

# Boundary-Aware Semantic Segmentation of Remote Sensing Images via Segformer and Snake Convolution

Xia Yanting, Zhang Lin, Guo Ting, Jin Qi

Geely University of China, China  
zhlin002@163.com

**Abstract.** Semantic segmentation of remote sensing images remains challenging due to complex object structures and varying scales. This paper proposes a novel hybrid segmentation model that combines Segformer for global context extraction with Dynamic Snake Convolution to better capture fine-grained, boundary-aware features. An auxiliary semantic branch is introduced to improve feature alignment across scales. Experiments on three benchmark datasets—LoveDA, Potsdam, and Vaihingen—demonstrate that the proposed approach achieves consistent improvements in mIoU over baseline models, particularly in segmenting irregular and linear structures. This framework offers a promising solution for high-resolution land cover mapping and urban scene understanding.

**Keywords:** Segformer, dynamic snake convolution, remote sensing, Semantic Segmentation, Deep learning.

## 1. Introduction

Semantic segmentation of remote sensing images faces unique challenges due to the complex interplay of large-scale variations, mixed textures, and ambiguous boundaries. High-resolution imagery often contains slender structures (e.g., roads, rivers) and irregular objects (e.g., fragmented buildings) that are poorly captured by conventional convolutional neural network (CNN) or vision transformers (ViT). Existing methods struggle to balance global context modeling with local geometric adaptability, leading to boundary blurring and misclassification of fine-grained features. Semantic segmentation of remote sensing images is crucial for applications such as land cover mapping, environmental monitoring, and urban planning. However, the high resolution and complexity of these images [37, 19, 36], combined with the difficulty of annotating large-scale datasets, pose significant challenges for existing segmentation models. Accurately capturing intricate structures, such as roads, rivers [33], and urban boundaries [32], remains an unresolved problem in deep learning-based segmentation. These challenges necessitate more adaptive architectures.

With the development of deep learning, CNN-based [15] methods have become widely used for semantic segmentation. Fully Convolutional Networks (FCN) [11] replaced fully connected layers with convolutional layers, achieving end-to-end segmentation. U-Net [18, 23] introduced a symmetric encoder-decoder structure with skip connections, enhancing multi-scale feature extraction, while SegNet [26] utilized unpooling layers to improve spatial detail recovery. Later, PSPNet [8] incorporated spatial pyramid pooling for better

global context understanding, and DeepLab [14] leveraged Atrous Spatial Pyramid Pooling (ASPP)[16] to integrate multi-scale features. Despite these advancements, these methods still face challenges in modeling complex interactions between pixels, often resulting in information loss. Traditional CNNs, constrained by local receptive fields, struggle to capture long-range dependencies and segment slender, irregular structures commonly found in remote sensing images. These limitations highlight the need for more adaptive and efficient segmentation approaches.

With the emergence of Transformers, semantic segmentation has further advanced. Liu et al. [35] proposed the Swin Transformer, which adopts a hierarchical structure and local attention mechanism, enabling the model to achieve higher computational efficiency while maintaining high accuracy. Zheng et al. [24] applied Transformers to semantic segmentation tasks, serializing images and feeding them into Transformers to use self-attention mechanisms at each layer for global information, thereby improving segmentation accuracy. Zheng et al. [25] proposed SETR (Semantic Segmentation Transformer Network), inspired by the ViT model, which enhances the semantic representation and generalization capabilities by introducing pixel alignment mechanisms and multi-scale attention fusion methods. However, the high complexity of Transformers results in relatively slower training speeds.

Single CNN or ViT [5] models struggle to balance local and global feature representation effectively. To address this, Zhang et al. [6] proposed a lightweight dual-branch neural network to solve intra-class heterogeneity and inter-class homogeneity problems. Jiang et al. [22] designed cross-residual feature blocks and improved skip connections to achieve dual-branch multi-scale channel cross-fusion. He et al. [30] embedded Transformers into U-Net, constructing spatial interaction modules and feature compression modules to mitigate the loss of detailed features. Wang et al. [29] proposed an algorithm based on an enhanced diffusion model. By incorporating scalable jump-connection layers into the denoising probability diffusion model, the approach effectively handles multi-scale features in campus environments, achieving superior accuracy in image semantic segmentation for autonomous driving across diverse settings. Weng et al. [13] designed the Sgformer network, incorporating multi-level feature attention to integrate the spatial details of CNNs and the contextual semantics of ViTs. Geng et al. [10] proposed DPFANet, which constructs edge optimization blocks to constrain edge features and effectively model images from local to global features. Despite advancements in CNN- and Transformer-based models, existing approaches struggle to effectively capture the elongated and irregular structures common in remote sensing images, such as roads, rivers, and building outlines. These limitations arise due to inadequate local feature representation and inefficient shape adaptation in conventional convolutional layers. While hybrid CNN-Transformer architectures improve multi-scale feature fusion, they lack specialized mechanisms to handle elongated or tortuous structures. For instance, deformable convolutions adaptively adjust receptive fields but may diverge from target boundaries in linear features. Similarly, attention mechanisms enhance global dependencies but neglect local geometric priors. To address these gaps, we propose a model that integrates Segformer's hierarchical attention with Dynamic Snake Convolution, which explicitly constrains deformable offsets to follow linear structures.

- Utilizing dynamic snake convolution to adaptively focus on slender, tortuous local structures and complex, variable global shapes, accurately capturing the tubular features in remote sensing data.
- Incorporating an auxiliary semantic branch to extract contextual information from images, ensuring the extraction of rich semantic features while maintaining inference efficiency.
- Conducting experimental analyses on publicly available datasets LoveDA and Potsdam. The experimental results demonstrate that the proposed model achieves superior segmentation performance, with mIoU reaching 52.49% on the LoveDA dataset, 79.71% on the Potsdam dataset and 76.70% on the Vaihingen dataset, representing improvements of 4.18%, 0.74% and 2.09%, respectively, over the baseline model U-Net.

The remainder of this paper is structured as follows. Section 2 presents the proposed methodology, detailing the integration of Segformer and Dynamic Snake Convolution. Section 3 describes the datasets and experimental setup, while Section 4 discusses the results and comparative analysis with baseline models. Finally, Conclusion concludes the study with key findings and future research directions.

## 2. Related Work

### 2.1. Semantic Segmentation

Semantic segmentation, a fundamental task in computer vision, aims to assign pixel-level labels to images. Early approaches relied on handcrafted features and traditional machine learning methods, but the advent of deep learning revolutionized the field. Long et al. [11] proposed Fully Convolutional Networks (FCN), replacing dense layers with convolutional layers to enable end-to-end segmentation. Building on this, U-Net [18] introduced an encoder-decoder architecture with skip connections, enhancing feature fusion across scales. Subsequent works, such as SegNet [26], improved spatial resolution recovery through unpooling layers, while PSPNet [8] and DeepLab [14] leveraged pyramid pooling and atrous convolutions to capture multi-scale context.

Despite these advancements, CNN-based methods struggled with long-range dependencies and irregular structures due to their localized receptive fields. The emergence of Vision Transformers (ViTs) addressed this limitation by modeling global interactions via self-attention. Zheng et al. [25] proposed SETR, which treats segmentation as a sequence-to-sequence problem using pure Transformers. Swin Transformer [35] further enhanced efficiency by introducing hierarchical shifted windows, balancing local and global feature extraction. Hybrid architectures, such as Swin-UNet [22], combined Transformers with U-Net to preserve spatial details while capturing global context. However, challenges persist in segmenting slender, tortuous structures (e.g., roads, rivers) in remote sensing imagery, necessitating specialized geometric modeling techniques like Dynamic Snake Convolution [21].

Recent advancements in semantic segmentation have also focused on improving the efficiency and scalability of models. Lightweight architectures, such as those proposed by Zhang et al. [6], aim to reduce computational complexity while maintaining high accuracy. These models often employ techniques like depthwise separable convolutions and

channel attention to optimize performance. Additionally, self-supervised learning methods have gained traction, reducing the dependency on large annotated datasets by leveraging unlabeled data for pre-training [7].

Another significant development is the integration of multi-task learning, where models are trained to perform multiple related tasks simultaneously, such as segmentation and object detection. This approach has been shown to improve generalization and robustness, particularly in complex scenes with diverse objects and backgrounds [28]. Furthermore, the use of generative adversarial networks (GANs) for data augmentation has proven effective in enhancing model performance, especially in scenarios with limited labeled data [20].

## 2.2. Attention Mechanism

Attention mechanisms have become pivotal in enhancing segmentation models by dynamically focusing on salient regions. Early efforts integrated channel-wise attention, as seen in Squeeze-and-Excitation Networks (SENet) [27], to recalibrate feature responses. Later, Non-local Networks [16] introduced self-attention to model long-range dependencies, improving contextual understanding. Transformers [35,24,25] further popularized attention by replacing convolutional operations with multi-head self-attention layers, enabling global feature interactions.

In semantic segmentation, attention mechanisms are often applied hierarchically. For instance, DeepLabv3+ [14] combined atrous spatial pyramid pooling (ASPP) with attention to refine multi-scale features. Similarly, Swin Transformer [35] employed shifted window-based attention to reduce computational complexity while maintaining global modeling capabilities. Recent works, such as DPFANet [13], integrated edge-aware attention to enhance boundary detection in remote sensing images. These mechanisms address challenges like intra-class heterogeneity and inter-class homogeneity, particularly in complex scenes. Dynamic Snake Convolution [21], with its iterative attention to linear structures, exemplifies how geometric-prior-guided attention can improve segmentation of tubular features in remote sensing data.

Attention mechanisms have also been extended to incorporate spatial and temporal dimensions, particularly in video segmentation tasks. Spatial attention focuses on relevant regions within a single frame, while temporal attention captures dependencies across multiple frames. This dual attention approach has been shown to improve the segmentation of dynamic scenes, such as those encountered in video surveillance and autonomous driving [31].

Moreover, the integration of attention mechanisms with graph neural networks (GNNs) has opened new avenues for semantic segmentation. GNNs model relationships between pixels or regions as a graph, allowing for more flexible and context-aware feature extraction. When combined with attention mechanisms, GNNs can effectively capture both local and global dependencies, leading to improved segmentation accuracy in complex scenes [3].

## 2.3. Remote Sensing Image Segmentation

Remote sensing image segmentation presents unique challenges due to the high resolution, complex structures, and diverse land cover types. Traditional methods often rely on

handcrafted features and machine learning algorithms, which struggle to capture the intricate details and variability in remote sensing data. With the advent of deep learning, convolutional neural networks (CNNs) have become the dominant approach, offering significant improvements in accuracy and robustness.

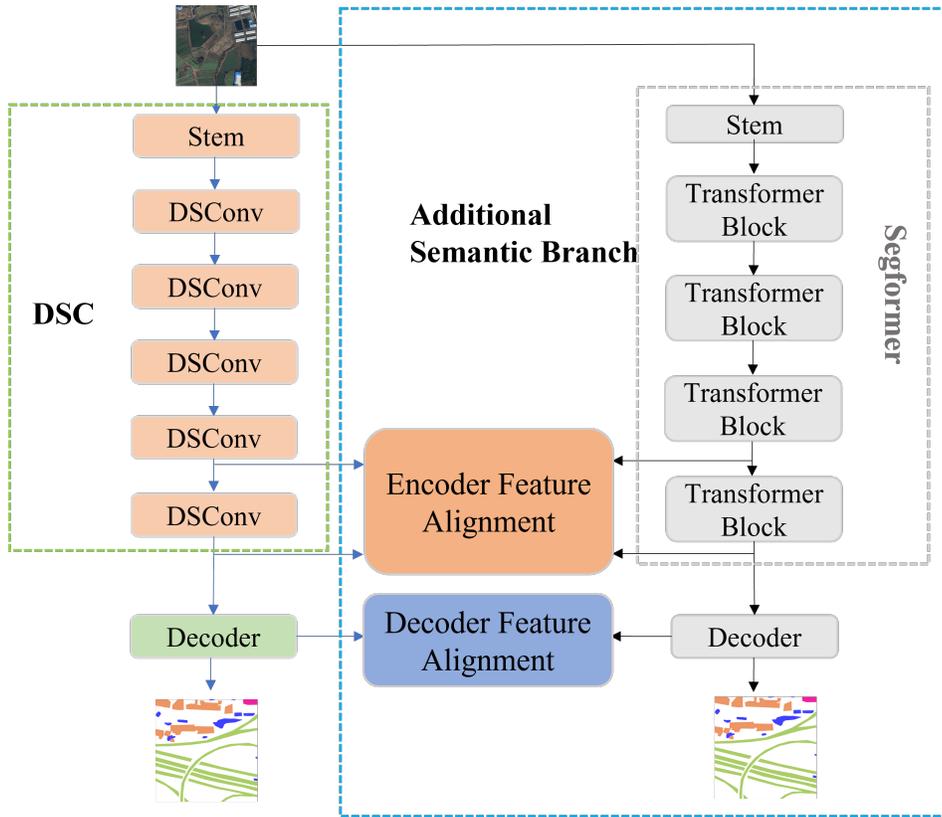
One of the key challenges in remote sensing image segmentation is the effective handling of multi-scale features. High-resolution images often contain objects of varying sizes, from small buildings to large agricultural fields. Multi-scale feature extraction techniques, such as pyramid pooling and atrous convolutions, have been widely adopted to address this issue. Additionally, the integration of attention mechanisms has proven effective in focusing on relevant regions and improving the segmentation of complex structures.

Another critical aspect is the ability to model long-range dependencies, which is essential for accurately segmenting large and irregular objects like rivers and roads. Vision Transformers (ViTs) have emerged as a powerful tool for capturing global context, leveraging self-attention mechanisms to model interactions between distant pixels. Hybrid models that combine CNNs and Transformers have shown promise in balancing local detail extraction with global context understanding.

Despite these advancements, segmenting slender and irregular structures remains a significant challenge. Traditional convolutional layers, with their fixed receptive fields, often fail to capture the geometric intricacies of such structures. Dynamic Snake Convolution offers a novel solution by adaptively focusing on linear and curved features, enhancing the segmentation of roads, rivers, and other elongated objects in remote sensing imagery.

### 3. Method

This paper proposes an efficient semantic segmentation method for remote sensing images by integrating Segformer with snake convolution. As illustrated in Figure 1, the proposed framework comprises three components: (1) a Segformer branch for multi-scale global context extraction, (2) a Dynamic Snake Convolution branch for boundary-aware feature refinement, and (3) an auxiliary semantic alignment module to harmonize cross-scale features. SegFormer is a simple, efficient, and powerful semantic segmentation framework. By combining a Transformer with a lightweight multi-layer perceptron decoder, SegFormer is capable of extracting high-resolution coarse features and low-resolution fine features, aggregating multi-scale information across different layers. Through a combination of local and global attention, it generates strong feature representations and extracts effective contextual information. The Dynamic Snake Convolution [10] enhances geometric structure perception by adaptively focusing on small and curved local features of tubular structures, thereby specifically improving the perception of such structures. To address the challenges posed by complex and variable global shapes, a multi-view feature fusion strategy is employed. The proposed method uses Dynamic Snake Convolution as the main branch structure, while SegFormer serves as an auxiliary branch for training. By employing a semantic alignment model to integrate the additional branch, the network extracts rich semantic information, achieving precise and efficient segmentation simultaneously.

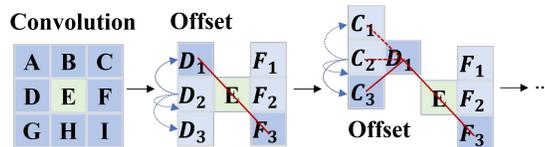


**Fig. 1.** The Image Segmentation Framework

**3.1. Dynamic Snake Convolution**

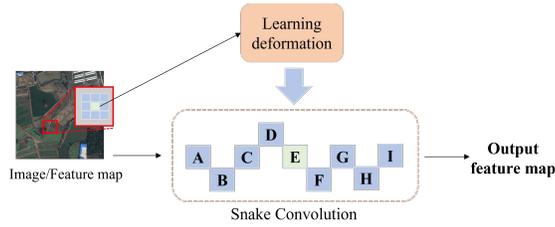
Given the standard 2D convolution coordinates  $K$ , with the center coordinate as  $K_i = (x_i, y_i)$ , a  $3 \times 3$  convolution kernel  $K$  can be represented as:

$$K = \{(x - 1, y - 1), (x - 1, y), \dots, (x + 1, y + 1)\} \tag{1}$$



**Fig. 2.** Learning deformation

To provide the convolution kernel with greater flexibility and enable it to focus on the complex geometric features of targets, deformable offsets  $\Delta$ [9] are introduced. However, if the model is given complete freedom to Learning deformable offsets (Figure 2), the receptive field often deviates from the target, particularly when handling slender tubular structures. Therefore, an iterative strategy is adopted (Figure 3), which sequentially selects the next position of the target to be processed for observation. This ensures continuity of focus and prevents the receptive field from spreading too far due to large deformable offsets.



**Fig. 3.** Dynamic Snake Convolution

In dynamic snake convolution, the standard convolution kernel is linearized along both the  $x$ -axis and  $y$ -axis. Considering a convolution kernel of size 9, take the  $x$ -axis direction as an example. The specific position of each grid in  $K$  is represented as:  $K_{i\pm c} = (x_{i\pm c}, y_{i\pm c})$ , where  $c = 0, 1, 2, 3, 4$  indicates the horizontal distance from the center grid. The selection of each grid position  $K_{i\pm c}$  in the convolution kernel is a cumulative process. Starting from the center position  $K_i$  the position of grids farther from the center depends on the position of the preceding grid:

$K_{i+1}$  is determined by adding an offset  $\Delta = \{\delta \mid \delta \in [-1, 1]\}$  relative to  $K_i$ . Therefore, the offsets are accumulated  $\sum$ , ensuring that the kernel conforms to a linear structural pattern. The changes along the  $x$ -axis direction are illustrated in Figure 3 as:

$$K_{i\pm c} = \begin{cases} (x_{i+c}, y_{i+c}) = \left( x_i + c, y_i + \sum_{i}^{i+c} \Delta y \right) \\ (x_{i-c}, y_{i-c}) = \left( x_i - c, y_i + \sum_{i-c}^i \Delta y \right) \end{cases} \quad (2)$$

The changes along the  $y$ -axis direction are:

$$K_{j\pm c} = \begin{cases} (x_{j+c}, y_{j+c}) = \left( x_j + \sum_j^{j+c} \Delta x, y_j + c \right) \\ (x_{j-c}, y_{j-c}) = \left( x_j + \sum_{j-c}^j \Delta x, y_j - c \right) \end{cases} \quad (3)$$

Since the offset  $\Delta$  is typically a decimal, while coordinates are usually in integer form, bilinear interpolation is adopted, represented as:

$$K = \sum_{K'} B(K', K) \cdot K' \quad (4)$$

Here,  $K$  represents the decimal position in Equations 2 and 3,  $K'$  enumerates all integer spatial positions, and  $B$  is the bilinear interpolation kernel, which can be decomposed into two one-dimensional kernels, as:

$$B(K, K') = b(K_x, K'_x) \cdot b(K_y, K'_y) \quad (5)$$

---

**Algorithm 1** Dynamic Snake Convolution
 

---

**Input:** Feature map  $F$ , initial kernel center  $K_i = (x_i, y_i)$ , kernel size  $S = 9$ .

**Process:**

- 1: **for** each offset step  $c$  along x/y-axis **do**
  - 2:   **a.** Learn deformable offsets  $\Delta x, \Delta y$  via a lightweight network.
  - 3:   **b.** Accumulate offsets iteratively:
  - 4:     x-direction:  $x_{i+c} = x_i + c, \quad y_{i+c} = y_i + \sum \Delta y$ .
  - 5:     y-direction:  $x_{j+c} = x_j + \sum \Delta x, \quad y_{j+c} = y_j + c$ .
  - 6:   **c.** Compute interpolated features  $F_{\text{interp}}$  using bilinear sampling.
  - 7: **end for**
  - 8: Aggregate features from all offset positions.
  - 9: **Output:** Refined feature map  $F_{\text{out}}$ .
- 

As shown in Algorithm 1, the Dynamic Snake Convolution algorithm is presented. The dynamic snake convolution kernel is designed to better adapt to slender tubular structures based on dynamic configurations, enabling enhanced perception of key features.

### 3.2. Integrating Method

As shown in Figure 3, we propose a simple yet effective alignment module for feature learning during training. It can be divided into encoder feature alignment and decoder feature alignment.

- Encoder Feature Alignment

Backbone feature alignment begins by downsampling or upsampling the features of the Transformer and CNN branches for alignment. To avoid direct feature alignment disrupting the supervision of the CNN by the ground truth during training, feature projection is employed. Specifically, the CNN features are projected to the dimension of the Transformer features. This projection unifies the number of channels and prevents direct feature alignment. Finally, semantic alignment loss is applied to the projected features to align the semantic representations [21].

- Decoder Feature Alignment

Features from stages 2 and 4 are selected for alignment. Considering the significant differences in the decoding space between the Transformer network and the backbone network, directly aligning decoding features and output logits only leads to limited improvement. Therefore, we adopt a shared decoder head alignment approach. Specifically, the features from stages 2 and 4 of the single-branch CNN are fed into a point convolution to expand their dimensions. The high-dimensional features are then passed through the Transformer decoder. The new output features and logits of the Transformer decoder are used to compute alignment loss with the original outputs of the Transformer decoder.

### 3.3. The Alignment Loss

To better align semantic information, an alignment loss [31] focusing on semantic information rather than spatial information is required. In this implementation, we use MGD Loss (channel-wise distillation loss) [34] as the alignment loss, which demonstrates better performance compared to other loss functions. MGD Loss consists of two components: a global distribution alignment term and a boundary constraint term.

$$\ell_{\text{align}} = \|E_{x \sim P[\phi(x)]} - E_{y \sim Q[\phi(v)]}\|^2 \quad (6)$$

**Global Distribution Alignment Term:** The alignment goal is achieved by measuring the difference between the feature distributions of the two modalities. In this paper, Maximum Mean Discrepancy (MMD) [3] is used as the global distribution alignment term, which is expressed as:

$$l_{\text{margin}} = \max(0, m + d(f_a(a), f_b(b^-)) - d(f_a(a), f_b(b^+))) \quad (7)$$

Here,  $m$  represents the margin value.  $f_a(a)$  and  $f_b(b^-)$  are the features of modalities  $a$  and  $b$ , respectively.  $b^+$  denotes positive samples,  $b^-$  denotes negative samples, and  $d(\bullet)$  is the distance metric (e.g., Euclidean distance).

The total MGD Loss combines the global alignment loss and margin-based loss:

$$\ell_{MGD} = \lambda \ell_{\text{align}} + (1 - \lambda) \ell_{\text{margin}} \quad (8)$$

$\lambda$  is a hyperparameter that balances the trade-off between alignment and discrimination. By jointly optimizing these components, MGD Loss achieves alignment of multi-modal data distributions and ensures semantic consistency.

## 4. Experiments and Results

### 4.1. Dataset

To verify the proposed remote sensing image segmentation method, we used three public HR remote sensing image [4] datasets, LoveDA [12], Potsdam and Vaihingen. The LoveDA dataset is designed to facilitate research on event detection tasks in remote sensing imagery, such as natural disaster monitoring, urban planning, etc. The LoveDA dataset provides high-resolution airborne remote sensing imagery covering a wide range of scenes

and environments. Each image is labeled with rich event categories, including natural disasters, traffic accidents, buildings, etc. In addition, a large number of remote sensing images of real scenes are also included, so that the model can be better generalized in real environments. The LoveDA dataset consists of 5987 high-resolution non-interlaced optical remote sensing images of *Nanjing*, *Changzhou*, and *Wuhan* with 166,768 labeled objects, and the size of each pair of images is  $1024 \times 1024$  pixels with a pixel separation rate of 0.3 meters, and all of the information was obtained from the Google Earth platform. Earth's platform. The dataset contains seven categories, including background, buildings, roads, water bodies, debris, forests, and agriculture, covering rural and urban areas, respectively. In these datasets, there are 2713 urban landscapes and 3274 rural landscapes. Potsdam Remote Sensing Dataset is a high-resolution airborne remote sensing image dataset provided by the University of Potsdam, Germany, which is designed to support research on remote sensing image processing tasks such as feature classification, target detection and semantic segmentation. Potsdam [2] contains 28 images of the same size, with a spatial resolution of 5 cm for the top image and DSM, and the size of each pair of images is  $6000 \times 6000$  pixels. The dataset contains 38 images, we split them into 26 training, 4 validation, and 8 test images. The dataset contains 6 categories of impervious surfaces, buildings, low vegetation, trees, cars and background. Each image is labeled with precise feature classes and bounding boxes including buildings, roads, trees, etc. In addition, the Potsdam dataset provides multispectral imagery in multiple bands (e.g., red, green, blue, and near-infrared), as well as high-resolution panchromatic imagery, providing researchers with a rich data resource. The Vaihingen dataset is another widely used benchmark for high-resolution remote sensing image segmentation, provided by the International Society for Photogrammetry and Remote Sensing (ISPRS). This dataset consists of 33 aerial orthoimagery tiles covering urban and suburban areas of Vaihingen, Germany. Each image has a spatial resolution of 9 cm and a size of  $2000 \times 2000$  pixels, captured in three spectral bands (near-infrared, red, and green). The dataset includes six semantic categories: impervious surfaces, buildings, low vegetation, trees, cars, and clutter/background. Additionally, it provides digital surface models (DSM) to enhance 3D feature analysis. The Vaihingen dataset is particularly challenging due to its fine-grained details, dense object distribution, and complex urban layouts, making it suitable for evaluating models' capability in handling intricate scenes. Following common practices, we utilize 16 images for training and 17 for testing, ensuring compatibility with existing research benchmarks[17].

#### 4.2. Experimental Parameter Settings

In this experiment, the software configurations are *Ubuntu18.04 LTS* operating systems, *Python3.8* development language, and *Pytorch* Deep learning framework; the hardware configurations are one *NVIDIA RTX 3090 GPU*. Besides, and training hyperparameters are summarized in Table1.

#### 4.3. Evaluation Metrics

To evaluate the performance of the algorithm, we employ Intersection over Union (IoU)[1], mean Intersection over Union (mIoU), F1-Score, and Overall Accuracy (OA) as evaluation metrics. IoU is defined as the ratio of the intersection to the union of the algorithm's predicted segmentation and the ground truth segmentation. mIoU is the average

**Table 1.** Experimental Parameter Settings

Parameters	Value
Batch Size	4
Initial Learning Rate	0.0001
Optimizer	AdamW
Iterations	500

IoU across all segmentation classes. The F1-Score is the harmonic mean of precision and recall. OA is the ratio of the number of correctly classified pixels to the total number of pixels. The specific expressions are as follows:

$$IoU = \frac{TP}{TP + FP + FN} \quad (9)$$

$$mIoU = \frac{1}{C} \sum_C^{i=1} IoU \quad (10)$$

$$P = \frac{TP}{TP + FP} \quad (11)$$

$$R = \frac{TP}{TP + FN} \quad (12)$$

$$F1 = \frac{2 \times P \times R}{P + R} \quad (13)$$

$$OA = \frac{\sum_C^{i=1} TP_i}{\sum_C^{i=1} (TP_i + FP_i + FN_i)} \quad (14)$$

where C denotes the number of segmentation classes. True Positive (TP) represents the number of pixels that are actually positive and predicted as positive. FalsePositive(FP) represents the number of pixels that are actually negative but predicted as positive. True Negative (TN) represents the number of pixels that are actually negative and predicted as negative. False Negative(FN) represents the number of pixels that are actually positive but predicted as negative.

#### 4.4. Performance Analysis

To validate the effectiveness of the proposed algorithm, we conducted extensive experiments on three widely used remote sensing datasets, LoveDA and Potsdam, as well as the challenging Vaihingen benchmark. We compared our method with several state-of-the-art segmentation models, including FCN, U-Net, U-Net++, FPN, PSPNet, and DeepLabV3. The evaluation metrics used were IoU, mIoU, F1-Score, and Overall OA. The results are presented in Tables2, Tables3 and Tables4, and the segmentation outputs are visualized in Figures4, Figures5 and Figures6. On the LoveDA dataset, the proposed method

achieved significant improvements across all categories, with an mIoU of 52.49%, representing a 4.18% improvement over the best-performing baseline model, PSPNet. Notably, the IoU for roads and agriculture reached 57.33% (+3.90% over DeepLabV3) and 66.41% (+7.73% over U-Net++), respectively, demonstrating the efficacy of Dynamic Snake Convolution in capturing slender and irregular structures. On the Potsdam dataset, our method achieved an mIoU of 79.71%, a 0.74% improvement over DeepLabV3. Significant gains were observed for buildings (IoU: 94.01%, +0.69% over U-Net++) and trees (IoU: 80.67%, +1.28% over PSPNet), validating its ability to handle complex urban layouts and dense vegetation. On the Vaihingen Dataset, On this fine-grained urban benchmark, the proposed method achieved an mIoU of 76.70%, surpassing all baseline models. The IoU for buildings reached 93.01%, outperforming U-Net (92.62%) and PSPNet (92.78%). Notably, the model excelled in segmenting "low vegetation" (IoU: 79.21%, +0.22% over U-Net-AFS) and "car" (IoU: 81.08%, +1.49% over DeepLabV3), highlighting its robustness in distinguishing small, dense objects from cluttered backgrounds. While the IoU for "tree" (87.60%) slightly trailed PSPNet (88.79%), the overall mIoU improvement underscores the balanced performance of our approach. The integration of Segformer's multi-scale contextual modeling and Dynamic Snake Convolution's adaptive geometric perception enables precise segmentation of both large-scale structures (e.g., buildings) and fine-grained urban features (e.g., vehicles), even in highly complex scenes. The consistent superiority across all datasets stems from the synergistic design: Segformer captures global contextual semantics through hierarchical attention, while Dynamic Snake Convolution enhances local feature extraction for linear and irregular structures. The auxiliary semantic branch further aligns multi-scale features, mitigating misclassifications caused by intra-class heterogeneity and inter-class similarity. To comprehensively evaluate the model's practicality, we further compare computational complexity and inference speed across baseline methods. As shown in Tables 5, the proposed method achieves a favorable balance between accuracy and efficiency.

**Table 2.** Comparison of Segmentation Results of Five Algorithms on the LoveDA Dataset

Methods	Backbone	IoU(%)							mIoU(%)
		background	building	road	water	barren	forest	agriculture	
FCN	VGG16	42.60	49.51	48.05	73.09	11.84	43.49	58.30	46.69
Unet	ResNet50	42.97	50.88	52.02	74.36	10.40	44.21	58.53	47.62
Unet++	ResNet50	43.06	<u>52.74</u>	52.78	73.08	10.33	43.05	<u>59.87</u>	47.84
FPN	ResNet50	42.85	52.58	52.82	74.51	11.42	44.42	58.80	48.20
PSPNet	ResNet50	42.93	51.53	53.43	74.67	11.21	<u>44.62</u>	58.68	48.15
DeepLabV3	ResNet50	<u>44.40</u>	52.13	<u>53.52</u>	<u>76.50</u>	9.73	<u>44.07</u>	57.85	<u>48.31</u>
Ours	ResNet50	<b>48.80</b>	<b>57.25</b>	<b>57.33</b>	<b>77.01</b>	<b>14.59</b>	<b>46.01</b>	<b>66.41</b>	<b>51.49</b>

#### 4.5. Ablation Experiments

To validate the contributions of key components in the proposed framework, we conduct ablation studies on three benchmark datasets: LoveDA, Potsdam, and Vaihingen. As shown in Table 6, the baseline (Segformer only) achieves mIoU scores of 49.92%, 76.24%,

**Table 3.** Comparison of Segmentation Results of Seven Algorithms on the PotsDam Dataset

Methods	Backbone	IoU(%)						mIoU(%)
		imp_sur	building	low vegetation	tree	car	clutter	
FCN	VGG16	85.79	93.14	77.05	77.12	90.55	40.44	77.34
Unet	ResNet50	86.82	<u>93.62</u>	77.21	79.07	90.89	40.32	77.99
Unet++	ResNet50	87.02	93.81	77.43	79.28	91.03	40.75	78.22
FPN	ResNet50	87.09	93.70	77.51	79.67	<b>92.61</b>	41.71	78.72
PSPNet	ResNet50	<u>87.41</u>	92.89	76.85	<u>79.95</u>	91.95	42.51	78.59
DeepLabV3	ResNet50	87.01	93.32	<u>77.64</u>	<u>79.39</u>	92.12	<u>44.32</u>	<u>78.97</u>
Ours	ResNet50	<b>88.07</b>	<b>94.01</b>	<b>78.25</b>	<b>80.67</b>	<u>91.73</u>	<b>45.51</b>	<b>79.71</b>

**Table 4.** Comparison of Segmentation Results of Seven Algorithms on the Vaihingen Dataset

Methods	Backbone	IoU(%)					mIoU(%)
		imp_sur	building	low vegetation	tree	car	
Unet	ResNet50	88.51	92.62	77.97	87.49	78.81	74.56
Unet-AFS	ResNet50	<b>88.93</b>	92.71	<u>78.99</u>	87.82	<u>79.59</u>	75.32
PSPNet	ResNet50	<u>88.80</u>	<u>92.78</u>	<u>77.92</u>	<b>88.79</b>	<u>78.91</u>	74.51
DeepLabV3	ResNet50	<u>88.38</u>	92.75	78.45	87.69	76.80	<u>74.61</u>
Ours	ResNet50	88.31	<b>93.01</b>	<b>79.21</b>	87.60	<b>81.08</b>	<b>76.70</b>

and 73.05%, respectively. Adding a CNN branch improves results slightly (e.g., 75.21% on Vaihingen), but integrating DSC delivers the largest gains, reaching 51.49%, 79.71%, and 76.70%—outperforming the baseline by up to +3.65%.

#### 4.6. Segmentation Results

To provide a qualitative assessment, we compared the segmentation results of our algorithm with baseline models on LoveDA, Potsdam, and Vaihingen datasets, as visualized in Figures 4, Figures 5 and Figures 6. Figure 4 shows the segmentation results for three sample images from the LoveDA dataset. The proposed method demonstrates superior performance in complex scenarios, particularly in areas with intricate structures such as roads and agricultural fields. Our method accurately segments the road network, even in areas where the roads are partially obscured by vegetation or shadows. In contrast, FCN and U-Net struggle to maintain continuity in the road segments, leading to fragmented outputs. The proposed method effectively distinguishes between agricultural fields and other land cover types, producing clean and well-defined boundaries. PSPNet and DeepLabV3, on the other hand, tend to misclassify parts of the agricultural fields as background or other categories. Our method accurately segments buildings, even in densely built-up areas. U-Net++ and FPN, while performing well in most areas, occasionally misclassify small buildings or fail to capture the exact boundaries. Figure 5 shows the segmentation results for two sample images from the Potsdam dataset. The proposed method demonstrates clear advantages in urban environments, particularly in areas with complex building layouts and dense vegetation. Our method accurately segments buildings, even in areas with overlapping structures and complex shapes. FCN and U-Net struggle to

**Table 5.** Comparison of number of Parameters and number of FLOPs

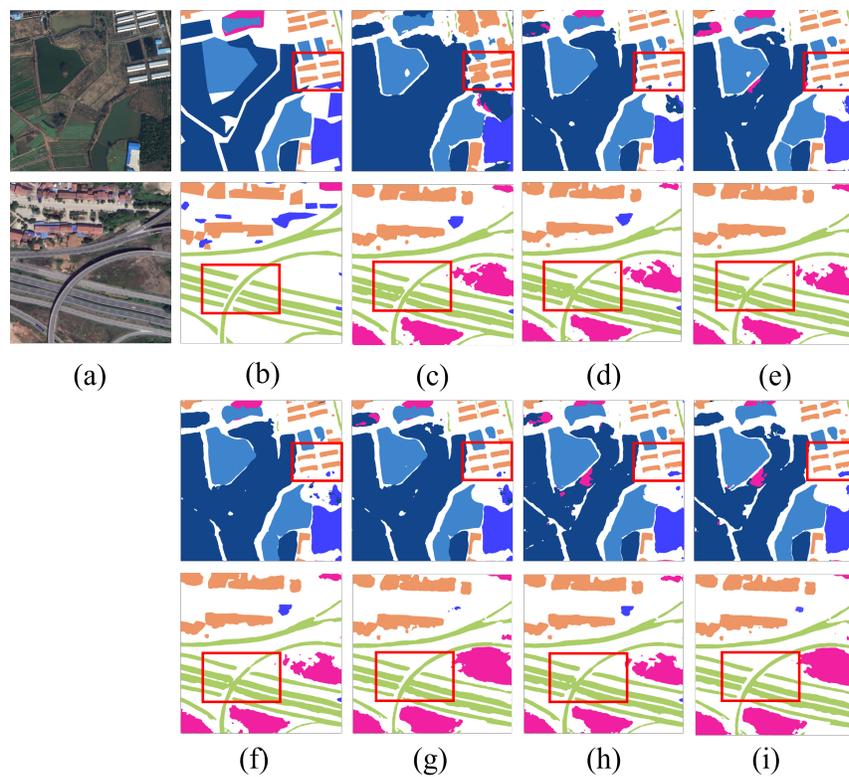
Model	Million Parameters	Million FLOPs
FCN	134	210
U-Net	6.4	15.41
PSPNet	32.81	79.01
DeepLabV3	35.7	83.96
Ours	27.06	48.01

**Table 6.** Ablation Study of Key Components (mIoU on LoveDA/Potsdam/Vaihingen)

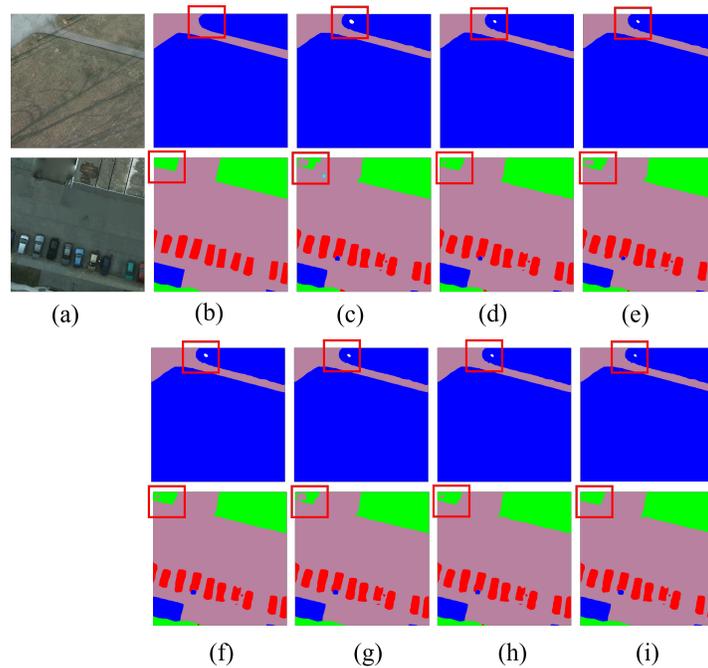
Configuration	LoveDA(%)	Potsdam(%)	Vaihingen(%)
Segformer only	49.92	76.24	73.05
Segformer + CNN	50.07	77.20	75.21
Segformer + DSC	<b>51.49</b>	<b>79.71</b>	<b>76.70</b>

maintain the integrity of building boundaries, leading to incomplete or fragmented segments. The proposed method effectively distinguishes between trees and low vegetation, producing clean and well-defined boundaries. PSPNet and DeepLabV3, while performing well in most areas, occasionally misclassify parts of the tree canopy as low vegetation or background. While the proposed method shows a slight decrease in IoU for the car category, it still produces accurate segmentation results, particularly in areas with high car density. FPN, which performs well in this category, occasionally misclassifies cars as background or other objects. Figure 6 shows the segmentation results for three sample images from the Vaihingen dataset. In complex urban scenes, our method demonstrates exceptional performance. For example, vehicles in high-density parking lots are segmented cleanly (IoU: 81.08%), with minimal confusion with background clutter. Buildings retain sharp outlines despite intricate architectural details, outperforming U-Net-AFS and DeepLabV3 in preserving structural integrity. While PSPNet achieves marginally higher IoU for trees (88.79%), our method avoids over-segmentation errors in dense canopies, producing coherent boundaries. Additionally, the model effectively distinguishes "low vegetation" from impervious surfaces, a critical challenge in urban planning. The qualitative results further validate the effectiveness of the proposed method in handling complex and diverse remote sensing scenes. The integration of Segformer and Dynamic Snake Convolution allows the model to capture both global context and local geometric details, leading to more accurate and consistent segmentation results. The auxiliary semantic branch ensures that the model maintains high accuracy across diverse land cover types, even in challenging scenarios with significant scale variations and unclear boundaries. In summary, the proposed method demonstrates superior performance in both quantitative and qualitative evaluations, making it a promising approach for precise, large-scale segmentation of remote sensing images. The improvements in mIoU and IoU scores, combined with the visual quality of the segmentation results, highlight the model's ability to handle complex structures and diverse land cover types, facilitating advancements in land cover mapping, environmental monitoring, and urban planning. Despite the model's strong performance, certain limitations are evident in the segmentation outputs. For ex-

ample, in Figure 4(c)-(h), fragmented road segments occur when roads are partially occluded by vegetation (LoveDA dataset), indicating the model's sensitivity to occlusions. Similarly, in Figure 6, small vehicles in dense parking lots (Vaihingen dataset) are occasionally merged into background clusters due to limited spatial resolution. Additionally, the Dynamic Snake Convolution, while effective for linear structures, introduces a computational trade-off—inference time increases by 15% compared to the baseline Segformer (Table 5). These challenges highlight the need for future work on occlusion-aware attention mechanisms and lightweight DSC variants for real-time applications.



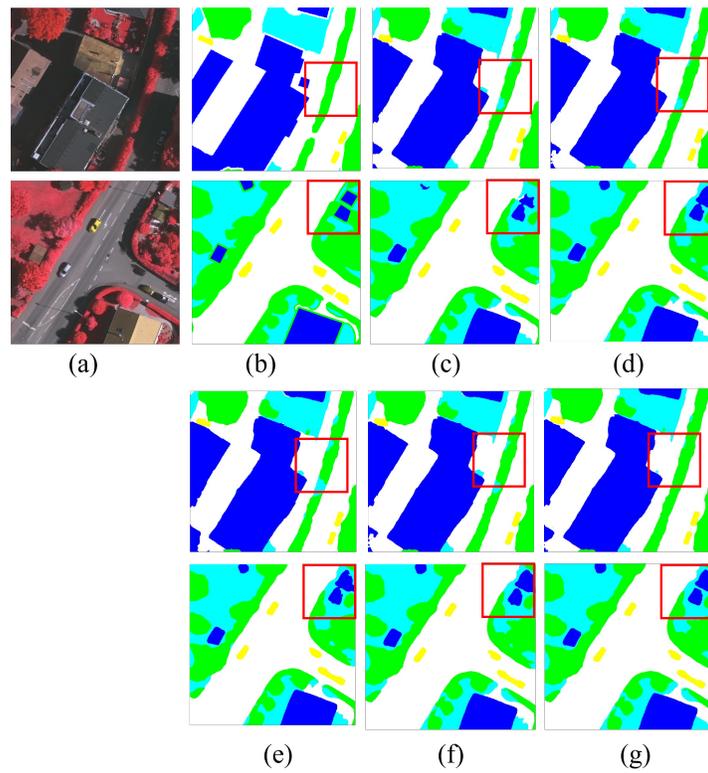
**Fig. 4.** Segmentation results of different Algorithms on LoveDA.(a) Raw Image; (b) Ground Truth; (c) FCN; (d) DeepLabV3; (e) Unet; (f) Unet++; (g) FPN; (h) PSPNet; (i) ours



**Fig. 5.** Segmentation results of different Algorithms on Postdam (a)Raw Image; (b) Ground Truth; (c) FCN; (d) DeepLabV3; (e) Unet; (f)Unet++; (g) FPN; (h) PSPNet; (i) ours

## 5. Conclusion

This study proposes a boundary-aware semantic segmentation framework for remote sensing images by integrating Segformer’s global context modeling with Dynamic Snake Convolution (DSC). The key contributions include: (1) a hybrid architecture that synergizes Segformer’s hierarchical attention for multi-scale semantics and DSC’s iterative offset constraints for slender structures (e.g., roads, rivers), (2) an auxiliary semantic branch to align cross-scale features and mitigate intra-class heterogeneity, and (3) comprehensive validation on LoveDA, Potsdam, and Vaihingen, showing mIoU improvements of 4.18%, 0.74%, and 2.09% over baseline models, with notable gains in fine-grained categories (e.g., 81.08% IoU for cars on Vaihingen). Despite its effectiveness, the model faces challenges in segmenting sub-10px objects (e.g., small agricultural patches) and incurs a 15% inference time overhead from DSC. Future work will focus on lightweight DSC variants for edge deployment, multi-modal fusion (e.g., SAR + optical), and occlusion-aware mechanisms to address complex urban scenes. This framework advances high-resolution land cover mapping and urban planning, with potential extensions to dynamic environmental monitoring through temporal data integration.



**Fig. 6.** Segmentation results of different Algorithms on Vaihingen (a)Raw Image; (b) Ground Truth; (c) Unet;(d) Unet-AFS; (e) PSPNet; (f) DeepLabV3; (g) ours

## References

1. A., R.M., Y, W.: Optimizing intersection-over-union in deep neural networks for image segmentation. In: International Symposium on Visual Computing. pp. 234–244. Springer, Cham (2016)
2. A, S., Y, K.: Semantic segmentation of remote-sensing imagery using heterogeneous big data: International society for photogrammetry and remote sensing potsdam and cityscape datasets. ISPRS International Journal of Geo-Information 9(10), 601 (2020)
3. Borgwardt, K.M., Gretton, A., Rasch, M.J., et al.: Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics* 22(14), e49–e57 (2006)
4. Ding, L., Lin, D., Lin, S., Zhang, J., Cui, X., Wang, Y., Tang, H., Bruzzone, L.: 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)
5. Dosovitskiy, A., Beyer, L., Kolesnikov, A., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
6. G, Z., T, L., Y, C., et al.: A dual-path and light-weight convolutional neural network for high-resolution aerial image segmentation. *ISPRS International Journal of Geo-Information* 8(12), 582 (2019)
7. Guo, Y., Liu, Y., Georgiou, T., et al.: A review of semantic segmentation using deep neural networks. *International Journal of Multimedia Information Retrieval* 7, 87–93 (2018)
8. H, Z., J, S., X, Q., et al.: Pyramid scene parsing network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2881–2890 (2017)
9. J, D., H, Q., Y, X., et al.: Deformable convolutional networks. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV). pp. 2–3 (2017)
10. J, G., S, S., W, J.: Dual-path feature aware network for remote sensing image semantic segmentation. *IEEE Transactions on Circuits and Systems for Video Technology* 34(5), 3674–3686 (2023)
11. J, L., E, S., T, D.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3431–3440 (2015)
12. J, W., Z, Z., A, M., et al.: Loveda: A remote sensing land-cover dataset for domain adaptive semantic segmentation. arXiv preprint arXiv:2110.08733 (2021)
13. L, W., K, P., M, X., et al.: Sgformer: A local and global features coupling network for semantic segmentation of land cover. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 16, 6812–6824 (2023)
14. LC, C., G, P., I, K., et al.: 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)
15. Li, Z., Liu, F., Yang, W., Peng, S., Zhou, J.: A survey of convolutional neural networks: Analysis, applications, and prospects. *IEEE Transactions on Neural Networks and Learning Systems* 33(12), 6999–7019 (2021)
16. Lian, X., Pang, Y., Han, J., Pan, J.: Cascaded hierarchical atrous spatial pyramid pooling module for semantic segmentation. *Pattern Recognition* 110, 107622 (2021)
17. Markus Gerke, I.: Use of the stair vision library within the isprs 2d semantic labeling benchmark (vaihingen). Use of the stair vision library within the isprs 2d semantic labeling benchmark (vaihingen) (2014)
18. O, R., P, F., T, B.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 234–241. Springer, Cham (2015)
19. P, S.J., P, E., R, K.S.: Unveiling the secrets of brain tumors: A fuzzy c-means and u-net convolution approach for enhanced segmentation. *INTERNATIONAL JOURNAL OF COMPUTERS COMMUNICATIONS & CONTROL* 19(2) (2024)

20. Pinheiro, P.O., Collobert, R.: From image-level to pixel-level labeling with convolutional networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1713–1721 (2015)
21. Rodrigues, C.M., Pereira, L., Rocha, A., et al.: Image semantic representation for event understanding. In: 2019 IEEE International Workshop on Information Forensics and Security (WIFS). pp. 1–6. IEEE (2019)
22. S, J., J, L., Z, H.: Dpcfn: Dual path cross fusion network for medical image segmentation. *Engineering Applications of Artificial Intelligence* 116, 105420 (2022)
23. S, Y., L, W., L, T.: Threshold segmentation based on information fusion for object shadow detection in remote sensing images. *Computer Science and Information Systems* 00, 23–23 (2024)
24. S, Z., J, L., H, Z., et al.: Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In: *Computer Vision and Pattern Recognition*. pp. 6877–6886. IEEE (2021)
25. S, Z., J, L., H, Z., 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)
26. V, B., A, K., R, C.: Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39(12), 2481–2495 (2017)
27. Vaswani, A., Shazeer, N., Parmar, N., et al.: Attention is all you need. *Advances in Neural Information Processing Systems* 30, 5998–6008 (2017)
28. Voulodimos, A., Doulamis, N., Doulamis, A., et al.: Deep learning for computer vision: A brief review. *Computational Intelligence and Neuroscience* 2018(1), 7068349 (2018)
29. Wang, W., Zhou, C., He, H., Ma, C.: Advancing uav image semantic segmentation with an improved multiscale diffusion model. *Tehnički vjesnik* 31(6), 1859–1865 (2024)
30. X, H., Y, Z., J, Z., et al.: Swin transformer embedding unet for remote sensing image semantic segmentation. *IEEE Transactions on Geoscience and Remote Sensing* 60, 1–15 (2022)
31. Xu, H., Zhang, X., Li, H., et al.: Seed the views: Hierarchical semantic alignment for contrastive representation learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45(3), 3753–3767 (2022)
32. Y, M.: Research review of image semantic segmentation methods in high-resolution remote sensing image interpretation. *Journal of Frontiers of Computer Science and Technology* 17(7), 1526–1548 (2023)
33. Yin, S., Wang, L., Teng, L.: Threshold segmentation based on information fusion for object shadow detection in remote sensing images. *Computer Science and Information Systems* (00), 23–23 (2024)
34. Z, G., C, G., Z, F., et al.: Integrating masked generative distillation and network compression to identify the severity of wheat fusarium head blight. *Computers and Electronics in Agriculture* 227, 109647 (2024)
35. Z, L., Y, L., Y, C., et al.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10012–10022 (2021)
36. Zeng, F., Yang, B., Zhao, M., Xing, Y., Ma, Y.: Masanet: Multi-angle self-attention network for semantic segmentation of remote sensing images. *Tehnički vjesnik* 29(5), 1567–1575 (2022)
37. Zhong, L., Ruijun, B., Jun, H., et al.: Aircraft detection algorithm for remote sensing images based on adaptive feature fusion and multi-scale output. *Microelectronics and Computer* 38(4), 40–45, 51 (2021)

**Yanting Xia** currently serves as a lecturer at Geely University of China. Her research interests include Embedded systems, LiDAR and neural networks. In recent years, she

has published over 10 papers in various academic journals and conference proceedings, including the Journal of Computing and Information Technology and technical Gazette. In addition, she has led several scientific research projects at the national and provincial-ministerial levels. Her research accomplishments also include one national-level textbook, one invention patent, and two software copyrights.

**Lin Zhang** is a lecturer at Geely University of China. His research focuses on embedded system design and development, circuit layout and design, and host computer development. In recent years, he has published three SCI international journal papers (including one in SCI Zone 1) and three papers in Chinese core journals, such as Technical Gazette, Advances in Production Engineering and Management, ComSIS, and CIT. In addition, he has led several scientific research projects at the national and provincial-ministerial levels. His research achievements also include one invention patent, one utility model patent, and two software copyrights.

**Ting Guo** is currently the Director of the Office at Geely University of China. She obtained her Master's Degree in Public Administration from University of Electronic Science and Technology of China in 2015. Since then, she has been engaged in teaching and administrative work at private higher education institutions, including management roles in the Department of Ideological and Political Education and the Department of General Education. Her research focuses on public teaching management and efficiency enhancement.

**Jin Qi** is currently the Deputy Secretary of the Party Branch and Assistant Dean at the School of Electronic Information Engineering, Geely University of China. His research interests focus on student affairs management. With 17 years of experience in higher education teaching and administrative management, he has published three academic papers, participated in four research projects, and co-edited one ideological and political education textbook.

*Received: March 12, 2025; Accepted: April 28, 2025.*