}$  with class probabilities  $\mathbf{L} \in \mathbb{R}^{n \times l}$  for an input image, where  $(h, w)$  denotes the mask resolution and  $l$  is the total number of categories. During training, the predictions are matched with the corresponding ground truth labels using the set prediction loss [6, 10]. At test time, the segmentation results are merged in terms of the foreground (*i.e.*, *things*) and background (*i.e.*, *stuff*) classes. Concretely, the masks that correspond to the same background class are merged via union operation; the masks with foreground classes are treated as independent instances; if a pixel belongs to multiple classes, the class with the highest probability is assigned to that pixel.

**YOSO Framework.** As shown in Fig.1, YOSO is a compact framework designed for real-time panoptic segmentation, which consists of a feature pyramid aggrega-

tor and a separable dynamic decoder. The backbone network, such as ResNet [22], extracts multi-level feature maps from input images. The feature pyramid aggregator compresses and aggregates the multi-level feature maps into single-level. The separable dynamic decoder then generates panoptic kernels with the single-level feature maps for both mask prediction and classification.

#### 3.2. Feature Pyramid Aggregator

**Deformable Feature Pyramid.** Given the multi-level feature maps  $C_2, C_3, C_4, C_5$  extracted by the backbone, we utilize the deformable feature pyramid network [15, 33] to enhance the feature maps from different scales. As depicted in Fig. 2, we first apply  $1 \times 1$  convolutional layers to compress the channels of the multi-level feature maps. Then, the feature maps from  $C_3$  to  $C_5$  are fed into deformable convolutional networks (DCNs) and upsampled level-by-level, yielding the pyramid feature maps  $P_2 \in \mathbb{R}^{c_2 \times h \times w}$ ,  $P_3 \in \mathbb{R}^{c_3 \times h/2 \times w/2}$ ,  $P_4 \in \mathbb{R}^{c_4 \times h/4 \times w/4}$ , and  $P_5 \in \mathbb{R}^{c_5 \times h/8 \times w/8}$ , respectively, where  $c_2, c_3, c_4, c_5$  denote the channel dimensions, and  $h, w$  represent the  $1/4$  scale of the input images. After obtaining the pyramid feature maps, we explore two aggregation methods, namely interpolation-first aggregation (IFA) and convolution-first aggregation (CFA), to merge the multi-level feature maps.

**Interpolation/Convolution-First Aggregation.** In IFA, the pyramid feature maps are first upsampled to the scale of  $h \times w$  via bilinear interpolation. Then, the feature maps are concatenated and fused using a  $1 \times 1$  convolutional layer. In CFA, the pyramid feature maps are first fed to different  $1 \times 1$  convolutional layers. Then, the feature maps are bilinearly interpolated to the scale of  $h \times w$  and summed.

**Observation I:** *The output of IFA is exactly equal to that of CFA when using  $1 \times 1$  convolution without bias.*

This observation is attributed to the homogeneity and additivity properties of the bilinear interpolation function  $f(\cdot)$ , where  $f(\sum_i w^i \mathbf{v}_{x,y}^i) = \sum_i w^i f(\mathbf{v}_{x,y}^i)$  for the constant  $w^i$  and the value vector  $\mathbf{v}_{x,y}^i$  from the four positions  $(x_1, y_1), (x_1, y_2), (x_2, y_1), (x_2, y_2)$  of the  $i$ -th feature map. Specifically, the bilinear interpolation estimates the value at  $(x_0, y_0)$  using the four positions by:

$$f(\sum_i w^i \mathbf{v}_{x,y}^i) = \frac{1}{(x_2 - x_1)(y_2 - y_1)} \begin{bmatrix} x_2 - x_0 \\ x_0 - x_1 \end{bmatrix}^\top \begin{bmatrix} \sum_i w^i \mathbf{v}_{x_1, y_1}^i & \sum_i w^i \mathbf{v}_{x_1, y_2}^i \\ \sum_i w^i \mathbf{v}_{x_2, y_1}^i & \sum_i w^i \mathbf{v}_{x_2, y_2}^i \end{bmatrix} \begin{bmatrix} y_2 - y_0 \\ y_0 - y_1 \end{bmatrix} \quad (1)$$

$$= \sum_i w^i f(\mathbf{v}_{x,y}^i),$$

where  $w^i$  can be interpreted as a  $1 \times 1$  kernel that convolutes the values of the original positions in the feature maps. Eq. 1 implies that applying  $1 \times 1$  convolution (without bias) before or after bilinear interpolation does not affect the final results. Hence, it can be inferred that when using  $1 \times 1$

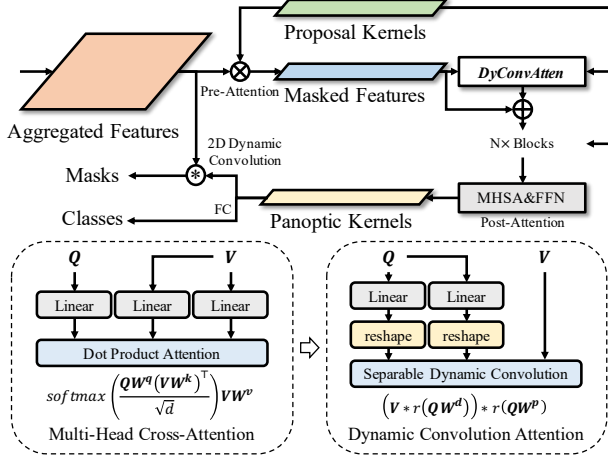


Figure 3. **Separable dynamic decoder.** Separable dynamic convolution learns to generate panoptic kernels for 2D dynamic convolutions, which performs multi-head cross-attention in a weight-sharing manner and achieves better accuracy.

convolution without bias in the aggregators, IFA and CFA produce identical outputs.

**Observation II:** CFA requires significantly fewer floating point operations (FLOPs) than IFA.

The reduction ratio of FLOPs between IFA and CFA is:

$$\frac{4(c_5 + c_4 + c_3)hw + (c_5 + c_4 + c_3 + c_2)dhw}{\left(\frac{c_5}{64} + \frac{c_4}{16} + \frac{c_3}{4} + c_2\right)dhw + 12dhw + 3dhw} > 1, \quad (2)$$

where  $d$  is the channel dimension of output feature maps. In the numerator of Eq. 2 (i.e., the FLOPs of IFA), the first term represents the number of FLOPs used in bilinear interpolation, and the second term represents the number of FLOPs used in  $1 \times 1$  convolution. In the denominator of Eq. 2 (i.e., the FLOPs of CFA), the terms represent the number of FLOPs for  $1 \times 1$  convolution, bilinear interpolation, and accumulation, respectively.

Given the above two observations, we adopt CFA in the proposed feature pyramid aggregator. It is noteworthy that the learned weights of the  $1 \times 1$  convolutional layer in IFA can be readily re-parameterized to CFA by dividing the weights into four  $1 \times 1$  convolutional layers. This can accelerate the pipeline without incurring any additional costs.

### 3.3. Separable Dynamic Decoder

In order to generate accurate kernels for segmentation, previous methods typically relied on dense predictors [31] or heavy Transformer decoders [58]. In contrast, we propose a lightweight kernel generator called the separable dynamic decoder, which speeds up kernel generation while maintaining high accuracy. The separable dynamic decoder, shown in Fig. 3, consists of three modules: a pre-attention module, a separable dynamic convolution module, and a post-attention module. Specifically, the separable dynamic

convolution efficiently performs multi-head cross-attention and achieves better accuracy. We describe each module in detail below.

**Pre-Attention.** Inspired by [58], the pre-attention module selectively extracts key information from the aggregated feature maps to diversify the proposal kernels. Concretely, the aggregated feature maps  $S \in \mathbb{R}^{d \times h \times w}$  are convolved by the learnable proposal kernels  $Q \in \mathbb{R}^{n \times d}$  to produce attention maps  $A \in \mathbb{R}^{n \times h \times w}$ . A hard sigmoid function  $\sigma(\cdot)$  is applied to the attention maps to activate the input values and discretize the values to either 0 or 1 using a threshold of 0.5. Then, with the attention maps, the masked features  $V \in \mathbb{R}^{n \times d}$  can be obtained by:

$$V = r(\sigma(S * Q))r(S)^T = r(\sigma(A))r(S)^T, \quad (3)$$

where  $r(\cdot)$  reshapes  $A$  and  $S$  to the sizes of  $(n, hw)$  and  $(d, hw)$ , respectively,  $*$  denotes 2D convolution operation.

**Vanilla Dynamic Convolution.** To increase the model capacity, multi-head cross-attention [46] has been widely used, which projects input features into multiple spaces to model their relationships. For instance, given  $Q, V \in \mathbb{R}^{n \times d}$  which contains  $n$  tokens with the hidden dimension of  $d$ , a  $t$ -head cross-attention can be defined as follows:

$$\text{MultiHeadCrossAtten}(Q, V) = \text{Concat}(H_1, H_2, \dots, H_t)W^o, \quad (4)$$

where  $W^o \in \mathbb{R}^{d \times d}$  is a linear transformation matrix. The attention operation inside each head is defined as:

$$\begin{aligned} H_i &= \text{softmax}\left(\frac{QW_i^q(VW_i^k)^T}{\sqrt{d}}\right)(VW_i^v) \\ &= K_i(VW_i^v) = K_iV_i, \end{aligned} \quad (5)$$

where  $W_i^q, W_i^k, W_i^v \in \mathbb{R}^{d \times d/t}$  are the projection matrices,  $V_i \in \mathbb{R}^{n \times d/t}$  denotes the projected features, and  $K_i \in \mathbb{R}^{n \times n}$  denotes the correlation matrix.

Although increasing the model capacity enhances performance, it also results in high computational burden. Intuitively, multi-head cross-attention involves three fundamental operations: multi-head projection, cross-token interaction, and cross-dimension interaction. In Eq. 4, the cross-dimension interaction, represented by  $W^o$ , learns to re-weight the importance of every hidden dimension. In Eq. 5, the multi-head projection, represented by  $VW_i^v$ , maps the hidden dimensions  $d$  into  $t$  different spaces with the size of  $d/t$ , and the cross-token interaction, represented by  $K_iV_i$ , uses the correlation matrix to interact between tokens. These operations motivated us to perform multi-head cross-attention using 1D convolution to make the process lightweight, which is defined as:

$$\text{Conv}_{1d}(K, V) = V * K, \quad (6)$$

where  $\mathbf{K} \in \mathbb{R}^{n \times n \times t}$  denotes the weights of kernels to convolute  $\mathbf{V} \in \mathbb{R}^{n \times d}$ ,  $*$  denotes 1D convolution operation. The 1D convolution computes the element of the output  $\mathbf{O} \in \mathbb{R}^{n \times d}$  at the position  $i, j$  by:

$$O_{i,j} = \sum_{p=1}^n \sum_{q=1}^t K_{i,p,q} \cdot V_{p,j+q-1}. \quad (7)$$

Correspondingly, the basic operations of multi-head cross-attention are also performed in 1D convolution in a weight-sharing manner. For the multi-head projection, the sliding window in 1D convolution densely splits the hidden dimensions into  $d$  groups of size  $t$ , and the  $t$  successive hidden dimensions in each group are projected with shared kernels. For the cross-token interaction, the first accumulation term in Eq. 7 interacts  $n$  tokens by  $n$  different kernels. For the cross-dimension interaction, the second accumulation term in Eq. 7 locally incorporates the information from  $t$  successive hidden dimensions, instead of using all hidden dimensions globally in multi-head cross-attention. Furthermore, inspired by [13, 26, 29, 48, 51, 57], we employ a dynamic approach to generate the kernel  $\mathbf{K}$  conditioned on  $\mathbf{Q}$ , which introduces the cross-attention mechanism into 1D convolution and defines a dynamic convolution attention as follows:

$$\text{DyConvAtten}(\mathbf{Q}, \mathbf{V}) = \mathbf{V} * r(\mathbf{Q}\mathbf{W}), \quad (8)$$

where  $\mathbf{W} \in \mathbb{R}^{d \times nt}$  is the projection matrix for generating the dynamic kernels.

**Separable Dynamic Convolution.** The standard convolution can be further decomposed into a depthwise convolution and a pointwise convolution as in [25]. Following this approach, we propose a separable form for vanilla dynamic convolution as follows:

$$\text{SepDyConvAtten}(\mathbf{Q}, \mathbf{V}) = (\mathbf{V} * r(\mathbf{Q}\mathbf{W}^d)) * r(\mathbf{Q}\mathbf{W}^p), \quad (9)$$

where  $\mathbf{W}^d \in \mathbb{R}^{d \times t}$ ,  $\mathbf{W}^p \in \mathbb{R}^{d \times n}$  project  $\mathbf{Q}$  to generate the kernels for depthwise convolution and pointwise convolution. The kernels are then reshaped to  $(n, 1, t)$  and  $(n, n, 1)$ , respectively. In Eq. 9, the depthwise convolution models the cross-dimension interaction, while the pointwise convolution models the cross-token interaction. Compared to multi-head cross-attention, the reduction in FLOPs is:

$$\frac{4nd^2 + 2n^2d}{2ndt + 2n^2d} = \frac{2d + n}{k + n} > 1. \quad (10)$$

In the numerator of Eq. 10 (*i.e.*, the FLOPs of MHCA), the first term denotes the number of FLOPs used in multi-head projection, while the second term represents the number of FLOPs used in cross-attention. In the denominator

of Eq. 10 (*i.e.*, the FLOPs of SDCA), the terms correspond to the number of FLOPs for 1D convolution operation and linear projection, respectively.

**Post-Attention.** In the post-attention module, we use a multi-head self-attention layer and a feed-forward network to generate the panoptic kernels. Subsequently, we produce the masks by means of 2D convolution and predict the associated classes using additional feed-forward networks. Inspired by [4, 7], we utilize the panoptic kernels to update the proposal kernels iteratively for improved accuracy.

## 4. Experiments

### 4.1. Datasets and Evaluation Metrics

**Datasets.** We evaluated the effectiveness and efficiency of YOSO on four widely used panoptic segmentation datasets: the COCO dataset [34], the Cityscapes dataset [14], the ADE20K [60] dataset, and the Mapillary Vistas [38] dataset. The COCO dataset gathers images of complex everyday scenes with common objects, containing 80 *things* categories and 53 *stuff* categories in 118k images for training and 5k images for validation. The Cityscapes dataset contains images of urban street-view scenes, which has 8 *things* categories and 11 *stuff* categories in 2.9k images for training and 0.5k images for validation. The ADE20K dataset is annotated in an open-vocabulary setting with 50 *things* categories and 100 *stuff* categories, including 20k images for training and 2k images for validation. The Mapillary Vistas dataset is a large-scale urban street-view dataset with 37 *things* categories and 28 *stuff* categories in 18k and 2k images for training and validation.

**Evaluation Metrics.** The panoptic segmentation results are assessed using the panoptic quality (PQ) metric [28], which can be further decomposed to the segmentation quality (SQ) and the recognition quality (RQ). The evaluation outcomes for *stuff* and *things* are represented by the superscripts  $s$  and  $t$ , respectively. The frame rate, *i.e.*, frames per second (FPS), of the YOSO models is evaluated on a single V100 GPU for the main results and a single 3090 GPU for the ablation study.

### 4.2. Implementation Details

For the COCO dataset, we used a batch size of 16 and set the learning rate to 0.0001. The models were trained for 370k iterations with large-scale jitter augmentation [20]. For the Cityscapes and Mapillary datasets, we trained the models with a batch size of 16, the learning rate set to 0.0001, and the training schedule set to 180k iterations. For the ADE20K dataset, we set the batch size to 16, the learning rate to 0.0001, and the models were trained for 30k iterations. The ResNet50 [22] (denoted as R50) pre-trained on the ImageNet [17] dataset is employed as our backbone for the four datasets. We used ResNet50 [22] (denoted as

Method	Backbone	Scale	PQ	PQ <sup>t</sup>	PQ <sup>s</sup>	FPS <sup>↑</sup>	GPU
BGRNet [52]	R50-FPN	800,1333	43.2	49.8	33.4	-	-
K-Net [58]	R50-FPN	800,1333	47.1	51.7	40.3	-	-
PanSegFormer [32]	R50	800,1333	49.6	54.4	42.4	-	-
Max-DeepLab [47]	Max-S	800,1333	48.4	53.0	41.5	7.6	V100
Mask2Former [10]	R50	800,1333	51.9	57.7	43.0	8.6	V100
UPSNNet [54]	R50-FPN	800,1333	42.5	48.5	33.4	9.1	V100
PanopticFCN [31]	R50-FPN	800,1333	44.3	50.0	35.6	9.2	V100
LPSNet [23]	R50-FPN	800,1333	39.1	43.9	30.1	9.3	V100
RealTimePan [24]	R50-FPN	800,1333	37.1	41.0	30.7	15.9	V100
PanopticFPN [27]	R50-FPN	800,1333	41.5	48.3	31.2	17.5	V100
MaskFormer [11]	R50	800,1333	46.5	51.0	39.8	17.6	V100
PanopticDeepLab [9]	R50	641,641	35.1	-	-	20.0	V100
<b>YOSO, ours</b>	R50	800,1333	48.4	53.5	40.8	23.6	V100
<b>YOSO, ours</b>	R50	512,800	46.4	50.7	40.0	45.6	V100

Table 1. Panoptic segmentation on the **COCO** validation set.

Method	Backbone	Scale	PQ	PQ <sup>t</sup>	PQ <sup>s</sup>	FPS <sup>↑</sup>	GPU
PanopticFPN [27]	R50-FPN	1024,2048	57.7	51.6	62.2	-	-
Seamless [43]	R50	1024, 2048	59.8	54.6	63.6	-	-
PanopticFCN [31]	R50-FPN	1024,2048	61.4	54.8	66.6	-	-
Mask2Former [10]	R50	1024,2048	62.1	54.9	67.3	4.1	V100
UPSNNet [54]	R50-FPN	1024,2048	59.3	54.6	62.7	7.5	V100
LPSNet [23]	R50-FPN	1024,2048	59.7	54.0	63.9	7.7	V100
PanopticDeepLab [9]	R50-FPN	1024,2048	59.7	-	-	8.5	V100
FPSNet [16]	R50-FPN	1024,2048	55.1	-	-	8.8	Titan
RealTimePan [24]	R50-FPN	1024,2048	58.8	52.1	63.7	10.1	V100
<b>YOSO, ours</b>	R50	1024,2048	59.7	51.0	66.1	11.1	V100
<b>YOSO, ours</b>	R50	512,1024	52.5	43.5	59.1	22.6	V100

Table 2. Panoptic segmentation on the **Cityscapes** validation set.

Method	Backbone	Scale	PQ	PQ <sup>t</sup>	PQ <sup>s</sup>	FPS <sup>↑</sup>	GPU
BGRNet [52]	R50-FPN	-	31.8	34.1	27.3	-	-
PanSegFormer [32]	R50	-	36.4	35.3	38.6	-	-
MaskFormer [11]	R50	640,2560	34.7	32.2	39.7	-	-
Mask2Former [10]	R50	640,2560	39.7	39.0	40.9	11.1	V100
<b>YOSO, ours</b>	R50	640,2560	38.0	37.3	39.4	35.4	V100

Table 3. Panoptic segmentation on the **ADE20K** validation set.

Method	Backbone	Scale	PQ	PQ <sup>t</sup>	PQ <sup>s</sup>	FPS <sup>↑</sup>	GPU
AdaptIS [44]	R50	-	32.0	39.1	26.6	-	-
Seamless [43]	R50	-	36.2	33.6	40.0	-	-
LPSNet [23]	R50-FPN	-	36.5	33.2	41.0	-	-
PanopticFCN [31]	R50-FPN	-	36.9	32.9	42.3	-	-
PanopticDeepLab [9]	R50	2176,2176	33.3	-	-	3.5	V100
Mask2Former [10]	R50	2048,2048	36.3	-	-	3.2	A100
<b>YOSO, ours</b>	R50	2048,2048	34.1	24.3	47.2	7.1	A100

Table 4. Panoptic segmentation on the **Mapillary** validation set.

R50) pre-trained on ImageNet [17] as the backbone for all four datasets, and set the hidden dimension  $d$  to 256 through experimentation. Since the COCO dataset contains images from various scenes, ranging from indoor to outdoor, we

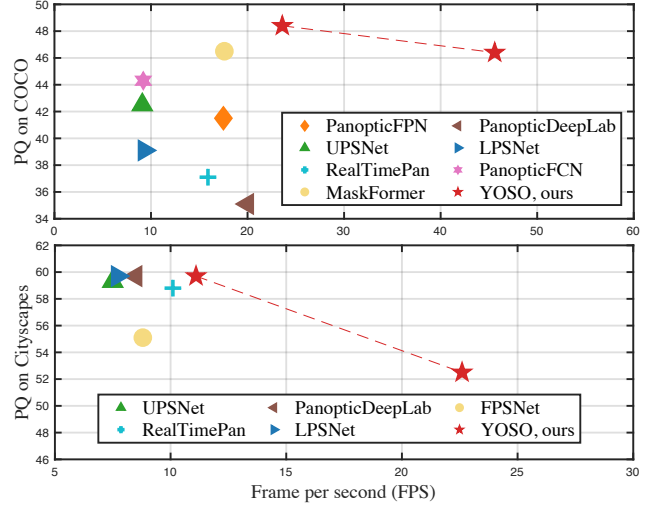


Figure 4. **FPS w.r.t. PQ** on the COCO and Cityscapes datasets.

conducted ablation studies on this dataset. In the ablation studies, the models were trained with 270k iterations.

### 4.3. Main Results

The results of panoptic segmentation on the COCO dataset are presented in Tab. 1. We have the following findings. First, YOSO is significantly faster than predominant efficient panoptic segmentation models such as PanopticFPN [27] and RealTimePan [24]. Specifically, YOSO achieves a PQ of 48.4 and an FPS of 23.6 with an input image scale of (800, 1333). This PQ is 11.3 higher than that of RealTimePan, while the speed is approximately  $1.5\times$  faster. Furthermore, when scaling the input image to (512, 800), YOSO is around  $2.3\times$  faster than the previous fastest model, PanopticDeepLab, while achieving an approximately 11.0-point higher PQ. Second, YOSO achieves comparable accuracy with state-of-the-art models such as MaskFormer [11], Mask2Former [10], Max-DeepLab [47], and K-Net [58]. For example, YOSO outperforms MaskFormer and K-Net by 1.9 and 1.3 PQ, respectively, and achieves the same PQ performance as Max-DeepLab. Although the PQ of YOSO is 3.5 points lower than that of Mask2Former, YOSO is  $2.7\times$  faster than Mask2Former.

The results of panoptic segmentation on the Cityscapes dataset are presented in Tab. 2. YOSO is the fastest model with competitive accuracy among state-of-the-art approaches. For example, YOSO achieves 59.7 PQ and 11.1 FPS with the input image size of (1024, 2048), which is 4.7 points higher than FPSNet. Moreover, when reducing the input image scale to (512, 1024), YOSO achieves an accuracy of 52.5 PQ with 22.6 FPS.

In Tab. 3 and Tab. 4, we show the panoptic segmentation results on the ADE20K and the Mapillary Vistas datasets, respectively, to evaluate the model generalization of YOSO.

Aggregator	PQ	SQ	RQ	PQ <sup>t</sup>	PQ <sup>s</sup>	FLOPs	Latency ( $\mu$ s) $\downarrow$	FPS $\uparrow$
IFA	47.5	82.2	56.9	52.7	39.7	16.6G	4871 $\pm$ 11	23.3
CFA	47.0	81.4	56.2	52.3	39.0	2.1G	1877 $\pm$ 52	29.2

Table 5. **Comparison of different aggregators.** The FLOPs and GPU latency were obtained from single modules with the setting of  $d=256$ ,  $c_2=128$ ,  $c_3=256$ ,  $c_4=512$ ,  $c_5=1024$ , and  $h=w=256$ .

Attention	PQ	SQ	RQ	PQ <sup>t</sup>	PQ <sup>s</sup>	FLOPs	Latency ( $\mu$ s) $\downarrow$	FPS $\uparrow$
MHCA	46.0	81.9	55.1	51.5	37.7	31.5M	2608 $\pm$ 210	27.3
SDCA	47.0	81.4	56.2	52.3	39.0	5.4M	2183 $\pm$ 279	29.2
DCA	46.9	82.0	55.8	51.9	38.7	15.5M	1701 $\pm$ 186	30.0
PDCA	43.7	81.3	52.3	49.3	35.3	5.2M	1450 $\pm$ 183	30.2
DDCA	46.6	82.3	55.9	52.3	38.4	0.3M	1242 $\pm$ 101	30.3

Table 6. **Comparison of different attention modules.** The FLOPs and GPU latency were obtained from single modules with the setting of  $n=100$ ,  $d=256$ , and  $t=3$ . The modules tested included MHCA (Multi-Head Cross-Attention), DCA (Dynamic Convolution Attention), SDCA (Separable Dynamic Convolution Attention), PDCA (Pointwise Dynamic Convolution Attention), and DDCA (Depthwise Dynamic Convolution Attention).

On the ADE20K dataset, YOSO outperforms most previous methods such as PanSegFormer [32] and MaskFormer [11] in terms of both speed and accuracy. On the Mapillary Vistas dataset, although YOSO has a good PQ<sup>s</sup>, the performance of PQ<sup>t</sup> lags behind that of state-of-the-art models. This suggests that YOSO still has the potential to be improved on the Mapillary Vistas dataset.

Additionally, we plot the PQ *w.r.t.* FPS results on the COCO and the Cityscapes datasets in Fig. 4, which shows that YOSO runs faster and achieves competitive accuracy among state-of-the-art models. In summary, the main results on the four datasets validate the good generalization and well-balanced speed-accuracy of YOSO.

#### 4.4. Ablation Study

In order to examine the impact of different components on the speed and accuracy of YOSO, we conducted several ablation studies focusing on the feature pyramid aggregator and the separable dynamic decoder. Specifically, we evaluated the effectiveness of the aggregation modules and the attention modules, which resulted in several interesting findings. Moreover, we analyzed how variations in the number of attention blocks, kernel size, iteration stages, and proposal kernels affected the performance. The details of our investigations are discussed below.

**Comparison of different aggregators.** The results of using different aggregators are presented in Tab. 5. Specifically, we trained YOSO with IFA and CFA, respectively. The PQ results indicate that IFA achieves higher accuracy than CFA, with PQ values of 47.5 and 47.0, respectively. However, IFA has much larger FLOPs than CFA, with 16.6G compared to 2.1G, and a slower GPU latency, with 4871 $\mu$ s compared to 1877 $\mu$ s.

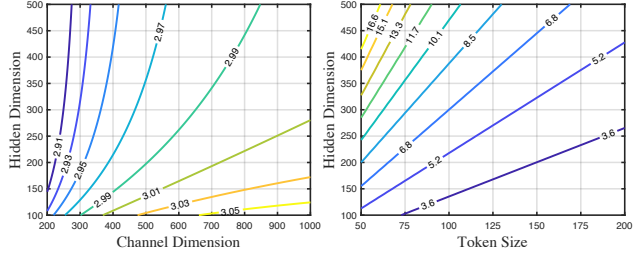


Figure 5. **Contours of FLOPs reduction ratio.** Left: CFA. Right: SDCA. The channel dimension of input feature maps for CFA and the hidden dimension of input tokens for SDCA are the key factors to reducing the number of FLOPs.

with 4871 $\mu$ s compared to 1877 $\mu$ s. Considering that the learned parameters of IFA can be directly re-parameterized to CFA, we can train YOSO with IFA and infer with the re-parameterized CFA for better speed and accuracy.

In Fig. 5 (left), we investigate the reduction ratio of FLOPs between CFA and IFA. Specifically, we set the channel dimensions of the input feature maps  $c_5$ ,  $c_4$ ,  $c_3$ , and  $c_2$  to be equal to  $c$ , and analyze how the input channel dimension  $c$  and the output hidden dimension  $d$  affect the FLOPs reduction ratio in Eq. 2. Our results indicate that increasing the input channel dimension and decreasing the output hidden dimension lead to an increase in the reduction ratio. This implies that CFA will be more efficient when the dimension of the input channel is large.

**Comparison of different attention modules.** We present a comparison of the effectiveness of different attention modules in Tab. 6. In terms of PQ performance, we have made two interesting findings. First, we were surprised to find that the MHCA module did not perform better than DCA and SDCA. The reason for this may be the difference between the basic operations in these two types of attention modules: DCA densely splits the hidden dimension into  $d$  groups, while MHCA sparsely splits it into  $t$  groups. Second, PDCA showed inferior performance compared to the other modules, with a PQ of 43.7. This is likely due to the fact that PDCA is the only module that does not apply cross-dimension interaction, which implies that interactions between hidden dimensions are significant for the attention module. This observation is supported by the performance of DDCA, which only performs cross-token interaction and achieved a PQ of 46.7. These two results suggest that the cross-dimension interaction may be more significant than the cross-token interaction in panoptic kernel generation, as the panoptic kernels are expected to be independent to represent different *stuff* or *things*.

In terms of FLOPs, we observed that the attention modules based on dynamic convolution require fewer FLOPs. Specifically, DDCA exhibits the lowest computational cost, requiring only 0.3M FLOPs. In terms of GPU latency, we made two interesting observations. First, although the

#Blocks	PQ	SQ	RQ	PQ <sup>t</sup>	SQ <sup>t</sup>	RQ <sup>t</sup>	PQ <sup>s</sup>	SQ <sup>s</sup>	RQ <sup>s</sup>	FPS↑
$N=1$	46.4	81.7	55.7	52.1	82.3	62.1	37.9	79.9	46.0	30.1
$N=2$	47.0	81.4	56.2	52.3	83.1	62.2	39.0	78.8	47.1	29.2
$N=3$	47.5	82.1	56.8	52.6	83.2	62.6	39.7	80.4	48.3	28.9
$N=4$	47.2	81.8	56.7	52.5	82.9	62.6	39.3	80.1	47.6	28.7

Table 7. **Influence of attention block number.** The PQ accuracy increases with more blocks, while the FPS decreases.

Kernel Size	PQ	SQ	RQ	PQ <sup>t</sup>	SQ <sup>t</sup>	RQ <sup>t</sup>	PQ <sup>s</sup>	SQ <sup>s</sup>	RQ <sup>s</sup>	FPS↑
$t=1$	46.8	82.3	55.9	52.4	83.4	62.3	38.5	80.6	46.4	30.6
$t=3$	47.0	81.4	56.2	52.3	83.1	62.2	39.0	78.8	47.1	29.2
$t=5$	47.3	81.3	56.8	52.2	82.0	62.2	40.0	80.3	48.6	28.3
$t=7$	47.1	81.2	56.3	52.4	83.2	62.3	39.8	80.2	48.4	27.6

Table 8. **Results of using different convolution kernel sizes.** The PQ performance saturates when  $t=5$ .

time complexity of MHCA, DCA, SDCA, and PDCA are all  $O(n^2d)$ , the speedup on the GPU latency is remarkable. Second, we found that SDCA runs slower than DCA, contrary to what the FLOPs analysis suggests. We speculate that the additional convolutional and fully connected layers in SDCA are executed serially rather than in parallel, which may result in longer execution times in practice.

Fig. 5 (right) presents the analysis of the FLOPs reduction ratio, as given by Eq. 10. We fix the kernel size  $t$  to 3, and investigate the effect of the token size  $n$  and the hidden dimension  $d$  on the reduction ratio. The results indicate that the reduction ratio enlarges with the token size decreases and the hidden dimension increases. Specifically, DCA performs better when the token size is smaller than the hidden dimension, indicating its suitability for vision tasks where the token size is smaller than the hidden dimension.

**Number of attention blocks.** To assess the effectiveness of the separable dynamic convolution attention in the separable dynamic decoder, we varied the number of blocks and analyzed the speed-accuracy trade-off, as presented in Tab. 7. The results show that the PQ accuracy improves with additional blocks, but at the expense of lower FPS. Consequently, we selected  $N = 2$  as a compromise between speed and accuracy for YOSO.

**Kernel size of dynamic convolution.** Tab. 8 shows the effectiveness of different kernel sizes for the separable dynamic convolution attention module. The PQ performance confirms our observation in Tab. 6, suggesting that the cross-dimension interaction plays a significant role in the attention modules. Specifically, reducing the kernel size from 5 to 1 leads to a degradation in the PQ performance from 47.3 to 46.8. Furthermore, the performance reaches saturation when the kernel size is increased from 5 to 7.

**Iteration stages.** To assess the impact of the number of stages, we conducted experiments with different numbers of stages and report the results in Tab. 9. Our results show

#Stages	PQ	SQ	RQ	PQ <sup>t</sup>	SQ <sup>t</sup>	RQ <sup>t</sup>	PQ <sup>s</sup>	SQ <sup>s</sup>	RQ <sup>s</sup>	FPS↑
$T=1$	45.7	80.4	53.6	50.3	82.1	59.6	36.5	77.3	45.2	30.1
$T=2$	47.0	81.4	56.2	52.3	83.1	62.2	39.0	78.8	47.1	29.2
$T=3$	47.5	81.9	56.9	52.8	82.9	63.0	39.4	80.4	47.7	28.9

Table 9. **Impact of iteration stages.** The PQ accuracy increases while the FPS decreases with more iteration stages.

#Proposals	PQ	SQ	RQ	PQ <sup>t</sup>	SQ <sup>t</sup>	RQ <sup>t</sup>	PQ <sup>s</sup>	SQ <sup>s</sup>	RQ <sup>s</sup>	FPS↑
$n=50$	44.2	80.1	52.6	49.8	81.2	59.2	36.1	78.4	44.1	33.7
$n=100$	47.0	81.4	56.2	52.3	83.1	62.2	39.0	78.8	47.1	29.2
$n=150$	46.9	81.0	55.7	52.0	82.7	61.5	38.9	78.6	46.8	26.3
$n=200$	46.7	81.1	55.2	51.5	82.5	61.2	38.7	78.5	46.2	24.3

Table 10. **Different number of proposals.** The PQ accuracy increases with more proposals and saturates when  $n=150$ .

that increasing the number of stages improves the PQ performance, but it comes at the cost of decreasing the FPS performance. We observed that the trade-off between speed and accuracy is best achieved when using  $T = 2$  stages. Hence, we selected this configuration for YOSO.

**Number of proposal kernels.** We investigate the impact of the number of proposal kernels in Tab. 10. The results indicate that the PQ performance improves when increasing the number of proposal kernels from 50 to 100, and saturates at 150. Meanwhile, the speed decreases when increasing the number of proposal kernels. The setting of  $n=100$  well balances the accuracy and speed for YOSO.

## 5. Conclusion

In this paper, we propose a real-time panoptic segmentation framework, termed YOSO. With YOSO, *you only need to segment once* for the masks of both foreground *things* and background *stuff*. YOSO includes a feature pyramid aggregator and a separable dynamic decoder to accelerate the pipeline. The CFA module in the feature pyramid aggregator re-parameters the IFA module, reducing FLOPs without extra costs. The SDC module in the separable dynamic decoder performs weight-sharing multi-head cross-attention, enhancing both speed and accuracy. Our extensive experiments demonstrate that YOSO is significantly faster than other prominent panoptic segmentation methods while maintaining competitive PQ performance. Given its effectiveness and simplicity, we hope YOSO can serve as a strong baseline and bring fresh insights for future research on real-time panoptic segmentation.

**Acknowledgements.** This work was supported by National Key R&D Program of China (2022ZD0118202), National Science Fund for Distinguished Young Scholars (No.62025603), National Natural Science Foundation of China (No.U21B2037, No.U22B2051, No.62176222, No.62176223, No.62176226, No.62072386, No.62072387, No.62072389, No.62002305 and No.62272401), and Natural Science Foundation of Fujian Province of China (No.2021J01002, No.2022J06001).



## 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*, 2017. 2
- [2] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019. 2
- [3] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact++: Better real-time instance segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2
- [4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018. 5
- [5] Jiale Cao, Rao Muhammad Anwer, Hisham Cholakkal, Fahad Shahbaz Khan, Yanwei Pang, and Ling Shao. Sipmask: Spatial information preservation for fast image and video instance segmentation. In *European Conference on Computer Vision*, 2020. 2
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, 2020. 2, 3
- [7] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 5
- [8] Xia Chen, Jianren Wang, and Martial Hebert. Panonet: Real-time panoptic segmentation through position-sensitive feature embedding. *arXiv preprint arXiv:2008.00192*, 2020. 2
- [9] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 1, 2, 6
- [10] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 1, 2, 3, 6
- [11] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems*, 2021. 2, 6, 7
- [12] Tianheng Cheng, Xinggang Wang, Shaoyu Chen, Wenqiang Zhang, Qian Zhang, Chang Huang, Zhaoxiang Zhang, and Wenyu Liu. Sparse instance activation for real-time instance segmentation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 2
- [13] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. In *International Conference on Learning Representations*, 2021. 5
- [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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016. 5
- [15] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2017. 3
- [16] Daan de Geus, Panagiotis Meletis, and Gijs Dubbelman. Fast panoptic segmentation network. *IEEE Robotics and Automation Letters*, 2020. 1, 2, 6
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2009. 5, 6
- [18] Manuel Diaz-Zapata, Özgür Ercent, and Christian Laugier. Yolo-based panoptic segmentation network. In *IEEE Computers, Software, and Applications Conference*, 2021. 2
- [19] Wentao Du, Zhiyu Xiang, Shuya Chen, Chengyu Qiao, Yiman Chen, and Tingming Bai. Real-time instance segmentation with discriminative orientation maps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 2
- [20] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 5
- [21] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2017. 2
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016. 3, 5
- [23] Weixiang Hong, Qingpei Guo, Wei Zhang, Jingdong Chen, and Wei Chu. Lpsnet: A lightweight solution for fast panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 1, 2, 6
- [24] Rui Hou, Jie Li, Arjun Bhargava, Allan Raventos, Vitor Guizilini, Chao Fang, Jerome Lynch, and Adrien Gaidon. Real-time panoptic segmentation from dense detections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 1, 2, 6
- [25] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 5
- [26] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, 2020. 5

- [27] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 6
- [28] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 5
- [29] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020. 5
- [30] Youngwan Lee and Jongyoul Park. Centermask: Real-time anchor-free instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [31] Yanwei Li, Hengshuang Zhao, Xiaojuan Qi, Liwei Wang, Zeming Li, Jian Sun, and Jiaya Jia. Fully convolutional networks for panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 1, 2, 4, 6
- [32] Zhiqi Li, Wenhai Wang, Enze Xie, Zhiding Yu, Anima Anandkumar, Jose M. Alvarez, Tong Lu, and Ping Luo. Panoptic segformer: Delving deeper into panoptic segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 6, 7
- [33] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017. 3
- [34] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, 2014. 5
- [35] Sachin Mehta, Mohammad Rastegari, Anat Caspi, Linda Shapiro, and Hannaneh Hajishirzi. Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation. In *European Conference on Computer Vision*, 2018. 2
- [36] Sachin Mehta, Mohammad Rastegari, Linda Shapiro, and Hannaneh Hajishirzi. Espnetv2: A light-weight, power efficient, and general purpose convolutional neural network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [37] Rohit Mohan and Abhinav Valada. Efficientpts: Efficient panoptic segmentation. *International Journal of Computer Vision*, 2021. 2
- [38] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulò, and Peter Kotschieder. The mapillary vistas dataset for semantic understanding of street scenes. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2017. 5
- [39] Adam Paszke, Abhishek Chaurasia, Sangpil Kim, and Eugenio Culurciello. Enet: A deep neural network architecture for real-time semantic segmentation. *arXiv preprint arXiv:1606.02147*, 2016. 2
- [40] Sida Peng, Wen Jiang, Huaijin Pi, Xiuli Li, Hujun Bao, and Xiaowei Zhou. Deep snake for real-time instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [41] Andra Petrovai and Sergiu Nedevschi. Real-time panoptic segmentation with prototype masks for automated driving. In *IEEE Intelligent Vehicles Symposium*, 2020. 2
- [42] Andra Petrovai and Sergiu Nedevschi. Fast panoptic segmentation with soft attention embeddings. *Sensors*, 2022. 2
- [43] Lorenzo Porzi, Samuel Rota Bulò, Aleksander Colovic, and Peter Kotschieder. Seamless scene segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 6
- [44] Konstantin Sofiiuk, Olga Barinova, and Anton Konushin. Adaptis: Adaptive instance selection network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019. 6
- [45] Zhi Tian, Bowen Zhang, Hao Chen, and Chunhua Shen. Instance and panoptic segmentation using conditional convolutions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 1, 2
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 4
- [47] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Max-deeplab: End-to-end panoptic segmentation with mask transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 6
- [48] Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020. 5
- [49] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *European Conference on Computer Vision*, 2020. 2
- [50] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. *Advances in Neural Information Processing Systems*, 2020. 2
- [51] Felix Wu, Angela Fan, Alexei Baevski, Yann N. Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. In *International Conference on Learning Representations*, 2019. 5
- [52] Yangxin Wu, Gengwei Zhang, Yiming Gao, Xiajun Deng, Ke Gong, Xiaodan Liang, and Liang Lin. Bidirectional graph reasoning network for panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 6
- [53] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 2021. 3
- [54] Yuwen Xiong, Renjie Liao, Hengshuang Zhao, Rui Hu, Min Bai, Ersin Yumer, and Raquel Urtasun. Upsnet: A

- unified panoptic segmentation network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2, 6
- [55] Tien-Ju Yang, Maxwell D Collins, Yukun Zhu, Jyh-Jing Hwang, Ting Liu, Xiao Zhang, Vivienne Sze, George Papandreou, and Liang-Chieh Chen. Deeperlab: Single-shot image parser. *arXiv preprint arXiv:1902.05093*, 2019. 2
- [56] Changqian Yu, Changxin Gao, Jingbo Wang, Gang Yu, Chunhua Shen, and Nong Sang. Bisenet v2: Bilateral network with guided aggregation for real-time semantic segmentation. *International Journal of Computer Vision*, 2021. 2
- [57] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *Advances in Neural Information Processing Systems*, 2020. 5
- [58] Wenwei Zhang, Jiangmiao Pang, Kai Chen, and Chen Change Loy. K-net: Towards unified image segmentation. *Advances in Neural Information Processing Systems*, 2021. 1, 2, 4, 6
- [59] Hengshuang Zhao, Xiaojuan Qi, Xiaoyong Shen, Jianping Shi, and Jiaya Jia. Icnnet for real-time semantic segmentation on high-resolution images. In *European Conference on Computer Vision*, 2018. 2
- [60] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017. 5