

ElegantSeg: End-to-End Holistic Learning for Extra-Large Image Semantic Segmentation

Wei Chen Yansheng Li* Bo Dang Yongjun Zhang
 Wuhan University, Wuhan, China
 {weichenrs, yansheng.li, bodang, zhangyj}@whu.edu.cn

Abstract

This paper presents a new paradigm for *Extra-large image semantic Segmentation*, called *ElegantSeg*, that capably processes holistic extra-large image semantic segmentation (ELISS). The extremely large sizes of extra-large images (ELIs) tend to cause GPU memory exhaustion. To tackle this issue, prevailing works either follow the global-local fusion pipeline or conduct the multi-stage refinement. These methods can only process limited information at one time, and they are not able to thoroughly exploit the abundant information in ELIs. Unlike previous methods, *ElegantSeg* can elegantly process holistic ELISS by extending the tensor storage from GPU memory to host memory. To the best of our knowledge, it is the first time that ELISS can be performed holistically. Besides, *ElegantSeg* is specifically designed with three modules to utilize the characteristics of ELIs, including the multiple large kernel module for developing long-range dependency, the efficient class relation module for building holistic contextual relationships, and the boundary-aware enhancement module for obtaining complete object boundaries. *ElegantSeg* outperforms previous state-of-the-art on two typical ELISS datasets. We hope that *ElegantSeg* can open a new perspective for ELISS. The code and models will be made publicly available.

1. Introduction

With the development of photography and sensor technologies, human beings can get access to extra-large images (ELIs) with millions or even billions of pixels, consisting of satellite images [9, 32], histopathological images [6, 33], and natural images [4], which will benefit a wild range of scientific applications, such as urban planning and designing [28, 46], disease monitoring and diagnosing [6, 33].

As one of the most fundamental low-level vision tasks, semantic segmentation (SS) provides a basic but important understanding of images. From the earliest fully con-

*corresponding author.

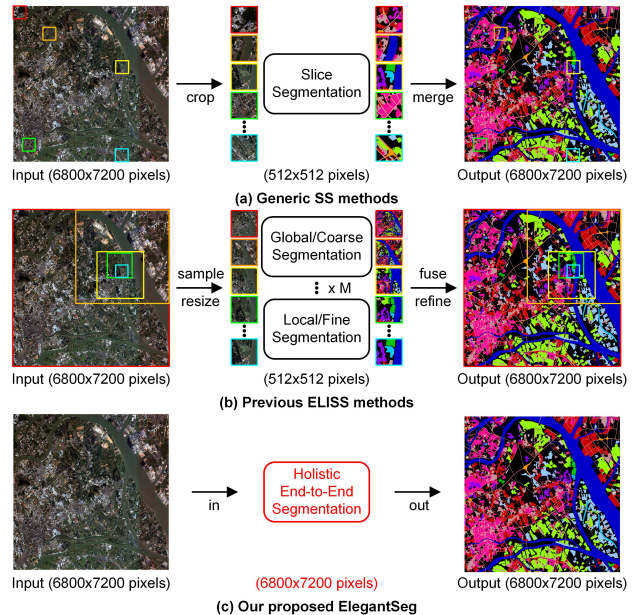


Figure 1. **Comparison of the pipelines for ELISS.** (a) Generic SS methods can only process cropped slices from the ELIs and merge the results of these slices into the whole output. (b) Previous ELISS methods tend to train multiple models with different-size samples, and then conduct either the global-to-local fusion or the coarse-to-fine refinement on multi-stage predictions. (c) Our proposed *ElegantSeg* can perform ELISS by **end-to-end holistic learning**, without merging, fusion, or refinement.

volutional networks [26] to the most recent Transformers [24, 25, 31, 38, 45], SS has witnessed an explosion of deep learning-based techniques of continuously growing capability and capacity [2, 18, 34, 37, 39, 43, 44]. However, despite significant interest in generic SS methods, the progress of extra-large image semantic segmentation (ELISS) lags behind the SS of generic-size images, with only preliminary explorations having been done [3, 15, 17, 20, 22]. We ask: *what makes ELISS different from generic SS?* We attempt to answer this question from the following perspectives:

(i) Compared to generic-size images with limited con-

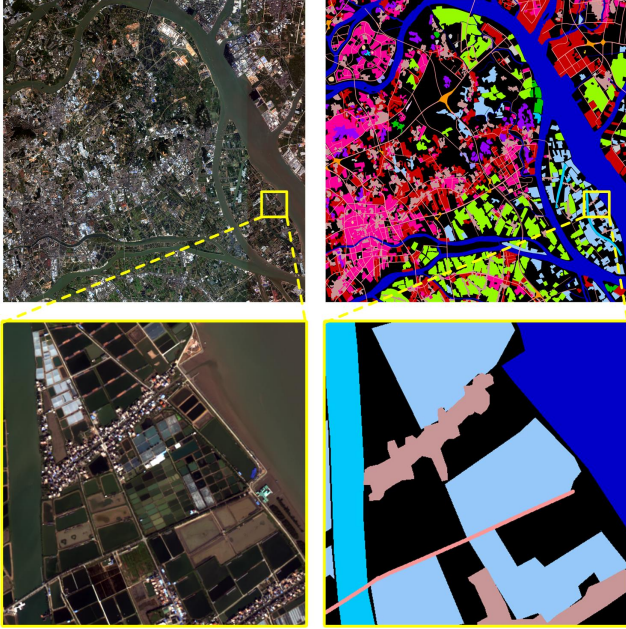


Figure 2. **The necessity of holistic ELISS.** One can see from the lower two images that the cropped slice may lead to ambiguous prediction. Even human beings can not promise to definitely distinguish the **pond** on the left side and the **river** on the right side. Nevertheless, aided by the holistic perspective of the ELI, they can be distinguished with complete structure and rich spatial context. Best viewed on a high-resolution display.

texts, ELIs are able to provide both macroscopic structures and microcosmic details. As illustrated in Fig. 2, the global view can offer complete boundaries of the ELI, while the local view helps to bring fine information for segmenting specific regions. The combination of the two views leads to accurate segmentation results, which can be seen from the results in Fig. 8. This is also similar to the mode of human vision [8, 16].

(ii) ELISS is not only allocating a category to each pixel in the image but also pursuing high-level semantics from a more general perspective. Taking satellite images as examples, rather than caring about object-level categories (cars, buildings, trees, etc.), ELISS requires methods to model the long-range relationship among various regions for SS with complicated categories. For example, groups of houses and trees in ELISS are going to be identified as an “urban area” [1]. However, the pixel-by-pixel SS may lead to confusing results under that circumstance.

(iii) ELIs tend to contain millions or even billions of pixels, which inevitably brings the unbearable cost of memory and computation under the constraint of limited GPU memory during ELISS. To fit the ELIs into the model and perform holistic learning is a huge challenge, which is far away from being solved [3, 15, 17, 20, 22].

Driven by this analysis, we present a simple but effec-

tive paradigm named ElegantSeg for ELISS. Unlike previous ELISS methods [3, 15, 17, 20, 22], our ElegantSeg can holistically process ELISS by utilizing both GPU memory and host memory. With the aid of three specially-designed modules including the multiple large kernel (MLK) module, the efficient class relation (ECR) module, and the boundary-aware enhancement (BAE) module, the characteristics of ELIs are thoroughly exploited. Extensive experiments on two typical ELISS datasets [9, 32] demonstrate the superiority of ElegantSeg over previous state-of-the-arts (SOTAs).

The main contributions are summarized as follows:

- To the best of our knowledge, we propose an end-to-end holistic learning framework for the first time to solve ELISS, named ElegantSeg.
- Considering the characteristics of ELIs, we design three specific modules to enhance the segmentation performance of ElegantSeg.
- Extensive experiments and analyses on two typical ELISS datasets indicate the efficacy of ElegantSeg and its superiority over previous SOTAs.

2. Related work

2.1. Generic semantic segmentation

Since FCN [26], generic SS methods have been explicitly exploiting multi-scale feature fusion to surpass the limit of local receptive field by convolution layer and to capture the contextual information at multiple scales [2, 34, 43]. DeepLabV3Plus [2] applies several parallel atrous convolutions with different rates, while PSPNet [43] performs pooling operations at different grid scales. HRNet [34] connects high-to-low resolution convolutions in parallel and conducts multi-scale fusions across parallel convolutions to produce strong and spatially precise high-resolution representations. Recently, the computer vision domain has been witnessing the breakthrough in Vision Transformers [13], which shows huge potential in SS for capturing the long-range contextual dependency by adapting self-attention to capture the global spatial context in generic SS tasks [24, 25, 38, 45]. However, these generic SS methods can not adjust to ELISS due to the constraint of GPU memory.

2.2. Extra-large image semantic segmentation

From GLNet [3], ELISS (also called ultra-high resolution image segmentation) has been studied for years [3, 10, 11, 15, 20, 22, 36]. Following a similar pattern, these methods train multiple models with different-sized samples and conduct the global-to-local fusion on the predictions. Others adopt the scheme of refinement [4, 17, 18, 21, 23, 29, 40] and perform coarse-to-fine refinement on multi-stage predictions for better segmentation results. Nevertheless, these

ELISS methods do not thoroughly exploit the abundant information in ELIs. Although these methods try to exploit the characteristics of ELIs, there is still a long way to go.

2.3. Potential techniques for holistic ELISS

Large kernel convolution DeepLabV3Plus [2] tries to use dilated convolution to reach a larger receptive field. However, it leads to severe grid effects. Recently, RepLKNet [12] rethinks the use of large convolutional kernels [27] for a larger effective receptive field, which is helpful for capturing long-range dependency. The depth-wise separable convolution [5] and group convolution [41] are used to reduce the cost of memory and computation.

Self-attention Pioneers [35, 39, 42] explore self-attention for capturing global contexts. However, the cost of self-attention is much too heavy for ELIs since they are quite large. Fortunately, [30] proposes efficient attention to dramatically reduce the cost of memory and computation. Which is required to be exploited in ELISS.

3. Method

3.1. Large model support

Before introducing the details of our ElegantSeg, it is necessary to talk about the core of holistic learning first, that is, the Large Model Support (LMS) toolbox. LMS¹ is a toolbox provided by IBM Watson Machine Learning Community Edition that allows the successful training of deep learning models that would otherwise exhaust GPU memory and abort with "out-of-memory" errors. LMS manages this oversubscription of GPU memory by temporarily swapping tensors to host memory when they are not needed.

The extremely large sizes of ELIs always lead to GPU memory exhaustion. Traditionally, the solution to this problem has been to modify the model until it fits in GPU memory. This approach, however, can negatively impact accuracy – especially if concessions are made by reducing data fidelity or model complexity. With LMS, deep learning models can scale significantly beyond what was previously possible and, ultimately, generate more accurate results.

3.2. Overview

With the aid of sufficient host memory, the limited GPU memory is no longer a constraint for ELISS. LMS helps the ELIs to fit in the model and makes it possible for our ElegantSeg to process holistic learning. Considering both efficacy and effectiveness, we choose HRNet [34] as the backbone model to build our ElegantSeg model. HRNet is a powerful and efficient model which can obtain multi-scale features. After that, we are going to design specific mod-

¹<https://github.com/IBM/pytorch-large-model-support>

ules to utilize the characteristics of ELIs. The overview of ElegantSeg can be seen in Fig. 3.

3.3. Multiple Large Kernel module

Early in Sec. 2.3, we mention that large convolutional kernels can provide a much larger effective receptive field than normal convolutional kernels. LKM [27] and RepLKNet [12] perform large convolutional kernels on SS for generic-size images for a larger spatial context. Inspired by them, we design the multiple large kernel (MLK) module in ElegantSeg. For ELIs, the effective receptive field is much more important for capturing global structures. Therefore, we use up to 31×31 convolution kernel in the MLK module. The MLK module consists of a 1×1 convolution for reducing the dimension, an activation function, L multiple large kernel units (MLKUs), and a 1×1 convolution for raising the dimension, as can be seen in Fig. 4. The MLKU is the core of the MLK module. It consists of a 5×5 convolution, three large-kernel depth-wise separable convolutions including 7×7 , 15×15 , and 31×31 , and a fusion function. The MLKU will be repeated L times in the MLK module.

3.4. Efficient Class Relation module

We argue that ELISS is not simply allocating a category to each pixel in the image, but pursuing high-level semantics. It is needed to take into account the relationship among pixels, regions, and categories, the combination of which forms the spatial context. For one target pixel, we think that all relevant pixels in the ELI can contribute to it holistically. Under this hypothesis, we propose the efficient class relation (ECR) module with an efficient attention unit. See Fig. 5 and Algorithm 1 for the detailed process. In general, we use the entire features of the ELI to enhance the features on the hard region. Besides, it is worth noting that the use of efficient attention can save about 99.99% of memory with features on the FBP dataset ($n \times 256 \times 1700 \times 1800$).

3.5. Boundary-Aware Enhancement module

The boundary region is quite important for SS. Typically, SegFix [40] and PointRend [18] propose effective models to refine the segmentation results on the boundary region. However, these methods require too much computation and memory, especially for ELIs. For us, it is important but not necessary in our proposed ElegantSeg. Thus, we use the simplest way to perform boundary-aware enhancement (BAE). We naively use the Sobel operator to extract the boundaries among object regions from the mask. Then boundaries are used as the loss weights to punish the model on the boundary region, for a better learning of pixel relationship on the object boundary. This BAE module is simple but effective, which is especially useful in holistic learning to obtain complete boundaries for better segmentation. See Algorithm 1 for the detailed process.

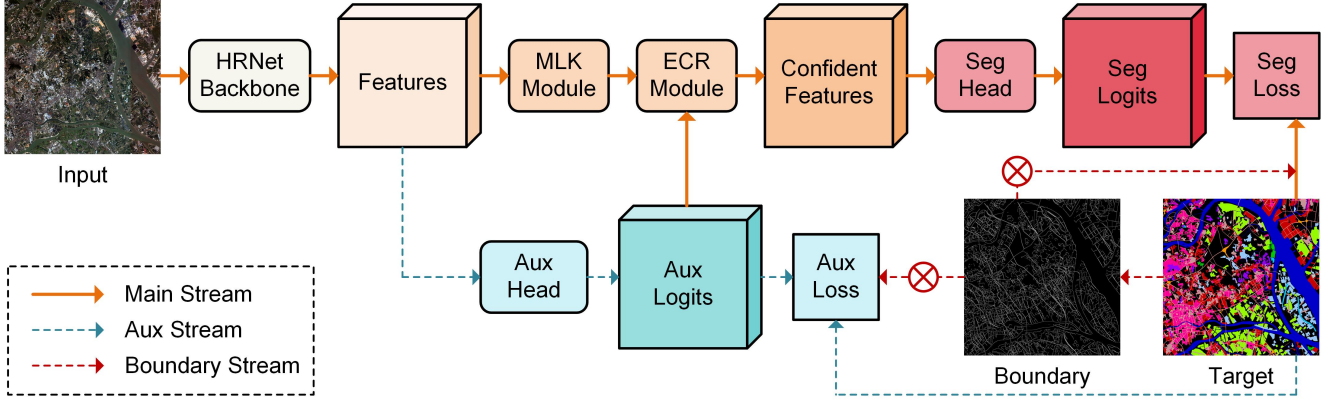


Figure 3. **Overview of the proposed ElegantSeg.** It consists of a HRNet backbone model, a MLK module, a ECR module, a BAE module, an auxiliary segmentation head and a main segmentation head.

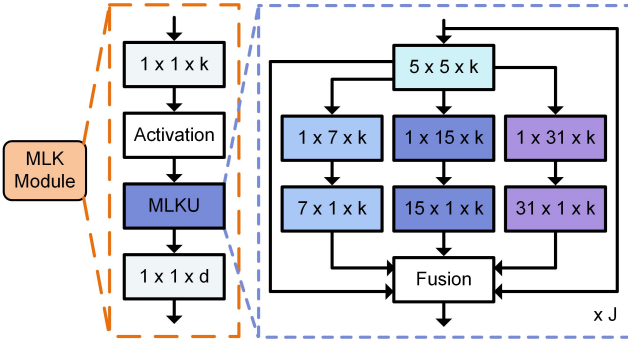


Figure 4. **Details of the proposed multiple large kernel module.** It can build long-range dependency and enlarge the effective receptive field of the ElegantSeg.

3.6. Loss function

We use auxiliary segmentation loss for deep supervision [19, 43] to better train the ElegantSeg model. Cross entropy loss is used for both the main segmentation loss and the auxiliary segmentation loss. The weights of losses are set to 1 and 0.4, respectively. See Algorithm 1.

4. Experiment

4.1. Datasets

DeepGlobe The DeepGlobe dataset is the most popular satellite image dataset in the computer vision domain. It contains 803 ELIs with 6 categories (2448×2448 pixels). We randomly split images into training, validation, and testing sets with 454, 142, and 207 images respectively, following the protocols of [3, 15, 17, 20].

FBP The FBP dataset contains more than 5 billion labeled pixels of 150 ELIs (6800×7200 pixels), annotated in a 24-category system covering artificial-constructed, agricultural, and natural classes. We follow the same testing sets

Algorithm 1 Pseudocode of ElegantSeg in PyTorch style.

```

# a: percentage for hard region selection
# b: kernel size of Sobel operator
# c: number of classes
# d: dimension of features
# n: mini-batch size
# h, w: height and width of features
# norm: softmax (or scaling)
# w_bd: weight for boundary enhancement loss
# w_aux: weight for auxiliary loss

# load a minibatch of images and labels
for images, labels in loader:
    feats = backbone.forward(images)
    aux_logits = aux_head.forward(feats)

# Multiple Large Kernel Module
feats = mlk_module.forward(feats)

# Efficient Class Relation Module
## Hard Region Selection
sort_logits = aux_logits.view(n, c, h*w)
feats = feats.view(n, d, h*w)
### compute the difference of logits
top1_logits = topk(aux_logits, dim=1, k=2)[:, 0, :]
top2_logits = topk(aux_logits, dim=1, k=2)[:, 1, :]
dif_logits = top1_logits - top2_logits
_, sorted_index = dif_logits.sort(descending)
### select a% of hard pixels
hard_region = sorted_index[:, 0:h*w*a]
### hard_feats: nxdx(h*w*a)
hard_feats = feats.gather(hard_region)
## Efficient Attention
query = norm(hard_feats, dim=1)
key = norm(feats, dim=2).view(n, d, h*w)
value = feats.view(n, h*w, d)
### holistic relation vector: nxdxd
relation = bmm(key, value)
att_value = bmm(query, relation)
### c_hard_feats: nxdx(h*w*a)
c_hard_feats = hard_feats + att_value
c_feats = feats.scatter(hard_region, c_hard_feats)
c_feats = c_feats.view(n, d, h*w)
seg_logits = seg_head.forward(c_feats)

# Boundary-Aware Enhancement Module
bd = Sobel(labels, ksize=b)
weight_bd = one_like(labels)+w_bd*binary(bd)
aux_loss = CELoss(aux_logits, labels, weight_bd)
seg_loss = CELoss(seg_logits, labels, weight_bd)
loss = seg_loss + w_aux*aux_loss
loss.backward()

```

bmm: batch matrix multiplication.

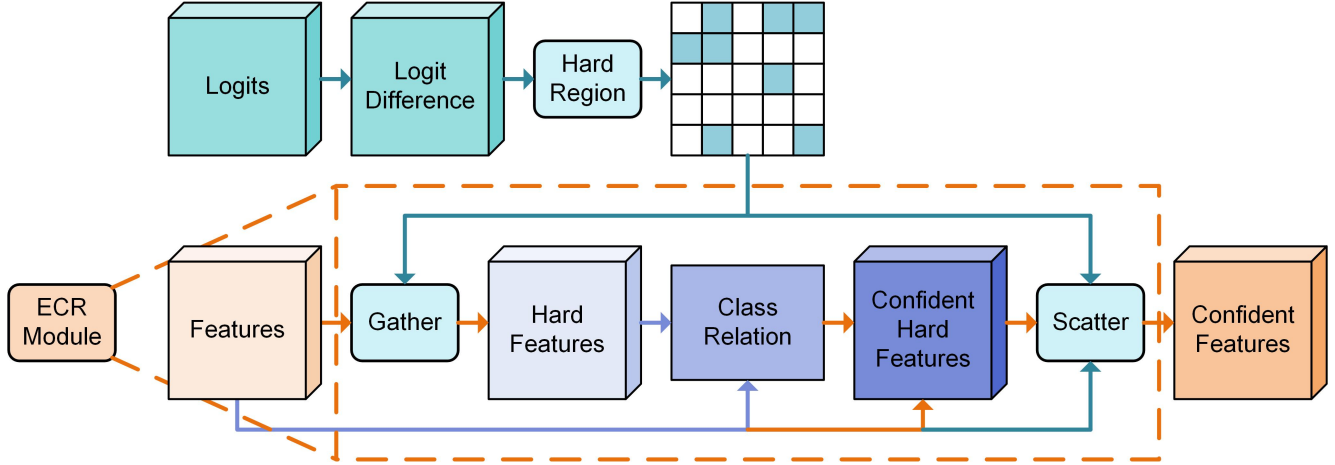


Figure 5. **Details of the proposed efficient class relation module.** It can help to develop holistic contextual relationships between regions.

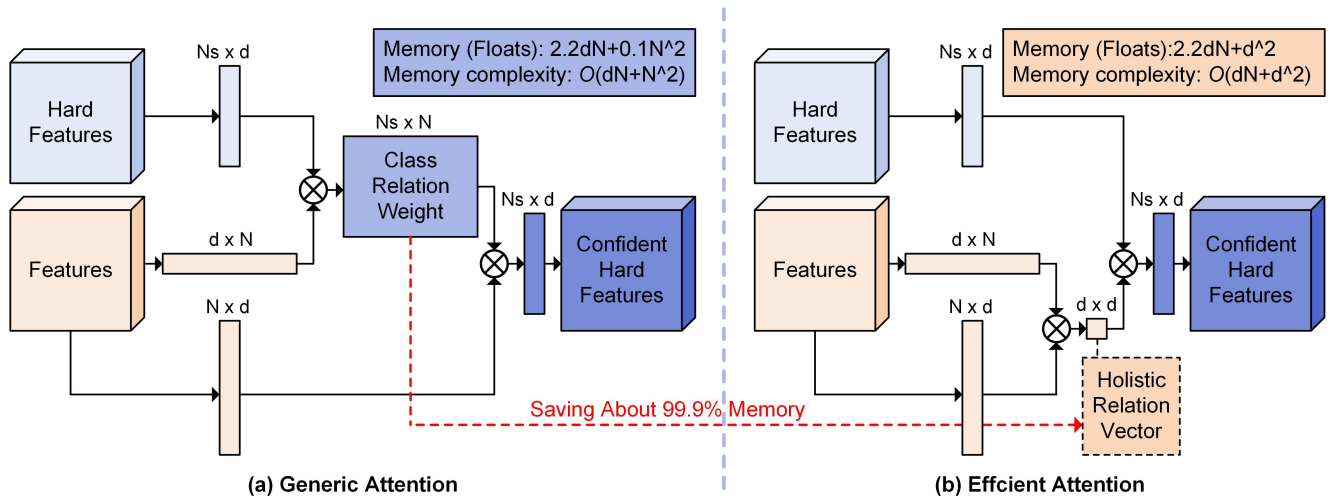


Figure 6. **Comparison between generic attention and efficient attention.** The two mechanisms are approximately equivalent according to the proof in [30]. The efficient attention can save about 99.99% memory on the FBP dataset.

with [32] and randomly split the rest images into training, validation with 90 and 30 images respectively.

4.2. Implementation details

We implement our framework using Pytorch on a server with three NVIDIA RTX Titan GPU, each of which has 24G GPU memory. The server has 1024G host memory to utilize the LMS toolbox. During training our local segmentation model, we adopt the SGD optimizer and a mini-batch size of 3 on the FBP dataset and 9 on the DeepGlobe dataset. The initial learning rate is set to 0.01 and it is decayed to 0.0001 by a step-based learning rate policy. In practice, it takes 10,000 iterations to converge our ElegantSeg model. For the first 1,500 iterations, the learning rate is set to 5e-5 for warming up. Following [18], we multiply the learning rate for the segmentation heads by 10 for better conver-

gence.

For the parameters, the number L for MLKU is set to 1. The percentage for hard region selection a is set to 10%. The kernel size of Sobel operator b is set to 5. The dimension of feature d is set to 256. The weights of the main segmentation loss and the auxiliary segmentation loss are 1 and 0.4, respectively. We use the mean intersection over union (mIoU) as the evaluation metric.

The implementations of existing generic SS methods are all based on the MMsegmentation [7] framework. All models are fine-tuned on the pre-trained models supported by MMsegmentation. As can be seen in Tab. 1, we select ten typical generic SS methods, and two popular transformer-based methods. All generic SS methods are trained for 10,000 iterations with the recommended hyper-parameters on the MMsegmentation [7] framework for a fair compari-

Methods	Publication	Backbone	End-to-End learning	Size of input	mIoU (%)	Size of input	mIoU (%)
<i>Generic semantic segmentation methods</i>				DeepGlobe [9]		FBP [32]	
PSPNet [43]	CVPR' 17	ResNet-101			71.59		53.33
DeepLabV3Plus [2]	ECCV' 18	ResNet-101			72.43		48.97
PSANet [44]	ECCV' 18	ResNet-101			72.71		54.22
UperNet [37]	ECCV' 18	ResNet-101			72.20		51.08
HRNet [34]	CVPR' 19	HRNet-W18			72.34		53.16
PointRend [18]	CVPR' 20	ResNet-101	✓	512×512	73.06	512×512	53.83
OCRNet [39]	ECCV' 20	HRNet-W48			73.04		55.71
STDC [14]	CVPR' 21	STDC			73.19		51.29
SwinTransformer [25]	ICCV' 21	UperNet-Swin-B			71.91		50.49
SegFormer [38]	NIPS' 21	MiT-B5			72.34		55.29
<i>Extra-large image semantic segmentation methods</i>							
GLNet [3]	CVPR' 19	FPN-ResNet-50		508×508	71.60	508×508	42.05
CascadePSP [4]	CVPR' 20	PSPNet-ResNet-50		512×512	68.50	512×512	-
PPN [36]	AAAI' 20	FPN-ResNet-50		512×512	71.90	512×512	-
MagNet [17]	CVPR' 21	FPN-ResNet-50	-	508×508	72.96	508×508	44.20
FCtL [20]	ICCV' 21	VGG-16		508×508	72.76	508×508	48.28
ISDNet [15]	CVPR' 22	DeepLabV3-ResNet-18		1224×1224	73.30	1224×1224	21.98
ElegantSeg-S		HRNet-W18		2448×2448*	73.95	6800×7200*	56.38
ElegantSeg-M	-	HRNet-W32	✓	2448×2448*	74.10	5000×5000	58.84
ElegantSeg-L		HRNet-W48		2448×2448*	74.32	4000×4000	61.62

* denotes the size for holistic learning.

Table 1. **Comparison with previous SOTAs on the DeepGlobe and the FBP dataset.** We evaluate all the generic SS methods on both datasets and some of the ELISS methods on the FBP dataset. The results of ELISS methods on the DeepGlobe dataset are collected from [15]. The codes released by [3, 15, 17, 20] are modified to adjust to the FBP dataset. We have tried our best to reproduce the good performance of these methods, but one of them (*i.e.*, [15]) fails to perform well in our re-implemented version.

	Holistic learning	BAE	MLK	ECR	mIoU (%)
ElegantSeg	-	-	-	-	72.40
		✓	-	-	72.25
		-	✓	-	72.02
		-	-	✓	72.13
		✓	✓	-	72.32
		-	✓	✓	71.98
		✓	✓	✓	72.02
ElegantSeg	✓	-	-	-	73.42
		✓	-	-	73.55
		✓	✓	-	73.71
		✓	✓	✓	73.95

Table 2. **Ablation study of proposed modules on the DeepGlobe dataset.** It shows that the specific effectiveness of the three proposed modules on the setting of holistic learning.

son.

For ELISS methods, we collect the results on the DeepGlobe dataset from [15] and re-implement four of them on the FBP dataset except for CascadePSP [4] and PPN [36] since CascadePSP [4] need a pre-trained model to predict initial segmentation for refinement and the scheme of PPN [36] is quite similar to GLNet [3] and FCtL [20]. As for GLNet [3] and FCtL [20], they cost too much time and memory with more than a week for training one stage, and about a month for the whole training process. The original version of [17] includes four-stage refinement. But we can only re-implement one-stage refinement according to their released code. Besides, the model of ISDNet [15] is specially designed for 3-band images, and we adjust it to fit the 4-band images on the FBP dataset. However, the result of ISDNet [15] is not promising.

4.3. Comparison with state-of-the-arts

The size of the DeepGlobe dataset is not that large and the annotations are coarse, which makes the ELISS much easier. While The FBP dataset is much harder compared to

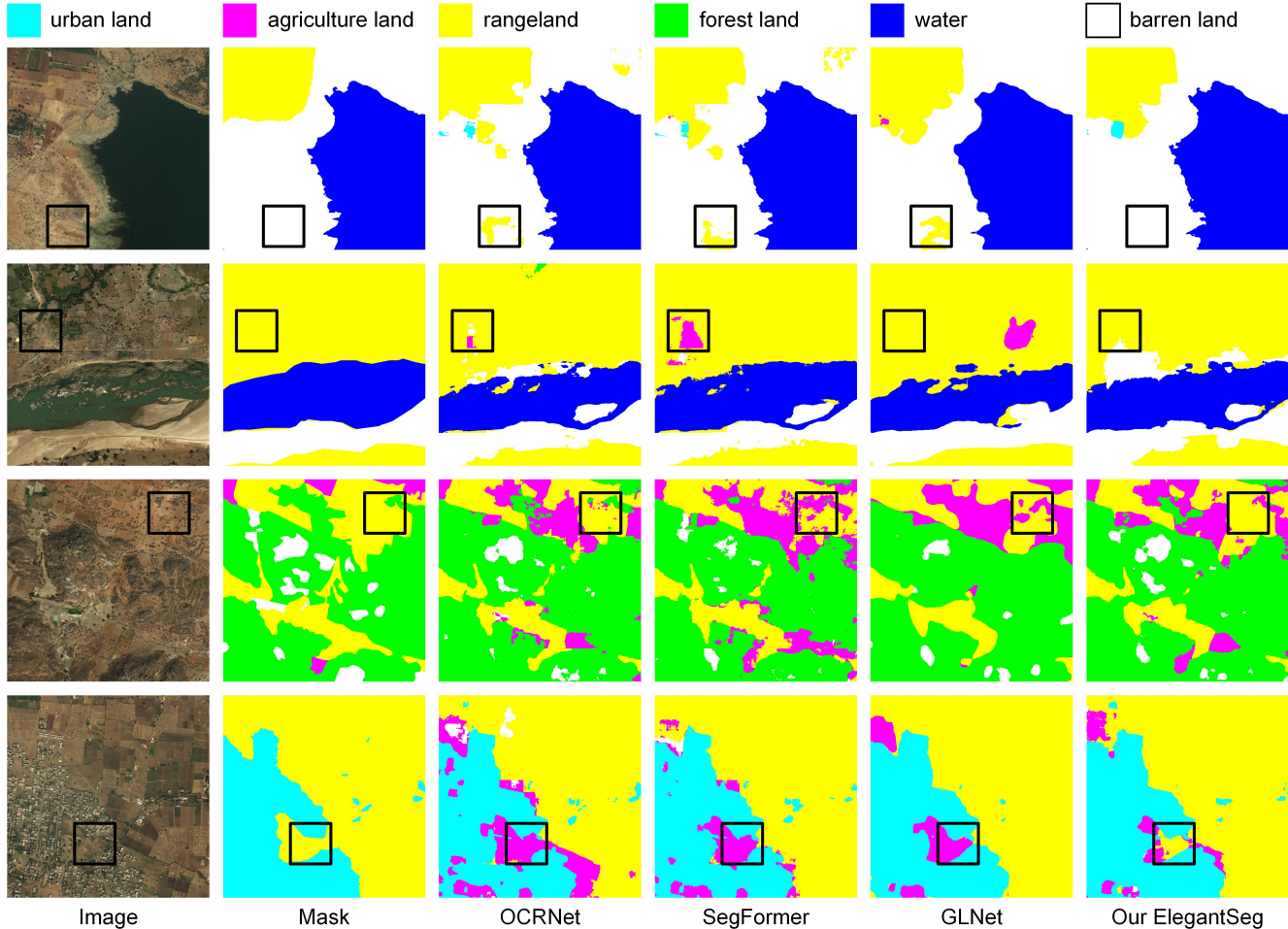


Figure 7. We illustrate some representative examples of the DeepGlobe dataset. Comparing with previous SOTAs, Our ElegantSeg shows superiority in obtaining complete object contours, as demonstrated in the black boxes.

the DeepGlobe dataset, containing 24 classes. The performances of [44], [39], and [38] on the FBP dataset are good. Based on the results, one can see that attention is important for extracting the spatial relationship in the FBP dataset.

Furthermore, existing methods fail to work well on the FBP dataset, while our ElegantSeg performs well on both datasets. Some qualitative results are presented in Fig. 7 and Fig. 8, which also demonstrate that our proposed ElegantSeg is superior to the three typical methods. These baseline method fails to perform well on some locations in Fig. 7 and Fig. 8. However, our proposed ElegantSeg can utilize complete information with holistic learning, combining the advantage of global and local views for better segmentation results.

4.4. Ablation study

We ablate the proposed three modules to see whether they are really effective for ELISS. To our surprise, the results turn out to be very interesting. The upper part of

Tab. 2 shows that when ElegantSeg is trained with sliced input (*i.e.*, without holistic learning), the proposed three modules including BAE, MLK, and ECR may be harmful to the performance.

However, they do help to increase the accuracy when trained with whole ELIs (*i.e.*, with holistic learning). We argue that the cropped slice is too small to provide enough spatial context. Thus, these three modules designed for ELIs bring too much noise for local pixels, which impairs performance.

5. Conclusion and Limitation

5.1. Conclusion

This paper presents a new scheme named **ElegantSeg** which can capably process holistic ELISS and overcome the GPU memory exhaustion caused by the extremely large sizes of ELIs. Different from existing global-local fusion methods or multi-stage refinement methods, ElegantSeg

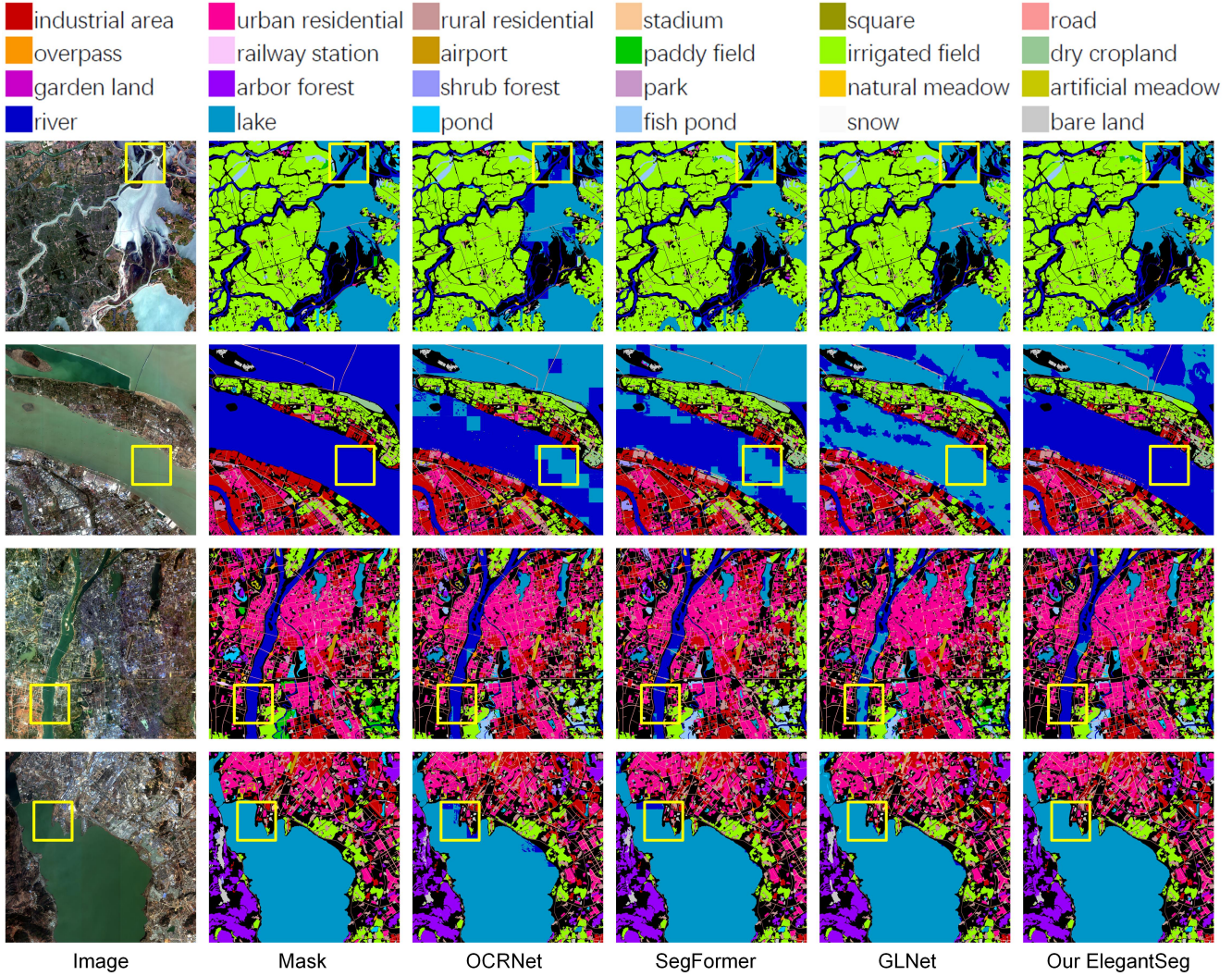


Figure 8. We illustrate some representative examples of the FBP dataset. Comparing with previous SOTAs, Our ElegantSeg shows superiority in obtaining complete object contours, as demonstrated in the yellow boxes.

can elegantly process holistic ELISS, and is able to thoroughly exploit the abundant information in ELIs. By extending the tensor storage from GPU memory to host memory, ElegantSeg breaks through the constraint of limited GPU memory. To the best of our knowledge, it is the first time that holistic ELISS can be done. Besides, three modules are specifically designed to utilize the characteristics of ELIs, including the multiple large kernel module, the efficient class relation module, and the boundary-aware enhancement module. Extensive experiments show that ElegantSeg outperforms the previous SOTAs on two typical ELISS datasets. We hope that ElegantSeg can open a new perspective for ELISS.

5.2. Limitation

Optimization ElegantSeg utilize the SGD optimizer which is designed for generic SS tasks. Nevertheless, the Mini-batch size is quite small while the image size is extremely large under the setting of ELISS. We believe that new optimization methods need to be explored for ELISS.

Efficiency For now, the efficiency of ElegantSeg is not fast enough. Some efficient techniques for accelerating the process of ELISS are supposed to be further developed, which are not limited to the attention function.

Generalization In this paper, the experiments are conducted on two ELI datasets with satellite images. It is required to conduct more experiments on ELI datasets with natural or medical images to verify the generalization of our proposed ElegantSeg.

References

- [1] Javiera Castillo-Navarro, Bertrand Le Saux, Alexandre Boulch, Nicolas Audebert, and Sébastien Lefèvre. Semi-supervised semantic segmentation in earth observation: The minifrance suite, dataset analysis and multi-task network study. *Machine Learning*, pages 1–36, 2021. [2](#)
- [2] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, pages 801–818, 2018. [1](#), [2](#), [3](#), [6](#)
- [3] Wuyang Chen, Ziyu Jiang, Zhangyang Wang, Kexin Cui, and Xiaoning Qian. Collaborative global-local networks for memory-efficient segmentation of ultra-high resolution images. In *CVPR*, pages 8924–8933, 2019. [1](#), [2](#), [4](#), [6](#)
- [4] Ho Kei Cheng, Jihoon Chung, Yu-Wing Tai, and Chi-Keung Tang. Cascadepsp: toward class-agnostic and very high-resolution segmentation via global and local refinement. In *CVPR*, pages 8890–8899, 2020. [1](#), [2](#), [6](#)
- [5] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. [3](#)
- [6] Noel CF Codella, David Gutman, M Emre Celebi, Brian Helba, Michael A Marchetti, Stephen W Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, et al. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic). In *2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018)*, pages 168–172, 2018. [1](#)
- [7] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. <https://github.com/open-mmlab/mms Segmentation>, 2020. [5](#)
- [8] Mihai Datcu and Klaus Seidel. Human-centered concepts for exploration and understanding of earth observation images. *IEEE Transactions on Geoscience and Remote Sensing*, 43(3):601–609, 2005. [2](#)
- [9] Ilke Demir, Krzysztof Koperski, David Lindenbaum, Guan Pang, Jing Huang, Saikat Basu, Forest Hughes, Devis Tuia, and Ramesh Raskar. Deepglobe 2018: A challenge to parse the earth through satellite images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 172–181, 2018. [1](#), [2](#), [6](#)
- [10] Lei Ding, Dong Lin, Shaofu Lin, Jing Zhang, Xiaojie Cui, Yuebin Wang, Hao Tang, and Lorenzo Bruzzone. Looking outside the window: Wide-context transformer for the semantic segmentation of high-resolution remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–13, 2022. [2](#)
- [11] Lei Ding, Jing Zhang, and Lorenzo Bruzzone. Semantic segmentation of large-size vhr remote sensing images using a two-stage multiscale training architecture. *IEEE Transactions on Geoscience and Remote Sensing*, 58(8):5367–5376, 2020. [2](#)
- [12] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11963–11975, 2022. [3](#)
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. [2](#)
- [14] Mingyuan Fan, Shenqi Lai, Junshi Huang, Xiaoming Wei, Zhenhua Chai, Junfeng Luo, and Xiaolin Wei. Rethinking bisenet for real-time semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9716–9725, 2021.
- [15] Shaohua Guo, Liang Liu, Zhenye Gan, Yabiao Wang, Wuhao Zhang, Chengjie Wang, Guannan Jiang, Wei Zhang, Ran Yi, Lizhuang Ma, et al. Isdnet: Integrating shallow and deep networks for efficient ultra-high resolution segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4361–4370, 2022. [1](#), [2](#), [4](#), [6](#)
- [16] Bart M Ter Haar Romeny and Luc Florack. A multiscale geometric model of human vision. In *The Perception of Visual Information*, pages 73–114. Springer, 1993. [2](#)
- [17] Chuong Huynh, Anh Tuan Tran, Khoa Luu, and Minh Hoai. Progressive semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16755–16764, 2021. [1](#), [2](#), [4](#), [6](#)
- [18] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. Pointrend: Image segmentation as rendering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9799–9808, 2020. [1](#), [2](#), [3](#), [5](#)
- [19] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial intelligence and statistics*, pages 562–570. PMLR, 2015. [4](#)
- [20] Qi Li, Weixiang Yang, Wenxi Liu, Yuanlong Yu, and Shengfeng He. From contexts to locality: Ultra-high resolution image segmentation via locality-aware contextual correlation. In *ICCV*, pages 7252–7261, 2021. [1](#), [2](#), [4](#), [6](#)
- [21] Xiaoxiao Li, Ziwei Liu, Ping Luo, Chen Change Loy, and Xiaoou Tang. Not all pixels are equal: Difficulty-aware semantic segmentation via deep layer cascade. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3193–3202, 2017. [2](#)
- [22] Yansheng Li, Wei Chen, Xin Huang, Zhi Gao, Siwei Li, Tao He, and Zhang Yongjun. Mfvnet: Deep adaptive fusion network with multiple field-of-views for remote sensing image semantic segmentation. *SCIENCE CHINA Information Sciences*, 2022. [1](#), [2](#)
- [23] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1925–1934, 2017. [2](#)
- [24] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al.

- Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12009–12019, 2022. [1](#), [2](#)
- [25] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10012–10022, October 2021. [1](#), [2](#), [6](#)
- [26] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015. [1](#), [2](#)
- [27] Chao Peng, Xiangyu Zhang, Gang Yu, Guiming Luo, and Jian Sun. Large kernel matters—improve semantic segmentation by global convolutional network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4353–4361, 2017. [3](#)
- [28] Alexandre Robicquet, Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Learning social etiquette: Human trajectory understanding in crowded scenes. In *European conference on computer vision*, pages 549–565. Springer, 2016. [1](#)
- [29] Tiancheng Shen, Yuechen Zhang, Lu Qi, Jason Kuen, Xingyu Xie, Jianlong Wu, Zhe Lin, and Jiaya Jia. High quality segmentation for ultra high-resolution images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1310–1319, 2022. [2](#)
- [30] Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Efficient attention: Attention with linear complexities. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 3531–3539, 2021. [3](#), [5](#)
- [31] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7262–7272, 2021. [1](#)
- [32] Xin-Yi Tong, Gui-Song Xia, Qikai Lu, Huanfeng Shen, Shengyang Li, Shucheng You, and Liangpei Zhang. Land-cover classification with high-resolution remote sensing images using transferable deep models. *Remote Sensing of Environment*, 237:111322, 2020. [1](#), [2](#), [5](#), [6](#)
- [33] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9, 2018. [1](#)
- [34] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *PAMI*, 2020. [1](#), [2](#), [3](#)
- [35] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. [3](#)
- [36] Tong Wu, Zhenzhen Lei, Bingqian Lin, Cuihua Li, Yanyun Qu, and Yuan Xie. Patch proposal network for fast semantic segmentation of high-resolution images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12402–12409, 2020. [2](#), [6](#)
- [37] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European conference on computer vision (ECCV)*, pages 418–434, 2018. [1](#)
- [38] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M. Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 12077–12090. Curran Associates, Inc., 2021. [1](#), [2](#), [7](#)
- [39] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-contextual representations for semantic segmentation. In *European conference on computer vision*, pages 173–190. Springer, 2020. [1](#), [3](#), [7](#)
- [40] Yuhui Yuan, Jingyi Xie, Xilin Chen, and Jingdong Wang. Segfix: Model-agnostic boundary refinement for segmentation. In *European Conference on Computer Vision*, pages 489–506. Springer, 2020. [2](#), [3](#)
- [41] Ting Zhang, Guo-Jun Qi, Bin Xiao, and Jingdong Wang. Interleaved group convolutions. In *Proceedings of the IEEE international conference on computer vision*, pages 4373–4382, 2017. [3](#)
- [42] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10076–10085, 2020. [3](#)
- [43] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *CVPR*, pages 2881–2890, 2017. [1](#), [2](#), [4](#), [6](#)
- [44] Hengshuang Zhao, Yi Zhang, Shu Liu, Jianping Shi, Chen Change Loy, Dahua Lin, and Jiaya Jia. Psanet: Pointwise spatial attention network for scene parsing. In *Proceedings of the European conference on computer vision (ECCV)*, pages 267–283, 2018. [1](#), [7](#)
- [45] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xi Tian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *CVPR*, pages 6881–6890, 2021. [1](#), [2](#)
- [46] Zhuo Zheng, Yanfei Zhong, Junjie Wang, and Ailong Ma. Foreground-aware relation network for geospatial object segmentation in high spatial resolution remote sensing imagery. In *CVPR*, pages 4096–4105, 2020. [1](#)