

Feature Selective Transformer for Semantic Image Segmentation

Fangjian Lin^{1, 2*}, Tianyi Wu^{1, 2*}, Sitong Wu^{1, 2}, Shengwei Tian[†], Guodong Guo^{1, 2†}

¹Institute of Deep Learning, Baidu Research, Beijing, China

²National Engineering Laboratory for Deep Learning Technology and Application, Beijing, China

{linfangjian01, wusitong98}@gmail.com, {wutianyi01, guogudong01}@baidu.com

Abstract. Recently, it has attracted more and more attentions to fuse multi-scale features for semantic image segmentation. Various works were proposed to employ progressive local or global fusion, but the feature fusions are not rich enough for modeling multi-scale context features. In this work, we focus on fusing multi-scale features from Transformer-based backbones for semantic segmentation, and propose a Feature Selective Transformer (FeSeFormer), which aggregates features from all scales (or levels) for each query feature. Specifically, we first propose a Scale-level Feature Selection (SFS) module, which can choose an informative subset from the whole multi-scale feature set for each scale, where those features that are important for the current scale (or level) are selected and the redundant are discarded. Furthermore, we propose a Full-scale Feature Fusion (FFF) module, which can adaptively fuse features of all scales for queries. Based on the proposed SFS and FFF modules, we develop a Feature Selective Transformer (FeSeFormer), and evaluate our FeSeFormer on four challenging semantic segmentation benchmarks, including PASCAL Context, ADE20K, COCO-Stuff 10K, and Cityscapes, outperforming the state-of-the-art.

Keywords: Vision Transformer, Segmentation, Multi-scale Fusion

1 Introduction

Semantic image segmentation is an essential and challenging task with high potential values in a variety of applications, *e.g.*, human-computer interaction [15], augmented reality [1], and driverless technology [13]. The task aims to classify each pixel into a semantic category. Since Long *et al.* [33] proposed fully convolutional networks (FCN), it has attracted more and more attentions to model multi-scale context features for semantic segmentation.

Currently, many works [46,27,29,22,26,50,28,51,42,54,48,40] have explored how to fuse multi-scale features and demonstrated that modeling multi-scale context is very beneficial for semantic segmentation. The first family of works

* Equal contributions. † Corresponding author.

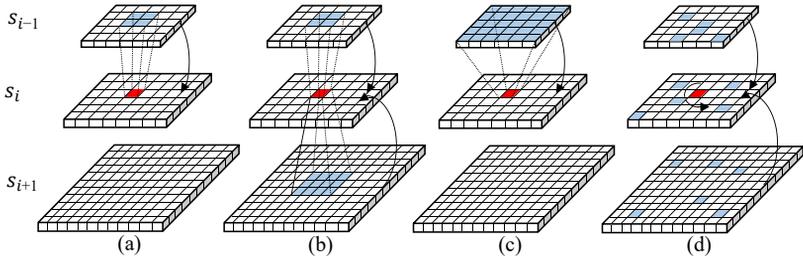


Fig. 1: Comparison with the existing feature fusion methods. A query feature (or position) is represented by the red dot. (a) Progressive local fusion, it fuses features of the local region (grey color) from its adjacent scale. (b) Non-progressive local fusion, it can fuse features of the local region from all levels or scales. (c) Progressive global fusion, it fuses features from its adjacent scale. (d) Ours, it can fuse features from all scales for each pixel.

utilized progressive local fusion [27,29,22,51]. For multi-scale features from CNN-based or Transformer-based backbone networks, up-samplings or convolutions are employed to transform per-scale features for local fusing. This local fusion is only conducted on two adjacent scales, and each pixel only collects information from the local region of its adjacent scale. As shown in Figure 1(a), the query (red) pixel on the features s_i only fuses features within the local region (of features s_{i-1}) where the center is the same as the red pixel. FPN [29] is a foundational work, which employed a top-down pathway to up-sample semantically stronger features and enhanced them with features from lower-level semantics by progressive local fusion. Following FPN, SETR-MLA [51] and Semantic-FPN [22] utilized this mechanism to fuse multi-scale features that share the same resolution. Substantial progress has been made by utilizing non-progressive local fusion. GFFNet[25] proposed to fuse features of local regions from all levels using gates. As shown in Figure 1(b), the query (red) pixel on the features s_i can fuse features on the local region from all levels. More recently, global fusion [48,54] is proposed, inspired by Non-local Networks [39] and Transformer [38]. This kind of method can fuse all features from another scale or level. As shown in Figure 1(c), the query (red) pixel on the features s_i can fuse all features of features s_{i-1} . ANN [54] proposed Asymmetric Fusion Non-local Block to fuse features from two different scales, while FPT [48] proposed Grounding Transformer to fuse higher-level feature maps to the lower-level ones. These methods achieved excellent performance, but these methods fused features from preset subset for queries, which may result in fusing irrelevant features for queries.

Different from the methods mentioned above, we propose a Feature Selective Transformer (FeSeFormer) for semantic segmentation, which can dynamically select informative subset from the whole multi-scale feature set to adaptively fuse them for queries. As shown in Figure 1(d), our method endows each pixel with the ability to choose informative features from all scales and positions, and use them

to enhance itself. Specifically, we first propose a Scale-level Feature Selection (SFS) module, which can select an informative subset from the whole multi-scale feature set for each scale, where those features that are important for the query scale (or features) are selected and the redundant ones are discarded. Secondly, we propose a Full-scale Feature Fusion (FFF) module, which can adaptively fuse features of all scales for queries.

Based on the proposed Scale-level Feature Selection and Full-scale Feature Fusion modules, we develop a Feature Selective Transformer (FeSeFormer) for semantic segmentation and demonstrate the effectiveness of our approach by conducting extensive ablation studies. Furthermore, we evaluate our FeSeFormer on four challenging semantic segmentation benchmarks, including PASCAL Context [35], ADE20K [53], COCO-Stuff 10K [2], and Cityscapes [9], achieving 58.91%, 54.43%, 49.80%, and 84.46% mIoU, respectively, which outperform the SOTA methods. Our main contributions include:

- We propose a Scale-level Feature Selection module, which selects an informative subset from the whole multi-scale feature set and discards the redundant for each scale.
- We propose a Full-scale Feature Fusion module, which can adaptively fuse features from all scales for modeling multi-scale contextual features.
- Based on the Scale-level Feature Selection module and the Full-scale Feature Fusion module, we develop a semantic segmentation framework, Feature Selective Transformer (FeSeFormer), outperforming the state-of-the-art on four challenging benchmarks, including PASCAL Context [35], ADE20K [53], COCO-Stuff 10K [2], and Cityscapes [9].

2 Related Work

In this section, we briefly review related works, including multi-scale features fusion and Transformer in semantic segmentation.

Multi-scale Features Fusion. There are various works exploring how to fuse multi-scale features for semantic segmentation. Inspired by FPN [29] that employed a top-down pathway and lateral connections for progressively fusing multi-scale features for object detection, Semantic-FPN [22] and SETR-MLA [51] extended this architecture to fuse multi-scale features for semantic segmentation. Based on this top-down fusion, ZigZagNet [28] proposed top-down and bottom-up propagations to aggregate multi-scale features, while FTN [40] proposed Feature Pyramid Transformer for multi-scale feature fusion. Differently, PSPNet [50] and DeepLab series [26,4,5] fused multi-scale features via concatenation at the channel dimension. Different from these methods that fused features on the local region, ANN [54] proposed an Asymmetric Fusion Non-local Block for fusing all features at one scale for each feature (position) on another scale, while FPT [48] proposed Grounding Transformer to ground the “concept” of the higher-level features to every pixel on the lower-level ones. Different from these methods that fuse features from preset subset for queries, we explore how

to dynamically select informative subset from the whole multi-scale feature set and fuse them for each query feature.

Transformer. Since Alexey *et al.* [12] introduced Visual Transformer (ViT) for image classification, it has attracted more and more attentions to explore how to use Transformer for semantic segmentation. These methods focused on exploring the various usages of Transformer, including extracting features [51,36,43] from input image, learning class embedding [44,37], or learning mask embedding [8]. For example, SETR [51] treated semantic segmentation as a sequence-to-sequence prediction task and deployed a pure transformer (i.e., without convolution and resolution reduction) to encode an image as a sequence of patches for feature extraction. DPT [36] reassembled the bag-of-words representation provided by ViT into image-like features at various resolutions, and progressively combined them into final predictions. Differently, Trans2Seg [44] formulated semantic segmentation as a problem of dictionary look-up, and designed a set of learnable prototypes as the query of Transformer decoder, where each prototype learns the statistics of one category. SegFormer [43] used Transformer-based encoder to extract features and the lightweight MLP-decoder to predict pixel by pixel. Segmenter [37] employed a Mask Transformer to learn a set of class embedding, which was used to generate class masks. Recent MaskFormer [8] proposed a simple mask classification model to predict a set of binary masks, where a transformer decoder was used to learn mask embedding. Different from these works, we explore how to use Transformer to fuse multi-scale features.

Most related to our work is GFFNet [25] and FPT [48]. GFFNet employed a gate mechanism to fuse features from multiple scales, and one feature (or pixel) fused the same position of all scales. However, our FeSeFormer can fuse features from all scales and positions simultaneously. Besides, our method is also different from FPT. The FPT proposed Grounding Transformer and Rendering Transformer to fuse the higher-level and lower-level features in a bidirectional fashion and merge two-scale features at a time. Besides, FPT required employing an extra Self-Transformer to fuse features within the same level or scale. However, our FeSeFormer simultaneously fuses features of all scales and positions for each query feature.

3 Method

We first describe our framework, Feature Selective Transformer (FeSeFormer). Then, we present the Scale-level Feature Selection (SFS) module, which is employed to select an informative subset from whole multi-scale feature set (from Transformer-based backbones) and discards redundant ones for each scale. Finally, we introduce the Full-scale Feature Fusion (FFF) module, which aims to adaptively fuse features from all scales and positions for modeling multi-scale context features.

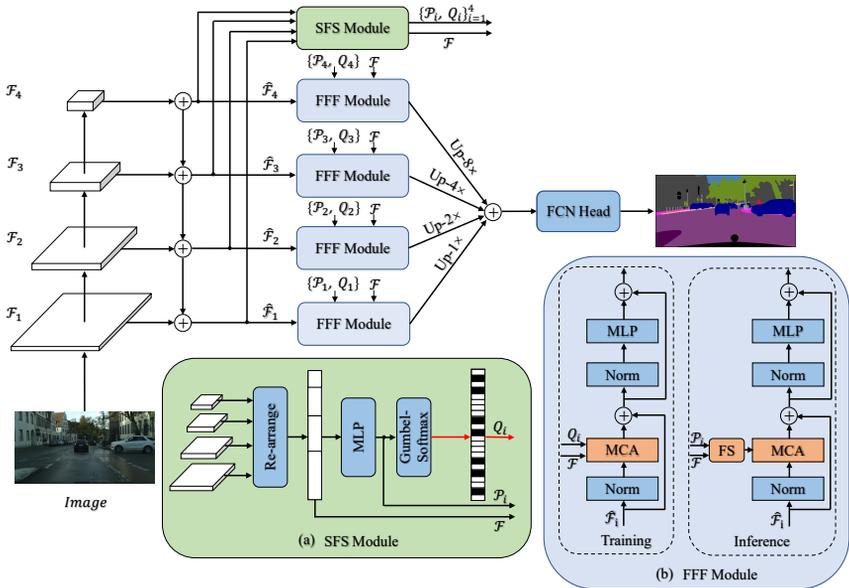


Fig. 2: Overall architecture of our Feature Selective Transformer (FeSeFormer). The red line does not need to be processed during inference. “MCA” and “FS” indicate multi-head cross attention and feature selection, respectively.

3.1 Framework

The overall framework of our FeSeFormer is shown in Figure 2, which consists of a Scale-level Feature Selection (SFS) module and a Full-scale Feature Fusion (FFF) module. Given an input image $\mathcal{I} \in \mathbb{R}^{3 \times H \times W}$, where H, W denotes the height and width, respectively. We first use a Transformer-based backbone (such as Swin Transformer [32]) to get multi-scale feature set $\{\mathcal{F}_i \in \mathbb{R}^{2^{i-1}C \times \frac{H}{2^{i+1}} \times \frac{W}{2^{i+1}}}\}_{i=1}^4$, where i indicates the scale or stage index of the backbone, and C is the base channel number. In order to select an informative subset of multi-scale features for each scale, we first employ a top-down pathway to inject the strongest semantics into all scales, getting an enhanced multi-scale representation set $\{\hat{\mathcal{F}}_i\}_{i=1}^4$, where $\hat{\mathcal{F}}_i \in \mathbb{R}^{D \times \frac{H}{2^{i+1}} \times \frac{W}{2^{i+1}}}$, and D indicates the channel number. Then, we employ the Scale-level Feature Selection module to choose an informative subset from the whole multi-scale feature set, which can adaptively select important features for each scale while discarding redundant features. After collecting important feature set for each scale, we utilize the proposed Full-scale Feature Fusion module to model multi-scale context features. Finally, we simply up-sample multi-scale feature maps to the same resolution, followed by FCN Head [33] (consisting of one 3×3 Conv. and $1 \times$ Conv.) for getting segmentation results.

3.2 Scale-level Feature Selection module

Previous works [27,22,50,28,51,54,48,40] have shown that fusing multi-scale features from multiple scales are critical for improving semantic segmentation, since the objects in the scene often present a variety of scales. Multi-scale features from the backbone networks usually have different spatial resolutions. High-resolution features contain more spatial details than low-resolution ones, while the latter have stronger semantics than the former. Besides, small-scale objects have no precise locations in the finer level (or higher-resolution), since they have been down-sampled several times. The large-scale objects have weak semantics at the coarser level since the receptive field is not enough. Based on this observation, for each query (or reference) scale, we want to select a feature subset from the whole feature set, which is beneficial for the query features. Those selected features are used to enhance the semantic-level or spatial-level information of the query features.

Specifically, we propose a Scale-level Feature Selection (SFS) module to adaptively select informative features subset for each reference scale. The architecture of the SFS module is illustrated in Figure 2 (a). First, we combine all multi-scale features into one-dimensional sequences representations via a re-arrange operation, which flattens all features into one-dimensional sequences and concatenates them along sequence dimensions. This process can be formulated as follow:

$$\mathcal{F} = [\phi(\widehat{\mathcal{F}}_1), \phi(\widehat{\mathcal{F}}_2), \phi(\widehat{\mathcal{F}}_3), \phi(\widehat{\mathcal{F}}_4)] \in \mathbb{R}^{L \times D}, \quad (1)$$

where $L = \sum_{i=1}^4 \frac{HW}{2^{2i+2}}$, and $[\cdot]$ denotes sequence-wise concatenation, and ϕ denotes the reshape operation.

Then, we employ an MLP module to predict the importance scores of each feature in a dynamic way, getting score vector $\mathcal{P} \in \mathbb{R}^{L \times 4}$. This process is formulated as follow:

$$\mathcal{P} = \text{Softmax}(\text{MLP}(\mathcal{F})). \quad (2)$$

Here, $\mathcal{P}_i^j \in [0, 1]$ ($j = 0, 1, \dots, L - 1$) is the element in the i -th row and j -th column of \mathcal{P} , which means the importance of pixel-level feature vector $\mathcal{F}_j \in \mathbb{R}^D$ to scale-level features $\widehat{\mathcal{F}}_i$. $\mathcal{P}_i \in \mathbb{R}^L$ is the i -th column of \mathcal{P} , and means the importance of the all feature \mathcal{F} to scale-level features $\widehat{\mathcal{F}}_i$. For any scale-level features $\widehat{\mathcal{F}}_i$, we generate binary decisions $Q_i \in \{0, 1\}^L$ to indicate whether to select each feature via sampling from \mathcal{P}_i . However, the sampling process is non-differentiable. Inspired by the work [20], we utilize the Gumbel-Softmax to sample from the probabilities \mathcal{P}_i , which can be formulated as follow:

$$Q_i = \text{Gumbel-Softmax}(\mathcal{P}_i) \in \{0, 1\}^L. \quad (3)$$

3.3 Full-scale Feature Fusion module

Different from previous feature fusion methods [27,29,22,51,25,48,54], our proposed Full-scale Feature Fusion module can fuse features from all scales and positions simultaneously. Inspired by ViT [12] which extended the standard

Transformer for image classification, we extend the Multi-head Self-Attention (MSA) into Multi-head Cross-Attention (MCA) for modeling the dependencies between the full-scale features and a reference scale’s ones. The architecture of the FFF module is illustrated in Figure 2 (b). For any scale features, we fuse them with the whole feature sets. Specifically, given query features $\widehat{\mathcal{F}}_i$, our FFF module takes it, features \mathcal{F} and decision Q_i as inputs, and employs Multi-head Cross-Attention to model dependencies. Here, we find it challenging to train our model. The decision Q_i has the various number of “1” for different images and scales in a batch, which prevents our model from efficiently parallel computing. To overcome this issue, we introduce an ingenious mask mechanism, which masks the attention score matrix during training. Specifically, we first compute the attention score matrix as follow:

$$\mathbb{A} = \frac{\widehat{\mathcal{F}}_i \mathcal{F}^T}{\sqrt{C}} \in \mathbb{R}^{N_i \times L}, \quad (4)$$

where N_i is the sequence length of features $\widehat{\mathcal{F}}_i$. Then, we repeat the binary decisions Q_i into a mask $M_i \in \mathbb{R}^{N_i \times L}$ via repeating N_i times. Furthermore, we compute normalized attention matrix as follow:

$$\widehat{\mathbb{A}}_{ij} = \frac{\exp(A_{ij})M_{ij}}{\sum_{k=1}^L M_{ik}}. \quad (5)$$

Then we compute the output of MCA as follow:

$$\mathcal{Y}_i = \widehat{\mathbb{A}}\mathcal{F} \in \mathbb{R}^{N_i \times D}. \quad (6)$$

For the Multi-head Cross Attention, the mask is shared by all heads.

During inference, our FFF module takes query features $\widehat{\mathcal{F}}_i$, features \mathcal{F} , and the corresponding importance scores \mathcal{P}_i as the input. First, we conduct a feature selection operation, given a certain proportion $\rho \in (0, 1]$, we select ρL feature vectors from the candidate set \mathcal{F} , according to the estimated importance score \mathcal{P}_i . For brevity, we denote the chosen features as $\mathcal{F}^s \in \mathbb{R}^{\rho L \times D}$. Then, we take features \mathcal{F}^s as Key and Value embedding and conduct Multi-Head Cross Attention without using mask mechanism in Eq. (5).

Efficient Variants. According to Eq. (4-6), the computational complexity of our MCA is $O(N_i L)$ during training, in order to further improve its efficiency, we employ a projection mechanism to map the input feature sequence into a shorter one. Especially, it first adaptively generates a project matrix Q_i for each query feature (or reference scales), which can be formulated as follow:

$$Q_i = f(\widehat{\mathcal{F}}_i) \in \mathbb{R}^{N_i \times N'_i}, \quad (7)$$

where $N'_i = \frac{N_i}{r}$, r is the ratio of reduction. After getting projection matrices, we map the input feature sequence $\widehat{\mathcal{F}}_i$ into a compact features $\widehat{\mathcal{F}}'_i \in \mathbb{R}^{N'_i \times D}$. This process can be formulated as follow:

$$\widehat{\mathcal{F}}'_i = Q_i^T \widehat{\mathcal{F}}_i, \quad (8)$$

Then, we conduct MCA, obtaining its output features:

$$\mathcal{Y}'_i = MCA(\widehat{\mathcal{F}}'_i, \mathcal{F}, \mathcal{P}_i) \in \mathbb{R}^{N'_i \times D}, \quad (9)$$

where $MCA(\cdot, \cdot, \cdot)$ indicates the computation process of Eq. (4–6). Finally, we employ the projection matrix to re-project the output \mathcal{Y}'_i to the original length, which can be formulated as follow:

$$\mathcal{Y}_i = \mathcal{Q}_i \mathcal{Y}'_i \in \mathbb{R}^{N_i \times D}. \quad (10)$$

3.4 Loss Function

First, we employ a ratio loss to constrain the ratio of the selected features. Given a target ratio $\rho \in (0, 1]$, the ratio loss is defined as follow:

$$\mathcal{L}_{ratio} = \frac{1}{S} \sum_{i=1}^S \|\rho - \frac{1}{L} \sum_{j=1}^L \mathcal{P}_i^j\|^2, \quad (11)$$

where $S = 4$ indicates the number of scales. Following previous works [50,49,47,41], we also add an auxiliary segmentation head attached to Stage 3 of the backbone network to promote the training of our model. Therefore, the objective of our FeSeFormer consists of a ratio loss \mathcal{L}_{ratio} , an auxiliary segmentation loss \mathcal{L}_{aux} , and main segmentation loss \mathcal{L}_{seg} , which can be formulated as:

$$\mathcal{L} = \mathcal{L}_{seg} + \alpha \mathcal{L}_{ratio} + \beta \mathcal{L}_{aux}, \quad (12)$$

where, α and β are hyper-parameters. The selection of α is discussed in the experiment section. Following previous work [50,49,47,41], we set the weight β of auxiliary loss to 0.4.

4 Experiments

We first introduce the datasets and implementation details. Then, we compare our method with the recent state-of-the-arts on four challenging semantic segmentation benchmarks. On the one hand, our experimental results shows that our method is very effective in fusing multi-scale features from Transformer-based backbones. On the other hand, our method can work for fusing multi-scale features from CNN-based backbones. Finally, extensive ablation studies and visualizations analysis are conducted to evaluate the effectiveness of our approach.

4.1 Datasets

PASCAL Context[35] is an extension of the PASCAL VOC 2010 detection challenge. It contains 4998 and 5105 images for training and validation, respectively. Following previous works, we evaluate the most frequent 60 classes (59 categories with background).

Table 1: Comparison with the state-of-the-art methods on the ADE20K dataset. “†” means the resolution of the image is 640×640 , otherwise 512×512 .

Method	Backbone	GFLOPs	Params	mIoU (%)
EncNet[49]	ResNet-101	219	55M	44.65
OCRNet[47]	HRNet-W48	165	71M	44.88
CCNet[19]	ResNet-101	278	69M	45.04
ANN[54]	ResNet-101	263	65M	45.24
PSPNet[50]	ResNet-269	256	68M	45.35
FPT[48]	ResNet-101	479	135M	45.90
DeepLabV3+[6]	ResNet-101	255	63M	46.35
FeSeFormer (ours)	ResNet-101	251	95M	46.56
SETR[52]	ViT-L	214	310M	50.28
DPT[?]	ViT-Hybrid	328	338M	49.02
MCIBI[21]	ViT-L	-	-	50.80
SegFormer[43]	MiT-B5	183	85M	51.80
FeSeFormer (ours)	Swin-L	348	237M	53.33
Swin-UpperNet[32]†	Swin-L	647	234M	53.50
Segmenter[37]†	ViT-L	380	342M	53.60
FeSeFormer (ours)†	Swin-L	587	240M	54.43

ADE20K[53] is a very challenging benchmark including 150 categories and diverse scenes with 1,038 image-level labels, which is split into 20000 and 2000 images for training and validation.

COCO-Stuff 10K[2] is a large scene parsing benchmark, which has 9000 training images and 1000 testing images with 182 categories (80 objects and 91 stuffs).

Cityscapes[9] carefully annotates 19 object categories of urban landscape images. It contains 5K finely annotated images, split into 2975 and 500 for training and validation.

4.2 Implementation details

We employ Swin Transformer [32] pretrained on ImageNet as the backbone. The channel D of features $\widehat{\mathcal{F}}_i$ is set to 256, the weight α of is set to 0.4, and the target ratio ρ is set to 0.6. The number of heads on MCA is set to 8. During training, data augmentation consists of three steps:(i) random horizontal flipping, (ii) we apply random resize with the ratio between 0.5 and 2, (iii) random cropping (480×480 for Pascal Context, 512×512 for ADE20K and COCO-Stuff-10K, and 768×768 for Cityscapes). For optimization, following prior works [32], we employ AdamW [34] to optimize our model with 0.9 momenta and 0.01 weight decay. The batch size is set to 8 for Cityscapes, and 16 for other datasets. We set the initial learning rate at 0.00006 on ADE20K and Cityscapes, and 0.00002 on Pascal Context and COCO-Stuff-10K. The total iterations are set to 160k, 60k, 80k, and 80k for ADE20K, COCO-Stuff-10k, Cityscapes, and PASCAL-Context, respectively. For evaluation, we follow previous works [32,51] to average

Table 2: Comparison with the state-of-the-art approaches on Cityscapes. “SS” and “MS” indicate single-scale inference and multi-scale inference, respectively. “†” means the input resolution is 1024×1024 , otherwise 768×768 on the Cityscapes dataset.

Method	Backbone	mIoU(SS)	mIoU(MS)
EncNet[49]	ResNet-101	76.10	76.97
PSPNet[50]	ResNet-101	78.87	80.04
GCNet[3]	ResNet-101	79.18	80.71
DNLNet[45]	ResNet-101	79.41	80.68
CCNet[19]	ResNet-101	79.45	80.66
DANet[14]	ResNet-101	80.47	82.02
ANN[54]	ResNet-101	-	81.30
MaskFormer[8]	ResNet-101	-	81.40
OCRNet[47]	HRNet-w48	80.70	81.87
FPT[48]	ResNet-101	81.70	-
FeSeFormer (ours)	ResNet-101	80.25	82.13
Segmenter[37]	DeiT-B	79.00	80.60
Segmenter[37]	ViT-L	-	81.30
SETR-PUP [52]	ViT-L	79.34	82.15
FeSeFormer (ours)	Swin-L	83.08	83.64
SegFormer[43] [†]	MiT-B5	82.40	84.00
FeSeFormer (ours)[†]	Swin-L	83.61	84.46

the multi-scale (0.5, 0.75, 1.0, 1.25, 1.5, 1.75) predictions of our model. The performance is measured by the standard mean intersection of union (mIoU) in all experiments. Considering the effectiveness and efficiency, we adopt the Swin-T [32] as the backbone in the ablation study, and report single-scale testing results.

4.3 Comparisons with the State-of-the-art

Results on ADE20K. Table 1 reports the comparison with the state-of-the-art methods on the ADE20K validation set. When Swin-Transformer is used as the backbone, our FeSeFormer is +1.53% mIoU higher (53.33% vs. 51.80%) than SegFormer with the same input size (512×512). While recent methods [37,32] showed that using a larger resolution (640×640) can bring more improvements, Our FeSeFormer is +0.73% and +0.83% mIoU higher than Segmenter [37] and Swin-UperNet [32], respectively. These results demonstrate that fusing multi-scale features from all scales is very effective for improving segmentation performance. Although our work focuses on how to fuse multi-scale features from Transformer-based backbone, we also conduct experiments with CNN-based backbone, e.g., ResNet-101 [17]. it can be seen that our FeSeFormer (ResNet-101) also outperforms DeepLabV3+ (ResNet-101) (46.56% vs. 46.35%), which is the best segmentation model among methods employing ResNet-101 as the backbone.

Table 3: Comparison with the state-of-the-art methods on PASCAL Context and COCO-Stuff 10K.

(a) Results on the PASCAL Context.

Method	Backbone	mIoU(%)
PSPNet[50]	ResNet-101	47.15
DeepLabV3+[6]	ResNet-101	48.26
DANet[14]	ResNet-101	52.60
ANN[54]	ResNet-101	52.80
EMANet[24]	ResNet-101	53.10
SVCNet[10]	ResNet-101	53.20
ACNet[11]	ResNet-101	54.10
GFFNet[25]	ResNet-101	54.20
Efficientfcn[31]	ResNet-101	54.30
APCNet[16]	ResNet-101	54.70
OCRNet[47]	ResNet-101	54.80
RecoNet[7]	ResNet-101	54.80
GINet[41]	ResNet-101	54.90
FeSeFormer (ours)	ResNet-101	55.23
SETR[52]	ViT-L	55.83
Swin-UperNet [40]	Swin-L	57.29
FeSeFormer (ours)	Swin-L	58.91

(b) Results on the COCO-Stuff 10K.

Method	Backbone	mIoU(%)
PSPNet[50]	ResNet-101	38.86
OCRNet[47]	ResNet-101	39.50
DANet[14]	ResNet-101	39.70
SVCNet[10]	ResNet-101	39.60
MaskFormer[8]	ResNet-101	39.80
EMANet[24]	ResNet-101	39.90
SpyGR[23]	ResNet-101	39.90
ACNet[11]	ResNet-101	40.10
GINet[41]	ResNet-101	40.60
OCRNet[47]	HRNetV2-W48	40.50
RegionContrast [18]	ResNet-101	40.70
RecoNet[7]	ResNet-101	41.50
FeSeFormer (ours)	ResNet-101	41.73
MCIBI[21]	ViT-L	44.89
Swin-UperNet [40]	Swin-L	47.71
FeSeFormer (ours)	Swin-L	49.80

Results on Cityscapes. Table 2 shows the comparative results on the Cityscapes validation set. Among methods that employed Transformer-based backbones, SETR-PUP achieved the best accuracy, which employed a huge backbone ViT-Large [12] and a progressive upsampling strategy for getting high-resolution predictions. Our FeSeFormer (Swin-L) is superior to it (83.64% vs. 82.15%). Furthermore, recent SegFormer [43] was trained with higher resolution (1024×1024), and achieved very stronger performance. For a fair comparison, we train our model with the same input size. We can see that our FeSeFormer[†] (Swin-L) is +1.21% and +0.46% mIoU higher than SegFormer (MiT-B5) under the single-scale and multi-scale testing, respectively. Besides, we also employed CNN-based backbone, e.g., ResNet-101 to extract features. It can be seen that our FeSeFormer (ResNet-101) is +0.26% mIoU higher (82.13% vs. 81.87%) than OCRNet (ResNet-101). We also compare with the closely related FPT [48]. According to the results in Table 2, our FeSeFormer is +1.38% mIoU higher than FPT (83.08% vs. 81.70%), which proposed Grounding Transformer and Rendering Transformer to fuse the higher-level and lower-level feature in a bidirectional fashion. Finally, we find an interesting phenomenon that our multi-scale test results are only slightly higher than the single-scale test results ($\sim 0.5\%$), while the former is usually $\sim 1.0\%$ higher than the latter for other methods. This also shows that our multi-scale feature fusion is better than other methods. These results demonstrate the effectiveness of fusing features of all scales for each query feature simultaneously.

Results on PASCAL Context. As shown in Table 3(a), we compare our method with the state-of-the-art models on PASCAL Context. From these re-

Table 4: Ablation study on PASCAL Context testing set. Our baseline model consists of Swin-T, Top-down Feature Pyramid Fusion module, FCN Head. FLOPs is measured with the input size of 480×480 . “PM” means projection mechanism.

Baseline	FFF	SFS	PM	FLOPs	Params	mIoU(%)
✓	-	-	-	54.0G	33M	46.75
✓	✓	-	-	75.9G	41M	48.35
✓	✓	✓	-	74.2G	42M	49.06
✓	✓	✓	✓	73.8G	39M	49.33

Table 5: Comparison of different weights of Ratio Loss on PASCAL Context testing set.

α	0.0	0.2	0.4	0.6	0.8	1.0
mIoU (%)	49.11	49.14	49.33	49.02	48.82	48.70

sults, it can be seen that our FeSeFormer is +3.08 mIoU higher (58.91 vs. 55.83) than the very famous SETR [52], which is the first work to use Transformer for semantic segmentation. Furthermore, our method also outperforms the recent work Swin-UperNet [32] (58.91 vs. 57.29) with the same backbone network. Besides, our FeSeFormer (ResNet-101) is +0.33 mIoU higher (55.23 vs. 54.90) than GINet [41], which achieved the best performance among CNN-based methods.

Results on COCO-Stuff 10K. Table 3(b) shows the comparison results on the COCO-Stuff 10K testing set. It can be seen that our FeSeFormer configured with Swin Transformer can achieve 49.80%, and outperforms the previous best Swin-UperNet (49.80 vs 47.71). Furthermore, our method is +4.91 mIoU higher than MCIBI [21] (49.80 vs. 44.89), which introduced a feature memory module to store the dataset-level contextual information of various categories and aggregated the dataset-level representations for each pixel. Besides, our FeSeFormer (ResNet-101) can achieve 41.73% mIoU, which outperforms RecoNet (ResNet-101), which is the best segmentation model among those methods that employed ResNet-101 as the backbone.

4.4 Ablation study

In this section, we give extensive experiments to show the efficacy of our method. We also give several design choices and show their effects on the results. Our baseline model takes Swin Transformer [32] as the backbone, and use FPN [29] to perform top-down feature fusion for getting feature maps with stride=8, followed by a FCN Head (consisting of two 3×3 Conv. and one 1×1 Conv). All the following experiments adopt Swin-T as the backbone, trained on PASCAL Context training set for 80K iterations.

Efficacy of SFS and FFF module. According to the results in Table 4, we can see that the baseline model can achieve 46.75% mIoU on the PASCAL

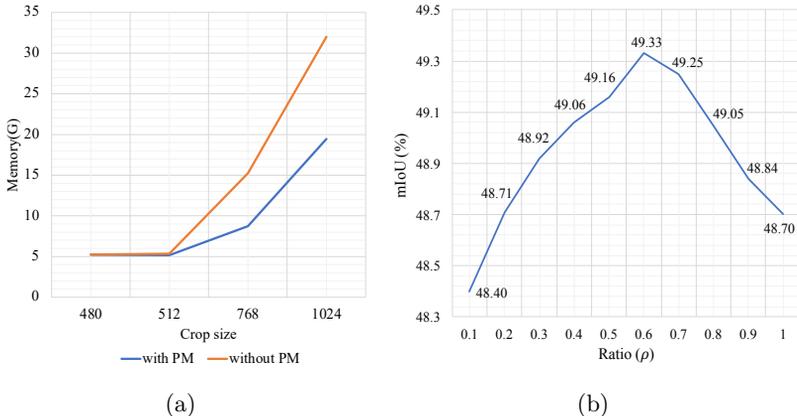


Fig. 3: (a) Comparison of memory overhead under the different input resolutions. “PM” indicates the project mechanism. (b) Comparisons of the performance of different target ratio on PASCAL Context testing set.

Table 6: Comparisons of efficiency and accuracy with other multi-scale feature fusion methods on the PASCAL Context testing set. We report the FLOPs and Params of decoders, relative to the backbone. The input resolution is set to 480×480 .

Decoder	GFLOPs	Params	mIoU (%)
UperNet [42]	187G	37M	45.06
Semantic-FPN [22]	112G	54M	46.72
SETR-MLA [52]	13G	3M	46.77
DPT [36]	97G	17M	46.79
GFFNet [25]	85G	17M	48.04
FPT [48]	414G	92M	48.31
Ours	52G	12M	49.33

testing set. By adding the FFF module, the performance is improved by +1.6% ($46.75\% \rightarrow 48.35\%$). When adding a SFS module, the performance is further improved by +0.71% ($48.35\% \rightarrow 49.06\%$). Furthermore, we deploy the projection mechanism on features $\hat{\mathcal{F}}_1$ for reducing computational costs, since its length is very long ($\frac{H}{8} \times \frac{W}{8}$). It can be seen that our model with the project mechanism can achieve 49.33% mIoU, which not only reduces the computational overhead but also brings a slight performance improvement (49.33% vs. 49.06%). Besides, we give a comparison of memory overhead under the different input resolutions in Fig. 3 (a). It can be seen that the larger the input resolution, the more significant our projection mechanism reduces the memory overhead.

Selection of the Weight α of Ratio Loss. To study the effect of Ratio Loss, we test different weights $\alpha = \{0.0, 0.2, 0.4, 0.6, 0.8, 1.0\}$. As shown in Table 5, it can be seen that $\alpha = 0.4$ can yield the best accuracy (49.33% mIoU).

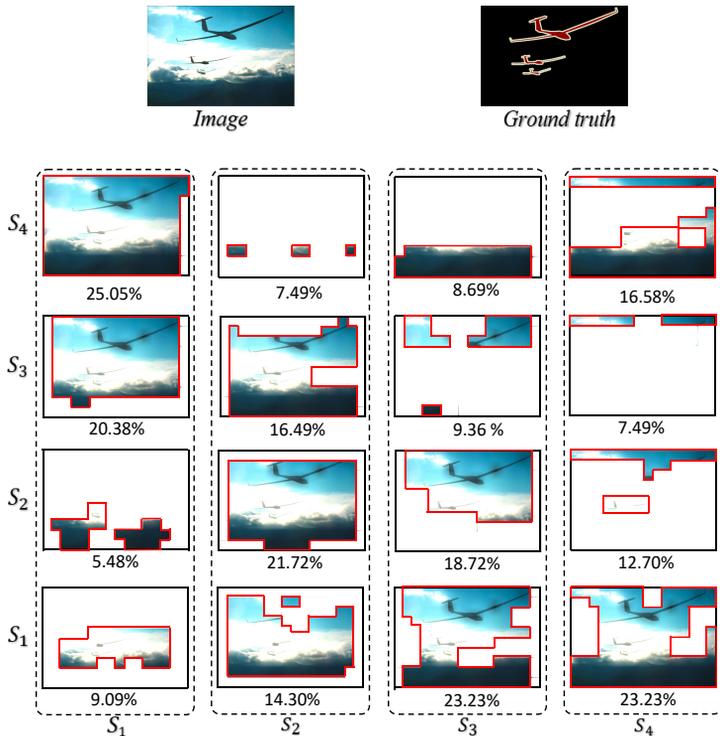


Fig. 4: Visualization of multi-scale feature selection. The red polygon represents the selection area. We draw the choosing position and ratio for each scale. The i -th column shows the feature selection of query scale S_i .

Selection of the Ratio of Choosing Features. Here, we study the effect of choosing different target ratios ρ . As shown in Fig. 3 (b), we can see that $\rho = 0.6$ can yield the best performance 49.33% mIoU, outperforming $\rho = 1.0$. Besides, decreasing the target ratio from 0.6 to 0.1, the segmentation accuracy shows a decreasing trend. These results demonstrate that it is necessary to select a subset of informative features. Furthermore, to better understand our method, we visualize the distribution of features selected of different scales in a sample in Fig. 4. It can be seen that query scale S_1 mainly select features from scale S_3 and S_4 , while query scale S_3 mainly select features from scale S_1 and S_2 . We provide more examples and analyses in the supplementary materials.

Comparisons with Related Multi-scale Feature Fusion Methods. Next, we compare our method with other multi-scale feature fusion methods in Table 6. Among progressive local fusion methods (UperNet, Semantic-FPN, SETR-MLA, and DPT), DPT achieved the best accuracy (46.79%). Our method is +2.54% mIoU higher (49.33% vs. 46.79%) than it. Besides, our FeSeFormer is +1.29% mIoU higher (49.33% vs. 48.04%) than non-progressive local fusion method, GFFNet. FPT is the most related to our work, which proposed Ground-

ing Transformer and Rendering Transformer to fuse the higher-level and lower-level features in a bidirectional fashion. Our method outperforms it by +1.02% mIoU with less computational overhead.

5 Conclusion

We have developed the Feature Selective Transformer (FeSeFormer) for semantic image segmentation. The core contributions of FeSeFormer are the proposed Scale-level Feature Selection (SFS) and Full-scale Feature Fusion (FFF) modules. The former chooses an informative subset from the multi-scale feature set for each scale. The latter fuses features of all scales in a dynamic way. Extensive experiments on PASCAL-Context, ADE20K, COCO-Stuff 10K, and Cityscapes have shown that our FeSeFormer can outperform the state-of-the-art methods in semantic image segmentation, demonstrating that our FeSeFormer can achieve better results than previous multi-scale feature fusion methods.

References

1. Alhaija, H.A., Mustikovela, S.K., Mescheder, L., Geiger, A., Rother, C.: Augmented reality meets deep learning for car instance segmentation in urban scenes. In: British machine vision conference. vol. 1, p. 2 (2017)
2. Caesar, H., Uijlings, J., Ferrari, V.: Coco-stuff: Thing and stuff classes in context. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1209–1218 (2018)
3. Cao, Y., Xu, J., Lin, S., Wei, F., Hu, H.: Gcnet: Non-local networks meet squeeze-excitation networks and beyond. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops. pp. 0–0 (2019)
4. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: 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)
5. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587* (2017)
6. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Proceedings of the European conference on computer vision (ECCV). pp. 801–818 (2018)
7. Chen, W., Zhu, X., Sun, R., He, J., Li, R., Shen, X., Yu, B.: Tensor low-rank reconstruction for semantic segmentation. In: European Conference on Computer Vision. pp. 52–69. Springer (2020)
8. Cheng, B., Schwing, A.G., Kirillov, A.: Per-pixel classification is not all you need for semantic segmentation. *arXiv preprint arXiv:2107.06278* (2021)
9. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3213–3223 (2016)
10. Ding, H., Jiang, X., Shuai, B., Liu, A.Q., Wang, G.: Semantic correlation promoted shape-variant context for segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8885–8894 (2019)

11. Ding, X., Guo, Y., Ding, G., Han, J.: Acnet: Strengthening the kernel skeletons for powerful cnn via asymmetric convolution blocks. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1911–1920 (2019)
12. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. ICLR (2021)
13. Feng, D., Haase-Schütz, C., Rosenbaum, L., Hertlein, H., Glaeser, C., Timm, F., Wiesbeck, W., Dietmayer, K.: Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. IEEE Transactions on Intelligent Transportation Systems (2020)
14. Fu, J., Liu, J., Tian, H., Li, Y., Bao, Y., Fang, Z., Lu, H.: Dual attention network for scene segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3146–3154 (2019)
15. Harders, M., Szekely, G.: Enhancing human-computer interaction in medical segmentation. Proceedings of the IEEE **91**(9), 1430–1442 (2003)
16. He, J., Deng, Z., Zhou, L., Wang, Y., Qiao, Y.: Adaptive pyramid context network for semantic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7519–7528 (2019)
17. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
18. Hu, H., Cui, J., Wang, L.: Region-aware contrastive learning for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 16291–16301 (2021)
19. Huang, Z., Wang, X., Huang, L., Huang, C., Wei, Y., Liu, W.: Cnet: Criss-cross attention for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 603–612 (2019)
20. Jang, E., Gu, S., Poole, B.: Categorical reparameterization with gumbel-softmax. arXiv preprint arXiv:1611.01144 (2016)
21. Jin, Z., Gong, T., Yu, D., Chu, Q., Wang, J., Wang, C., Shao, J.: Mining contextual information beyond image for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7231–7241 (2021)
22. Kirillov, A., Girshick, R., He, K., Dollár, P.: Panoptic feature pyramid networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6399–6408 (2019)
23. Li, X., Yang, Y., Zhao, Q., Shen, T., Lin, Z., Liu, H.: Spatial pyramid based graph reasoning for semantic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8950–8959 (2020)
24. Li, X., Zhong, Z., Wu, J., Yang, Y., Lin, Z., Liu, H.: Expectation-maximization attention networks for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9167–9176 (2019)
25. Li, X., Zhao, H., Han, L., Tong, Y., Tan, S., Yang, K.: Gated fully fusion for semantic segmentation. In: AAAI (2020)
26. Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, Alan L Yuille: Semantic image segmentation with deep convolutional nets and fully connected crfs. In: ICLR (2015)
27. Lin, D., Ji, Y., Lischinski, D., Cohen-Or, D., Huang, H.: Multi-scale context intertwining for semantic segmentation. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 603–619 (2018)

28. Lin, D., Shen, D., Shen, S., Ji, Y., Lischinski, D., Cohen-Or, D., Huang, H.: Zigzag-net: Fusing top-down and bottom-up context for object segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7490–7499 (2019)
29. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2117–2125 (2017)
30. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
31. Liu, J., He, J., Zhang, J., Ren, J.S., Li, H.: Efficientfcn: Holistically-guided decoding for semantic segmentation. In: European Conference on Computer Vision. pp. 1–17. Springer (2020)
32. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. arXiv preprint arXiv:2103.14030 (2021)
33. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)
34. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
35. Mottaghi, R., Chen, X., Liu, X., Cho, N.G., Lee, S.W., Fidler, S., Urtasun, R., Yuille, A.: 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. pp. 891–898 (2014)
36. Ranftl, R., Bochkovskiy, A., Koltun, V.: Vision transformers for dense prediction. ICCV (2021)
37. Strudel, R., Garcia, R., Laptev, I., Schmid, C.: Segmenter: Transformer for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 7262–7272 (October 2021)
38. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. arXiv preprint arXiv:1706.03762 (2017)
39. Wang, X., Girshick, R., Gupta, A., He, K.: Non-local neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7794–7803 (2018)
40. Wu, S., Wu, T., Lin, F., Tian, S., Guo, G.: Fully transformer networks for semantic image segmentation. arXiv preprint arXiv:2106.04108 (2021)
41. Wu, T., Lu, Y., Zhu, Y., Zhang, C., Wu, M., Ma, Z., Guo, G.: Ginet: Graph interaction network for scene parsing. In: European Conference on Computer Vision. pp. 34–51. Springer (2020)
42. Xiao, T., Liu, Y., Zhou, B., Jiang, Y., Sun, J.: Unified perceptual parsing for scene understanding. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 418–434 (2018)
43. Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., Luo, P.: Segformer: Simple and efficient design for semantic segmentation with transformers. arXiv preprint arXiv:2105.15203 (2021)
44. Xie, E., Wang, W., Wang, W., Sun, P., Xu, H., Liang, D., Luo, P.: Segmenting transparent object in the wild with transformer. arXiv preprint arXiv:2101.08461 (2021)

45. Yin, M., Yao, Z., Cao, Y., Li, X., Zhang, Z., Lin, S., Hu, H.: Disentangled non-local neural networks. In: European Conference on Computer Vision. pp. 191–207. Springer (2020)
46. Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:1511.07122 (2015)
47. Yuan, Y., Chen, X., Wang, J.: Object-contextual representations for semantic segmentation. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16. pp. 173–190. Springer (2020)
48. Zhang, D., Zhang, H., Tang, J., Wang, M., Hua, X., Sun, Q.: Feature pyramid transformer. In: European Conference on Computer Vision. pp. 323–339. Springer (2020)
49. Zhang, H., Dana, K., Shi, J., Zhang, Z., Wang, X., Tyagi, A., Agrawal, A.: Context encoding for semantic segmentation. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 7151–7160 (2018)
50. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2881–2890 (2017)
51. Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., Fu, Y., Feng, J., Xiang, T., Torr, P.H., et al.: Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. arXiv preprint arXiv:2012.15840 (2020)
52. Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., Fu, Y., Feng, J., Xiang, T., Torr, P.H., et al.: Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6881–6890 (2021)
53. Zhou, B., Zhao, H., Puig, X., Xiao, T., Fidler, S., Barriuso, A., Torralba, A.: Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision* **127**(3), 302–321 (2019)
54. Zhu, Z., Xu, M., Bai, S., Huang, T., Bai, X.: Asymmetric non-local neural networks for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 593–602 (2019)

Appendix

In this appendix, we first compare our method with related multi-scale fusion methods under the different backbones. Then, for demonstrating the generality of our method, we conduct the instance segmentation task on COCO [2]. Finally, we provide some visualization results and analysis.

A Comparisons with Related Multi-scale Feature Fusion Methods

As shown in Table 7, we compare the performance of several different multi-scale feature fusion approaches on the PASCAL context testing set with multiple Swin Transformer [32] variants, including Swin-S, Swin-B, and Swin-L. One can be seen that our method outperforms all related approaches under the different backbones.

Table 7: Comparisons of multi-scale feature fusion methods on the PASCAL Context testing set. The input resolution is set to 480×480 . Single-scale testing is adopted here.

	UperNet[42]	Semantic-FPN[22]	SETR-MLA[52]	DPT[36]	GFFNet[25]	FPT[48]	Ours
Swin-S	51.67	50.49	50.48	50.47	51.07	51.92	52.58
Swin-B	52.52	51.48	51.51	51.55	51.89	52.80	53.12
Swin-L	56.87	56.78	56.62	56.41	56.49	56.97	57.63

B Comparisons with Mask R-CNN on COCO

To demonstrate the generality of the proposed Scale-level Feature Selection (SFS) and Full-scale Feature Fusion (FFF) module, we conduct the instance segmentation task on COCO [30] using the competitive Mask R-CNN model as the baseline, and compare our method with FPN under the Mask R-CNN framework. We report the results in terms of mask AP in Table 8. We can see that our method outperforms Mask R-CNN (FPN) in all metrics.

C Visualizations

For better understanding our method, we visualize feature selection of query features. Examples from PASCAL Context [35], ADE20K [53], COCO-Stuff 10K [2], and Cityscapes[9] are shown in Figure 5, 6, 7, and 8, respectively. The i -th column shows the feature selection of query scale S_i . From these results, one can be seen that high-level semantic features tend to select low-level features with detailed spatial information, and vice versa.

Table 8: Comparisons of FPN and our method with Mask R-CNN on COCO dataset.

Method	backbone	AP^b	AP_{50}^b	AP_{75}^b	AP^m	AP_{50}^m	AP_{75}^m
FPN [29]	Swin-T	43.7	66.6	47.7	39.8	63.3	42.7
SFS+FFF(ours)	Swin-T	45.4	68.1	49.4	41.9	65.1	44.6
FPN [29]	Swin-S	46.5	68.7	51.3	42.1	65.8	45.2
SFS+FFF(ours)	Swin-S	47.7	70.0	52.8	43.5	67.1	46.6
FPN [29]	Swin-B	46.9	69.2	51.6	42.3	66.0	45.5
SFS+FFF(ours)	Swin-B	48.5	70.7	53.2	43.9	67.9	47.3

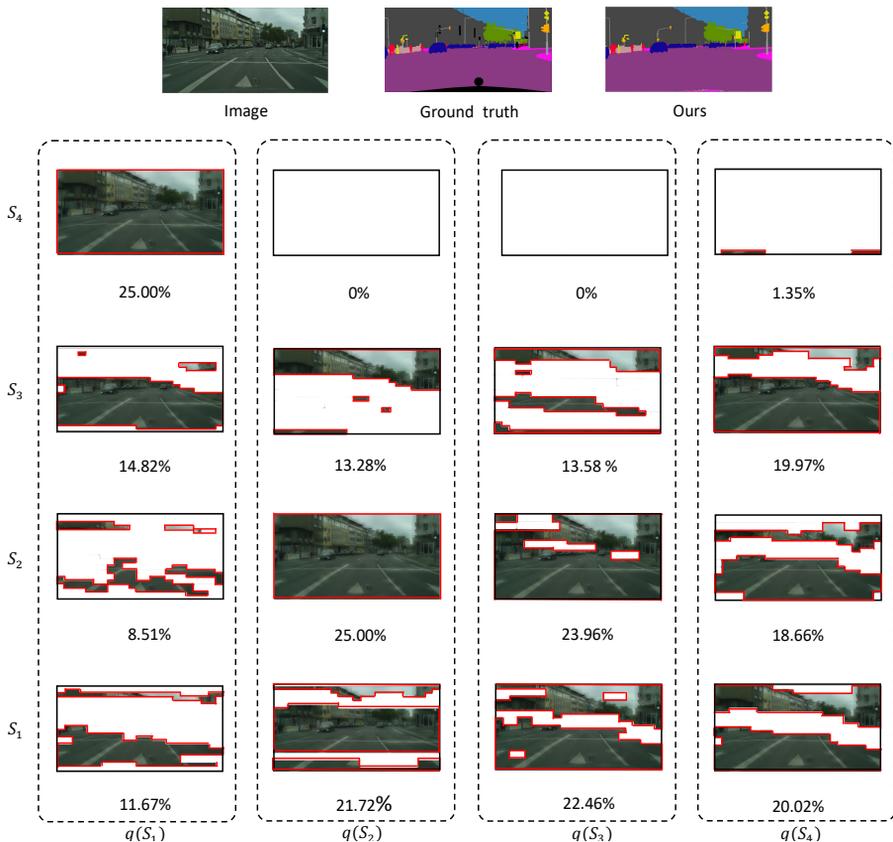


Fig. 5: Visualization of multi-scale feature selection on Cityscapes dataset. $q(S_i)$ indicates taking features from the stage or scale S_i as a query. The i -th column shows the feature selection of query scale S_i . The red polygon represents the selection area.



Fig. 6: Visualization of multi-scale feature selection on COCO-Stuff 10K dataset. $q(S_i)$ indicates taking features from the stage or scale S_i as a query. The i -th column shows the feature selection of query scale S_i . The red polygon represents the selection area.

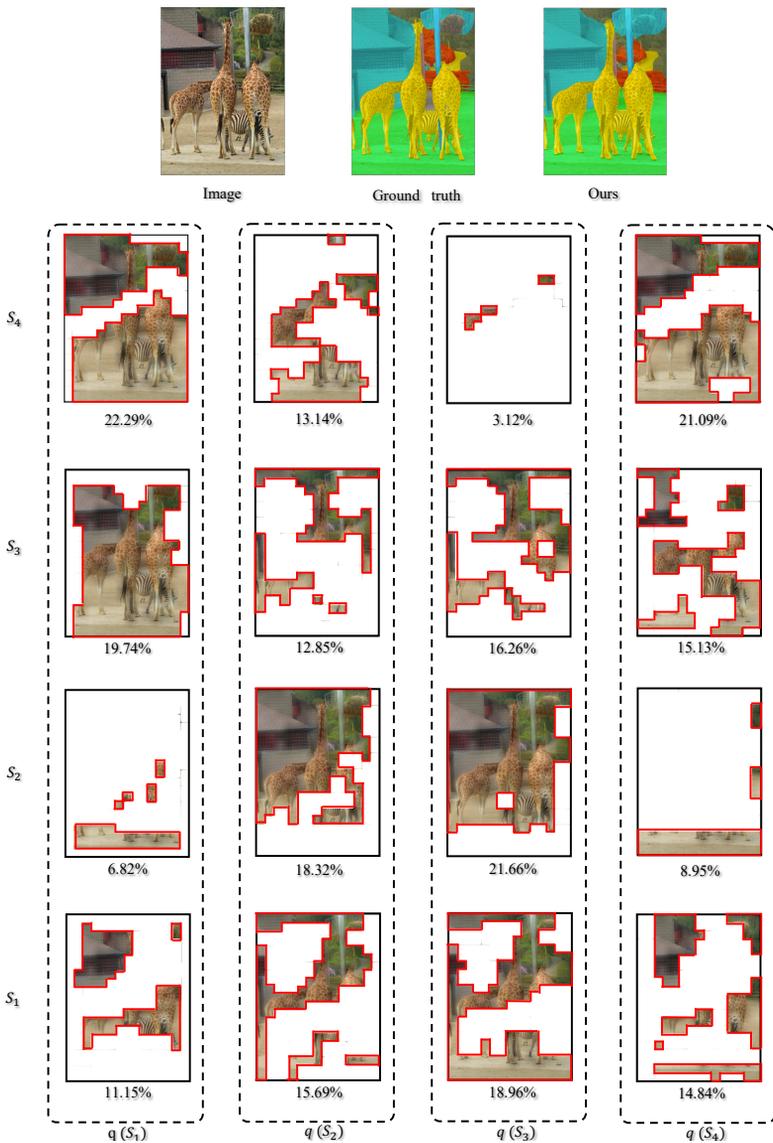


Fig. 7: Visualization of multi-scale feature selection on ADE20K dataset. $q(S_i)$ indicates taking features from the stage or scale S_i as a query. The i -th column shows the feature selection of query scale S_i . The red polygon represents the selection area.

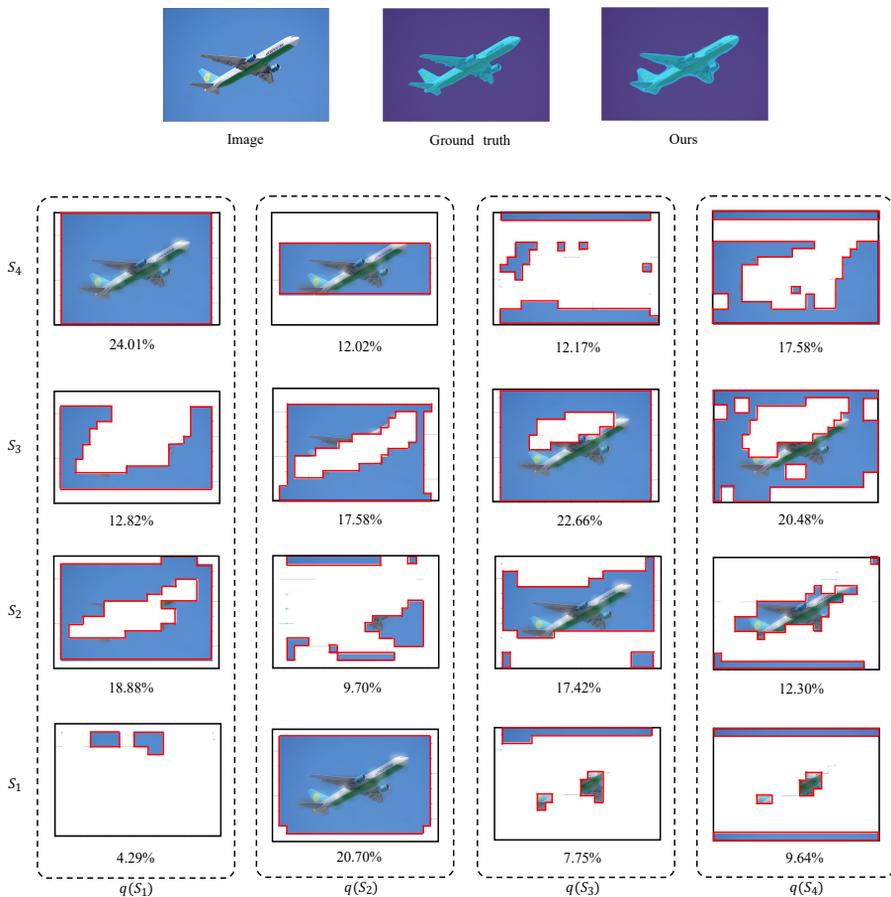


Fig. 8: Visualization of multi-scale feature selection on PASCAL Context dataset. $q(S_i)$ indicates taking features from the stage or scale S_i as a query. The i -th column shows the feature selection of query scale S_i . The red polygon represents the selection area.