

RESEARCH ARTICLE

A Lightweight High-Resolution RS Image Road Extraction Method Combining Multi-Scale and Attention Mechanism

RUI WANG¹, MINGXIANG CAI¹, AND ZIXUAN XIA²¹China Transport Telecommunications and Information Center, Beijing 100011, China²College of Art and Architectural Engineering, Heilongjiang University of Technology, Jixi 158100, China

Corresponding author: Mingxiang Cai (jamtasai@whu.edu.cn)

ABSTRACT Road information plays an indispensable role in human society's development. However, owing to the diversity and complexity of roads, it is difficult to obtain satisfactory road-extraction result. Some typical factors, such as discontinuity, loss of edge details, and long-time consumption, have negative impacts on obtaining accurate road information. These problems are particularly prominent during road extraction when high-resolution remote-sensing images are used. To obtain accurate road information, a novel lightweight deep learning neural network was proposed in this study by integrating a multiscale module and attention mechanisms. As an excellent multiscale segmentation module, the atrous spatial pyramid pooling was selected to enhance the road extraction ability of remote sensing images. In addition, an attention mechanism was employed to solve the problems of discontinuity and loss of edge details in road extraction, and MobileNet V2 was selected as the backbone of DeepLab V3+ because of its lightweight structure, which can help solve the problem of excessive training time consumption. The experimental verification was carried out on the Ottawa road dataset and the Massachusetts road dataset. Experimental results show that compared with U-Net, SegNet and MDeepLab v3+ networks, the proposed algorithm is the best in IoU, Recall, OA and Kappa. Among them, on the Ottawa road dataset, the OA and Kappa of the algorithm in this paper are 98.92 % and 95.02 %, respectively. On the Massachusetts road dataset, OA and Kappa 98.29% and 89.87%. In addition, the training time was significantly shorter than that of the other deep learning networks. The proposed method exhibited a good performance in road extraction.


INDEX TERMS Road extraction, deep learning, CAM, SAM, ASPP, lightweight.

I. INTRODUCTION

As important geo-information, roads are indispensable in modern society. Road information not only plays an important role in the transportation and insurance industry but also in domains such as vehicle navigation, unmanned driving, urban management planning, smart city construction, and geo-information database updating. Owing to their wide range of applications [1], [2], road extraction has become a re-search hotspot worldwide. With the development of remote sensing (RS) technology, the resolution of RS images has increased, and road extraction based on high-resolution remote sensing images can be effectively

and accurately achieved. However, obtaining road information from high-resolution RS images remains a challenging task owing to the presence of complex backgrounds, diverse topologies, and shadows.

In recent years, scholars have done a lot of research on road extraction by fully utilizing high-resolution RS images [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15]. They proposed a variety of methods for different road extraction tasks. For example, road extraction can be divided into two categories according to different emphasis points: (1) paying attention to local feature information. Local feature information of the road itself is used, such as spectral, texture, shape, and structural features. Road skeletons are extracted using template matching [4], [5] or knowledge-driven [6], [7], [8], [9] algorithms. However, these methods rely on

The associate editor coordinating the review of this manuscript and approving it for publication was Geng-Ming Jiang .

manual work and have a low degree of automation [10]. (2) Paying attention to the overall feature information. The road is regarded as the extraction object, and the road information in a high-resolution image is extracted through classical algorithms, such as image segmentation clustering [11], Support Vector Machine (SVM) [12], and conditional Random Field (RF) [13]. However, such methods generally suffer from poor stability and complicated operational steps. Extraction results from traditional methods often suffer from issues such as discontinuous extraction and rough edge details, resulting in poor extraction results [14]. Consequently, these methods cannot meet the current requirements for road formation extraction from high-resolution RS images [15].

At present, the rapid development of neural networks has a promoting effect on RS information extraction technology. Long et al. [16] proposed a Fully Convolutional neural network (FCN) in 2015. Realizing from simple object classification to pixel-level classification greatly improves the classification accuracy and efficiency. Subsequently, there has been an increasing number of RS image classifications based on the FCN, especially in the task of road extraction. First, the U-Net and SegNet networks are widely used in road extraction after being proposed. Kong et al. [17] proposed an improved U-Net network to extract road information from high-resolution RS images. This method adds a stripe pooling module to the down-sampling part of the coding layer to focus on local information. A hybrid pooling module was added to the convolution of the coding layer to enhance its ability to obtain the network context. The algorithm was verified using the Gaofen-2 high-resolution RS image dataset. The results showed that this algorithm could effectively extract the road, and the extraction effect was better than that of the traditional network. Li et al. [18] proposed an improved U-Net network that combines core and global attention to extract roads from high-resolution RS images. Experiments were performed on the Massachusetts road dataset and the DeepGlobe-CVPR 2018 road dataset, and the results showed that the proposed method can effectively extract the road area blocked by the tree canopy and improve the connectivity of the road network. Ai [19] proposed an improved SegNet network to extract road information. An edge feature pyramid module was added to the original SegNet network. Through multi-scale convolution of the edge pyramid module, the purpose of enhancing the recognition of small targets was achieved. The experimental results on the Cityscapes dataset showed that the accuracy was improved by 7.5 % and 6.2 %, respectively, compared with the original SegNet and U-Net networks. The effectiveness of the pyramid module in this method was demonstrated in the ablation experiment. Subsequently, the Google team proposed the DeepLab series of network models [20], [21], [22], [23], which scholars have widely used for road extraction. Han et al. [24] added dilated convolution based on the DeepLab V3 network to extract road information. The experiment verified that this method

has a higher extraction accuracy than the other methods. Liu et al. [25] proposed the DeepLab V3+ road extraction method with an attention module. This method obtains more spatial context information through spatial attention to enhance the extraction of road information. The effectiveness of this method was verified by using the Cityscapes dataset. The results show that the ability of this method to extract road information was significantly better than that of the original network.

Although these methods have their own advantages in road extraction, they still have some shortcomings. First, the FCN has a large number of parameters, which not only depends on superior hardware equipment but also takes a long time for network training. Second, SegNet and U-Net networks have the problem of incomplete context information extraction during road extraction. Similarly, the DeepLab V3 network still has problems, such as missing details and discontinuous extraction. To solve these problems, the accuracy of road extraction was further improved. This study is based on the DeepLab V3+ network, and a lightweight high-resolution RS image road extraction method is proposed that integrates multi-scale and attention mechanisms. First, atrous spatial pyramid pooling (ASPP) is an excellent multi-scale feature extraction module. By using a multiscale convolution kernel, the ASPP can extract roads of different sizes. Second, the reference channel attention enhances the shallow features, and the reference space attention enhances the deep features to solve the problems of edge detail loss and extraction discontinuity. Finally, the lightweight MobileNet V2 network replaced the Xception network used in the DeepLab V3+ network as the backbone network. The problem of a long training time caused by a large number of model parameters was solved.

II. GUIDELINES FOR MANUSCRIPT PREPARATION

A. EXPERIMENTAL DATA

The experimental data of this paper are Ottawa road dataset and Massachusetts road dataset. Among them, the Ottawa road dataset has a resolution of 0.2 meters. The Massachusetts road dataset has a resolution of 1 meter. The image data with dense roads in the two data sets is selected as the training samples of roads. The training sample is cropped to a size of 512×512 pixels, and the data set is expanded by rotating 90° , 180° horizontal mirroring and vertical mirroring. The Ottawa-Dataset and Massachusetts-Dataset images and labels are shown in Figure 1.

B. PARAMETER SETTINGS

The experimental environment in this study was a Windows 10 operating system with 128G memory. The processor is an Intel(R) CPU E5-2640 v3@2.60GHz. The graphics card was a Nvidia GTX 1600 super 6 GB graphics card. The deep learning framework is a high-level neural network application interface Keras 2.1.6 of Tensorflow 1.14.0. SegNet and

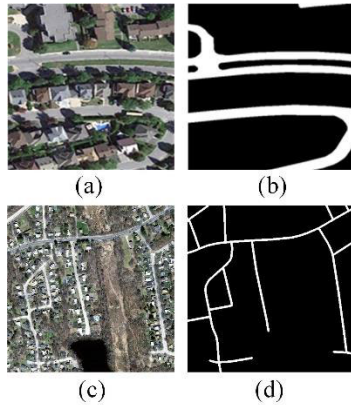


FIGURE 1. Dataset images and labels (a) Image of Ottawa-Dataset; (b) Corresponding label of Ottawa-Dataset; (c) Image of Massachusetts -Dataset; (d) Massachusetts -Dataset corresponding label.

U-Net were selected as comparison algorithms. According to the experimental environment of this study, the batch size was set to 2, and the epoch was set to 100 for the training network.

III. METHODS

The method in this study is based on the DeepLab V3+ network [26], [27], [28] to make two improvements. The first is that the lightweight MobileNet V2 network replaces the Xception of the DeepLab V3+ network as the backbone network, and the second is to enhance the deep features by referring to channel attention and spatial attention at the jump connection structure.

A. LIGHTWEIGHT ROAD EXTRACTION NETWORK

The method in this study adopts an encoder-decoder structure. First, in the encoder layer, the shallow road feature information in high-resolution RS images is extracted through the different channel-depth separable convolution layers of MobileNet V2. Shallow feature information is obtained using atrous spatial pyramid pooling (ASPP) to obtain deep semantic road feature information. In the ASPP module, shallow

features are processed through a 1×1 convolution, three dilated convolutions with dilation rates of 6, 12, and 18, and a global-average pooling layer. After processing, the deep feature information with multiscale features was obtained. Second, the deep feature information is input to the decoder layer together with the shallow feature information under the enhancement of the channel attention and spatial attention. Finally, in the decoder layer, the deep features are up-sampled by 4 times bilinear interpolation four times and combined with the shallow feature information of the remote sensing image road obtained by the encoder layer. The combined feature information passes through two 3×3 convolutional layers to restore detailed feature information and obtain fine target boundary information through four bilinear interpolation upsamplings. Through the above, the extraction results of the high-resolution remote sensing image road are finally obtained.

B. MibleNet V2

As a lightweight deep neural network, the MobileNet network [29] has the advantages of smaller volume, less computation, higher accuracy, and faster speed, and is suitable for a variety of application scenarios. Its core is a depth-wise separable convolution, and its structure is shown in Figure 3. First, a 3×3 convolution kernel separated the input feature channels by traversing the individual data in each channel. Second, a 1×1 convolution kernel was used to traverse each feature graph to integrate the feature information to obtain the output feature. Compared with standard convolution, depth-separable convolution can effectively reduce the number of model parameters.

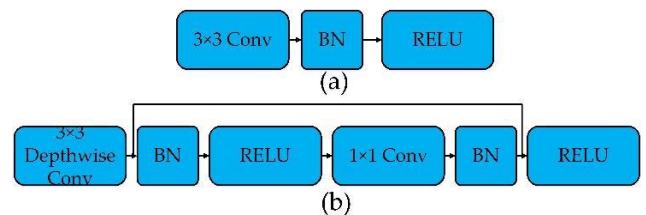


FIGURE 3. Deeply separable convolutional structures. (a) Standard convolution;(b) Depthwise sep-arable convolution.

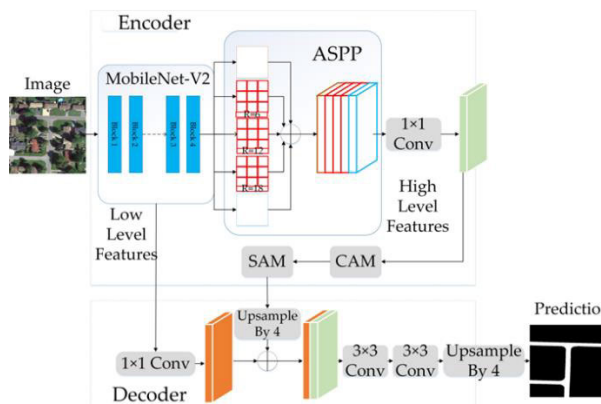


FIGURE 2. Lightweight DeepLab V3+ network.

MobileNet V2 [30], [31], [32], [33] is an improved version of the MobileNet model, as shown in Figure.4(a) in Figure.4(b). It has two distinct features. The first is to invert the residual structure. This structure first uses a 1×1 convolution before the 3×3 network structure to increase the dimensions. This structure then uses a 1×1 convolution after the 3×3 network structure to achieve dimensionality reduction and compression. Compared with the direct use of a 3×3 convolution, the effect is better, and the number of parameters is effectively reduced. The second is a linear bottleneck structure. The bottleneck structure of MobileNet V2 is shown in Figure 4. Because the linear ReLU function operation causes feature loss, to reduce the information loss,

the ReLU operation is no longer performed after the 1×1 convolution dimension reduction. However, the addition of the residual network was performed directly. The expansion coefficient was designed in MobileNet V2, the purpose of which is to achieve better control of the network size.

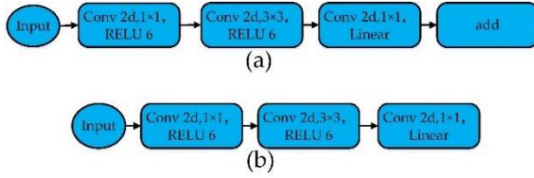


FIGURE 4. The bottleneck structure of MobileNet V2. (a) Step size is 1; (b) Step size is 2.

Compared with MobileNet, MobileNet V2 has a deeper network structure, a complete convolutional layer with 32 channels, and 17 bottleneck structures, as shown in Table 1. where t is a multiple of the internal dimensions of the bot-tleneck structure, c is the number of channels, n is the number of repetitions of the bot-tleneck structure, and s is the step size.

TABLE 1. MobileNet V2 network architecture.

Input	Operation	t	c	n	s
2242×3	Covn2d	-	32	1	2
1122×32	bottleneck	1	16	1	1
562×16	bottleneck	6	24	2	2
282×24	bottleneck	6	32	3	2
142×32	bottleneck	6	64	4	2
142×64	bottleneck	6	96	3	1
72×96	bottleneck	6	160	3	2
72×160	bottleneck	6	320	1	1
72×320	Covn2d 1×1	-	1280	1	1
72×1280	avgpool 7×7	-	-	1	-
$1 \times 1 \times 1280$	Covn2d 1×1	-	k	-	-

C. ASPP MODULE

Figure 5 shows the ASPP module [22], [34], [35]. The ASPP module in the DeepLab V3+ network contained a 1×1 convolution layer, hole convolution with hole rates of 6, 12, and 18, and a global average pooling layer. In addition, a 1×1 convolutional layer was added after the global average pooling layer and each dilated convolutional layer to adjust the dimensionality of the feature maps. The purpose of this method was to achieve the same dimensionality as the output feature map. Void convolution with different void-rate sizes can capture road feature information of various sizes and enhance the ability of the network to extract roads of different sizes. The 1×1 convolution in the ASPP is used to capture smaller targets, whereas global average pooling can integrate the information of the entire feature map. Finally, the concat operation is performed on the feature maps obtained by the 1×1 convolution, hole convolution with hole ratios of 6, 12, and 18, and global average pooling layer. The dimension of the output feature map is adjusted using a 1×1 convolution so that it is equivalent to the input.

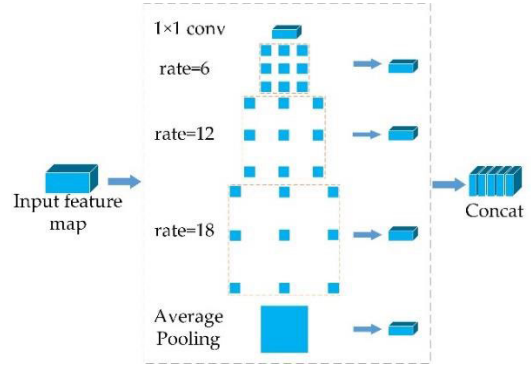


FIGURE 5. ASPP module.

D. ATTENTION MECHANISM

The core idea of the channel attention module (CAM) [36], [37], [38] is to learn the features of different feature channels to obtain weights. Then, CAM enhances the useful feature channels according to the weight, suppresses the useless feature channels, and realizes attention to the channel. The formula is as follows:

$$\begin{aligned} M_c(F) &= \sigma(MLP(AvgPool(F)) + MLP(MaxPool(F))) \\ &= \sigma(W_1(W_0(FC_{Avg})) + W_1(W_0(FC_{max}))) \end{aligned} \quad (1)$$

The specific implementation and structure of the CAM used in this study are shown in Figure 6. First, the feature graph F is processed by global average pooling and global maximum pooling. Global average pooling integrates the global spatial information of the input feature map F . Global maximum pooling reduces unnecessary information by extracting the maximum value of pixel points in the neighborhood. Second, the dimension is reduced through the first full connection layer to reduce the complexity. Subsequently, the number of channels is restored through the second full connection layer to build the correlation between channels. Finally, the weight of each channel was activated by the sigmoid function. It is applied to the feature channel corresponding to feature map F . Feature map MC is obtained after weighting processing by the channel attention module.

E. SPATIAL ATTENTION MODULE

The spatial attention module (SAM) [39], [40], [41] is to further improve the screening ability of salient features by focusing on spatial features and according to their importance. Its expression is shown in (2).

$$\begin{aligned} M_s(F) &= \sigma(f^{7 \times 7}([AvgPool(F), MaxPool(F)])) \\ &= \sigma(f^{7 \times 7}([F_{Avg}^C; F_{max}^S])) \end{aligned} \quad (2)$$

The implementation and structure of the SAM adopted in this study are shown in Figure 7. First, SAM compresses the channel-domain features of the input feature map F through global maximum pooling and global average pooling. Second, multiple channels are compressed into a single channel by convolution to reduce the influence of

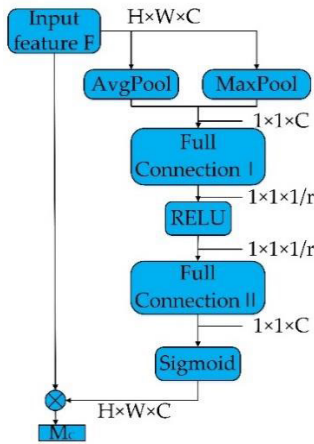


FIGURE 6. Channel attention module.

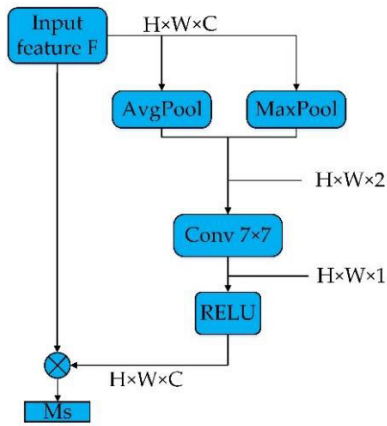


FIGURE 7. Spatial attention module.

inter-channel information on spatial attention. Finally, after the Sigmoid function was activated, the weight containing the spatial information feature was obtained. It is multiplied by the input feature map F pixel-by-pixel to obtain the feature map M_s .

F. ACCURACY EVALUATION

To effectively verify the effectiveness of the algorithm in this study, the average intersection-over-union ratio (IoU), Kappa coefficient, overall precision (OA), and Recall rate as the evaluation indicators of the experimental results. As shown in (3)–(6). The evaluation indices of network model complexity are model Parameters, Floating-Point operations (FLOPS), and training time. The number of floating-point operations represents the number of FLOPS operations during the forward inference process of the model.

$$IoU = TP / (TP + FN + FP) \tag{3}$$

$$Recall = TP / (TP + FN) \tag{4}$$

$$OA = (TP + TN) / (TP + FN + FP + TN) \tag{5}$$

$$Kappa = (P_0 - P_e) / (1 - P_e) \tag{6}$$

where TP represents the number of positive samples correctly extracted from the road. FN represents the number of positive samples directly extracted from the road. FP represents the number of negative samples incorrectly extracted from the road. TN represents the number of negative samples for which the road was extracted correctly. P_0 is the sum of the correctly extracted road pixels divided by the total number of pixels. P_e is the sum of the product of the real number of road samples and number of predicted samples divided by the square of the total number of samples.

IV. EXPERIMENT AND ANALYSIS

A. COMPARISON OF TRADITIONAL METHODS

To verify the superiority of the proposed method in road extraction, it was compared with two traditional RS image segmentation methods, namely support vector machine and random forest. SVM [41], [42] and RF [43], as shallow machine learning algorithms, exhibit good results in RS images [44]. In this study, through multiple verifications and referring to the relevant experimental results. In the SVM classification method in this study, the radial basis function was selected as the kernel function and the gamma value was set to 0.333. The penalty coefficient was maintained at a default value of 100. Finally, the number of decision trees in the RF classification was set to 100. The number of feature variables participating in the construction of decision trees was set as the square root of the number of all features of the classification. The image classification results are shown in Figure 8, where (a) is the image, (b) is the label map, (c) is the SVM extraction result, (d) is the RF extraction result, and (e) is the extraction result of the proposed method.

From Figure 8(c1) and (c2), the SVM road result extraction is discontinuous, and the road edge details are significantly lost. Buildings have a significant influence on the extraction of SVM. When the color of the roof is similar to that of the road, the extraction error becomes more serious. In places that are partially or completely blocked by obstacles, such as trees and buildings, the extraction effect is poor, and there are many misextraction parts. In Figure 8(d1) and (d2), the road extraction results of RF are more serious than those of SVM, and the mis-extracted part of SVM is slightly improved, but the problem of misallocated buildings still exists. In Figure 8(e1) and (e2), the algorithm in this study can effectively solve the problems of discontinuous road extraction and loss of road-edge details. There were no problems, such as false mentions. Buildings and roads can be effectively distinguished, and the occlusion problem does not affect the extraction results, which are basically consistent with the annotated data.

The precision calculation of pixel-by-pixel matching was performed on the experimental results and RS image label map. The accuracy evaluation results are shown in Table 2. It can be seen from Table 1 that the algorithm in this study is the best in terms of the accuracy indicators of IoU, Recall, OA, and Kappa. Among them, in the comparison experiment of Ottawa road dataset, compared with SVM, IoU and Recall

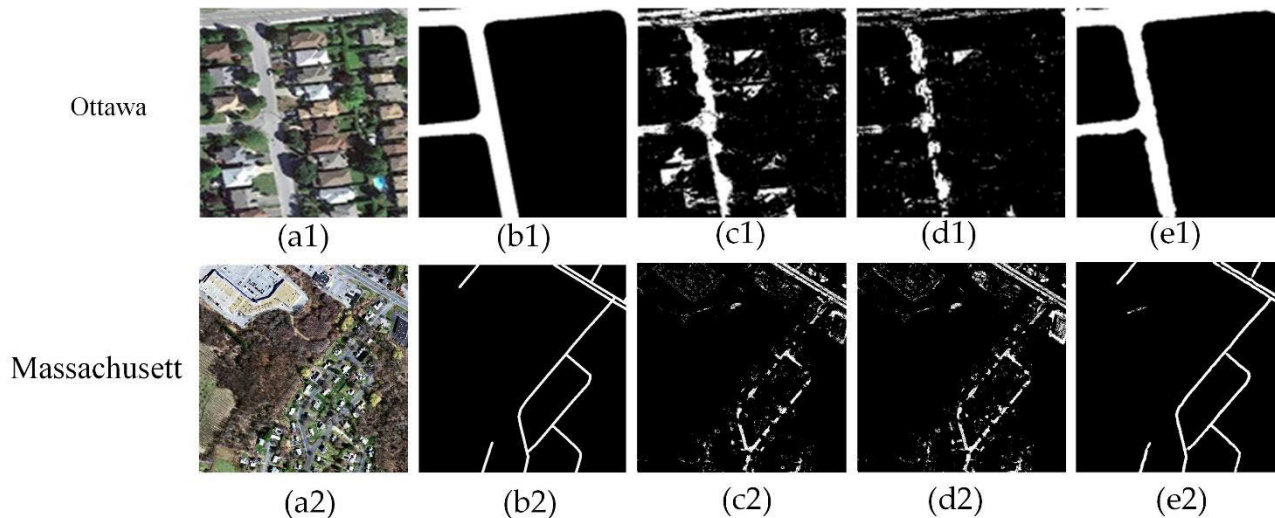


FIGURE 8. Traditional method results. (a) Image map; (b) Label map; (c) SVM extraction results; (d) RF extraction results; (e) Extraction results of this paper.

TABLE 2. Accuracy verification results of traditional methods.

DATA	Method	IoU/%	Recall/%	OA/%	Kappa
Ottawa	SVM	46.98	46.03	85.95	0.4601
	RF	43.29	42.53	77.87	0.4233
	Ours	91.63	95.07	98.92	0.9502
Massachusetts	SVM	64.28	49.10	96.06	0.4709
	RF	62.52	51.32	95.31	0.4352
	Ours	90.65	90.11	98.29	0.8987

increased by 44.65 and 49.04% respectively. OA increased by 12.97% and Kappa increased by 0.4904. Compared with RF, IoU increased by 48.34%, recall rate increased by 52.54%, OA increased by 12.97%, kappa increased by 0.5269. In the comparative experiment of the Massachusetts road dataset, IoU increased by 26.37 %, Recall increased by 41.01 %, OA increased by 2.23 %, and Kappa increased by 0.4278. Compared with RF, IoU increased by 28.13 %, recall rate increased by 41.01 %, OA increased by 2.98 %, kappa increased by 0.4635. It can be seen from the above indicators that the method in this study has more advantages than the traditional method for high-resolution road extraction. The method used in this study does not require post-processing and can achieve automatic and accurate extraction.

B. COMPARISON OF SEMANTIC SEGMENTATION METHODS

In order to verify the advantages of the algorithm in this paper in road extraction, the Ottawa road data set is used for experiments. Through U-Net, SegNet, the DeepLab V3+ network (hereinafter referred to as MDeepLab V3+) with MobileNet V2 as the backbone network and the algorithm in this paper are compared and tested. To ensure the comparability and validity of the results, the training and test datasets

are the same. The experimental results of the U-Net, SegNet, MDeepLab V3+ network, and algorithm used in this study are shown in Figure 9. Where (a) is the image, (b) is the label map, (c) is the U-Net extraction result, (d) is the SegNet extraction result, (e) is the MDeepLab V3+ extraction result, and (f) is the extraction result of the proposed method.

As shown in Figure 9(a1), the road was clear and regular. There are trees blocking it, and there are buildings with colors similar to the road in the figure. As shown in Figure 9(c1), U-Net can avoid the extraction discontinuity problem caused by obstacles such as trees and buildings in road results. However, the processing of edge details is not ideal, and holes remain in the middle. As shown in Figure 9(d1), the extraction result of SegNet is slightly worse than that of U-Net. There are discontinuities in the extraction results, and the detailed processing was slightly worse. However, the hole phenomenon is weaker. As can be seen in Figure 9(e1), the edge details of the extraction results of MDeepLab V3+ are better processed. The discontinuities and holes in the extraction results of U-Net and SegNet are weakened. This is owing to the multi-scale extraction of road details in the ASP module in the MDeepLab V3+ network. However, some problems remain, such as incomplete extraction. As shown in Figure 9(f1), the method proposed in this study can effectively extract roads. The edge details were handled well, and the extraction was complete and continuous. There were no problems such as false mentions or holes. As shown in Figure 9(a2), the road is irregular and there are many trees. There were no occlusions, such as buildings or trees. From Figure. 9(c2), it can be observed that the road extraction is discontinuous. There were large fractures associated with this problem. There are many holes in the middle of the road. There is a noticeable loss in the edge detail. In Figure 9(d2), SegNet shows obvious improvement compared with the U-Net extraction results, but the above problems still exist.

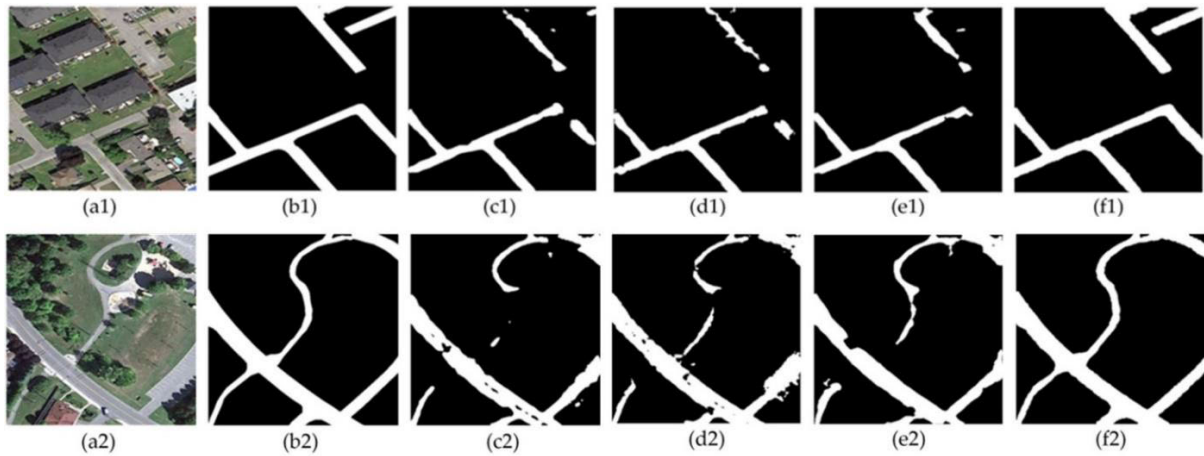


FIGURE 9. Compare the experimental results. (a) Image map; (b) Label map; (c) U-Net extraction results; (d) SegNet extraction results; (e) MDeepLab V3+ extraction results; (f) Extraction results of this paper.

TABLE 3. Comparison of experimental accuracy verification results.

Method	IoU/%	Recall/%	OA/%	Kappa	Params/MB	FLOPs/G	Training time/t
U-Net	61.67	64.57	95.34	0.7381	7.77	30.51	13.55
SegNet	77.60	80.76	97.40	0.8595	30.5	75.20	15.99
MDeepLabv 3+	80.21	81.80	97.23	0.8776	2.76	10.517	10.06
Ours	91.63	95.07	98.92	0.9502	2.85	12.76	10.08

It can be observed from Figure 9(e2) that the extraction result of MDeepLab V3+ is better. It can be accurately extracted, and there is no extraction discontinuity. However, the processing of edge details is still imperfect. After the method of this study integrates multiscale and attention mechanisms, the extraction effect is better than the extraction results of SegNet, U-Net, and MDeepLab V3+. As shown in Figure 9(f2), the method proposed in this study handles the edge details perfectly, and the extraction results are continuous. The problems with the extraction of the other three methods were solved. The extraction results are basically consistent with the labeled samples.

As can be seen from Table 3, the method presented in this paper is optimal for all indices. Among them, Compared with U-Net, IoU increased by 29.96%, Recall increased by 30.5%, OA increased by 3.58%, and kappa increased by 0.1821. Compared with SegNet, IoU increased by 14.03%, Recall increased by 14.31%, OA increased by 1.52%, and Kappa increased by 0.0907. Compared with MDeepLab V3+, all indicators have also been greatly improved, among which IoU and Recall have increased by more than 10%. The method in this study is much smaller than the U-Net and SegNet methods in terms of Params and FLOPs, which effectively reduces the training time. Compared with the MDeepLab V3+ network, the method in this study increases the attention mechanism. Although the performance is slightly improved in Params and FLOPs, the training time is not greatly increased. In addition to the shorter training time,

it is advantageous to obtain a higher accuracy. Based on the above accuracy indices, the proposed method can effectively extract high-resolution RS images. Compared with other methods, not only can it achieve the highest accuracy, such as Kappa, but it can also effectively solve the problem of training time, which shows that the method in this study performs well.

In order to verify the superiority and generalization ability of the algorithm in this paper in road extraction, the public dataset Massachusetts Roads road dataset is used for experimental verification. The model of U-Net, SegNet, MDeepLab V3+ and this paper are respectively trained. To ensure comparability and validity of results, the training and test data sets are the same. The experimental results of U-Net, SegNet and MDeepLab V3+ networks and algorithms are shown in Figure 10. Where (a) is the image, (b) is the label map, (c) is the U-Net extraction result, (d) is the SegNet extraction result, (e) is the MDeepLab V3+ extraction result, and (f) is the extraction result of the proposed method.

As shown in Figure 10(a1), the road is clear and regular, but there are also trees blocking it. As shown in Figure 9(c1), U-Net extracts roads incompletely, and occluded roads cannot be fully extracted. As shown in Figure 9(d1), the extraction results of SegNet are slightly better than U-Net. However, there are still discontinuities in the extraction results, and the detail processing is slightly poor. Figure 9(e1), The extraction results of MDeepLab V3+ are better in continuity than U-Net and SegNet. However, it is also affected

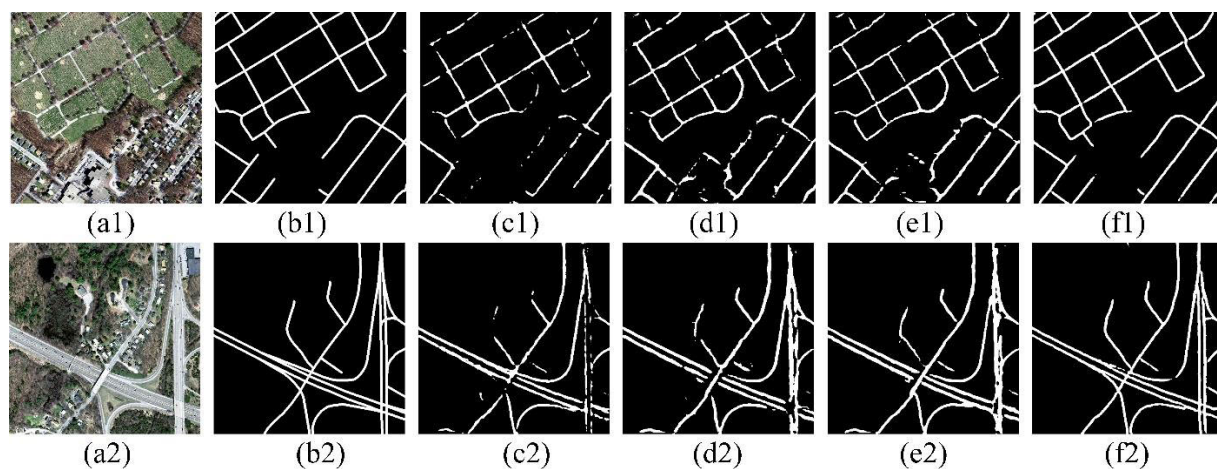


FIGURE 10. Compare the experimental results. (a) Image map; (b) Label map; (c) U-Net extraction results; (d) SegNet extraction results; (e) MDeepLab V3+ extraction results; (f) Extraction results of this paper.

TABLE 4. Comparison of experimental accuracy verification results.

Method	IoU/%	Recall/%	OA/%	Kappa	Params/MB	FLOPs/G	Training time/t
U-Net	74.95	64.27	95.09	0.6842	7.77	30.51	16.20
SegNet	76.88	64.88	96.13	0.7125	30.5	75.20	18.47
MDeepLabv 3+	77.25	83.12	96.83	0.7200	2.76	10.517	12.24
Ours	90.65	90.11	98.29	0.8987	2.85	12.76	12.30

by other ground features, resulting in inaccurate extraction. As shown in Figure 9 (f1), the method proposed in this paper can effectively extract roads. The edge details are well processed, and the extraction is complete and continuous. As shown in Figure 9 (a2), the road is irregular and the road structure is relatively complex. But the road is not blocked by buildings or trees etc. From Figure 9 (c2), the discontinuity of the extraction results of the U-Net network is serious, and the obvious road details are not accurately extracted. In Figure 9 (d2), SegNet showed obvious improvement compared with U-Net extraction results, but the above problems still existed. As can be seen from Figure 9 (e2), the extraction results of MDeepLab V3+ are improved compared with the previous two methods, but the extraction results are still inaccurate. After the method of this study integrates multi-scale and attention mechanism, the extraction effect is better than the extraction results of SegNet, U-Net and MDeepLab V3+. As shown in Fig. 9 (f2), the proposed method handles the edge details well and the extraction results are continuous. The extraction problems that existed with the other three methods are solved. The extraction results are basically consistent with the labeled samples.

It can be seen from Table 4 that the method proposed in this paper is optimal in all indicators. Among them, compared with U-Net, IoU increased by 15.7%, Recall increased by 25.84%, OA increased by 3.2% and kappa increased by 0.2145. Compared with SegNet, IoU increased

by 13.77%, Recall increased by 25.23%, OA increased by 2.16%, and Kappa increased by 0.1862. Compared with MDeepLab V3+, all indicators have also been greatly improved. The performance of the proposed method in Params and FLOPs is much lower than that of the U-Net and SegNet methods, which effectively reduces the training time. Compared with the MDeep Lab V3+ network, our method adds an attention mechanism. Although there is a slight improvement in the Params and FLOPs indicators, the training time does not increase to a large extent. Based on the above accuracy indicators, the method proposed in this paper can effectively extract high-resolution remote sensing images. Compared with other methods, it can not only achieve the highest accuracy, but also effectively solve the problem of training time. And the method in this paper has certain generalization ability.

C. ABLATION EXPERIMENT

To verify the effectiveness of the module in this study, an algorithm and its variants were used for experimental verification. MDeepLab V3+ is a DeepLab V3+ network with MobileNet V2 as the backbone network, and CDeepLab V3+ is an MDeepLab V3+ network that fuses channel attention. The experimental results are shown in Figure 11, where (a) is the image, (b) is the label, (c) is the extraction result of MDeepLab v3+, (d) is the extraction result of CDeepLab V3+, and (e) is the algorithm used in this study.

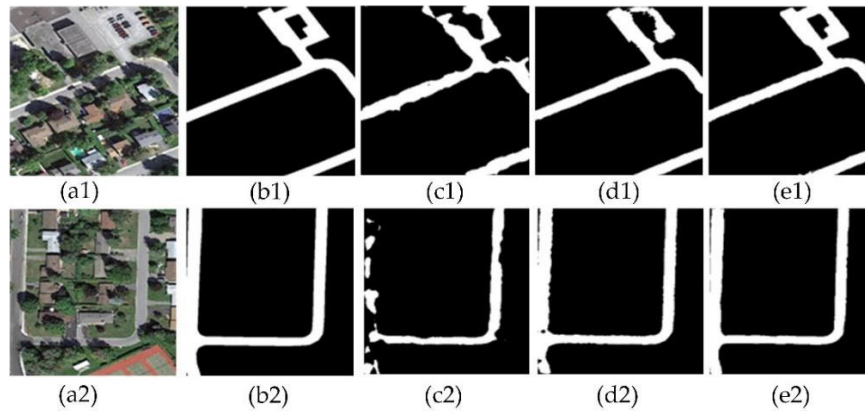


FIGURE 11. Ablation experiment results. (a) Image map; (b) Label map; (c) MDeepLab V3+ extraction results; (d) CDeepLab V3+ extraction results; (e) Extraction results of this paper.

TABLE 5. Accuracy verification results of ablation experiments.

Method	IoU/%	Recall/%	OA/%	Kappa	Params/M	FLOPs/G	Training time /t
MDeepLab V3+	80.21	81.80	97.23	0.8776	2.76	10.517	10.06
CDeepLab V3+	88.20	90.31	98.50	0.9288	2.83	11.74	10.08
Ours	91.63	95.07	98.92	0.9502	2.85	12.76	10.08

As shown in Figure 11 (c1), MDeepLab V3+ is not accurate for extraction around the parking lot. At the bend of the road, there is an inaccurate extraction problem caused by tree occlusion. The edge detail processing of the overall road is not accurate, and there is a loss of detail. Compared to the MDeepLab V3+ network, the CDeepLab V3+ network adds a channel attention mechanism. The problem in Figure 11(d1) is improved compared to that in Figure 11 (c1). The edge details are better and more accurate, but there are still inaccuracies in the extraction around the parking lot. In Figure 11 (e1), the proposed algorithm integrates the spatial attention and channel attention mechanisms, and the extraction results of the proposed algorithm are more accurate than those of the MDeepLab V3+ and CDeepLab V3+ networks. As shown in Figure 11 (c2), the extraction results of the MDeepLab V3+ network exhibit serious discontinuity problems, and the processing of edge details is poor. It can be seen from Figure 11 (d2) that the extraction results of CDeepLab V3+ are compared to those in Figure 11 (c2), and the discontinuity problem has been effectively improved, but there are holes. As can be seen in Figure 11 (e2), when the MDeepLab V3+ network simultaneously increases channel attention and spatial attention at the same time, the problem of extracting discontinuity and edge loss is optimized. The method presented in this paper obtained a good extraction effect. Based on the above, the modules of the algorithm in this study are effective, and the method in this study is more applicable.

As can be seen in Table 5, after MDeepLab V3+ increased the channel attention, IoU, Recall, OA and Kappa

increased significantly. Compared with MDeepLab V3+, the CDeepLab V3+ network has an 11.42% increase in IoU and a 13.27% increase in recall. OA increased by 1.69%, and Kappa reached 0.9288. When spatial attention and channel attention are increased at the same time, the IoU value of the proposed method reaches 91.63%. OA reached 98.92%, and Kappa reached 0.9502. It was shown that channel attention and spatial attention can effectively solve the problems of road extraction discontinuity and edge detail loss. However, there was not much improvement in the two individual factors of Params and FLOPs, resulting in no change in training time. The above indicators show that the improved part of the method in this study can effectively solve the problems of discontinuous extraction and loss of edge details. The algorithm in this study can effectively reduce training time and achieve accuracy and efficiency.

V. CONCLUSION

As important geo-information for a country, roads play an important role in the transportation and insurance industries. However, in the process of extracting road information from high-resolution RS images, problems such as incompleteness and inaccuracy of road information have become problems that need to be solved. This paper aims at the low precision of the traditional high-resolution RS image road extraction method, the post-processing requires a lot of manual participation. The FCN network method has problems, such as many parameters, time consumption, discontinuous extraction, and loss of edge details. A lightweight road extraction method that integrates multiscale and attention mechanisms is proposed.

Solve the problems of discontinuous extraction and loss of edge details through the ASPP and attention mechanism. By using the MobileNet V2 network instead of Xception as the backbone network of DeepLab V3+, the problem of a long training time is solved. The Ottawa Road dataset was verified. In the comparison test of traditional methods, SVM and RF were not suitable for road extraction, and they were easily affected by occluding objects. Compared with the two traditional methods, the method proposed in this paper is superior in terms of extraction. In the comparative experiment on the Ottawa road dataset, although the U-Net and SegNet networks are not affected by occlusions in the process of road extraction, there are discontinuities, edge detail losses, holes, and other phenomena in the road extraction results. MDeepLab V3+ integrates an ASPP module. This solves the phenomenon of holes in the extraction process, but there are still problems such as discontinuous extraction and loss of edge details. The method presented in this paper can effectively solve the above problems, and the road extraction problem can be effectively solved. Compared with the other three comparison methods, this study not only has a greater advantage in accuracy but also can save a lot of training time, effectively solving the problem of lengthy road extraction. In the comparison experiment on the Massachusetts road dataset, the extraction results of U-Net and SegNet both have discontinuous extraction. Although the extraction results of MDeepLab v3+ have been improved compared to the previous two methods, this phenomenon still exists. The method in this paper can effectively extract roads. It shows that the method in this paper has certain superiority and generalization ability. In addition, ablation experiments verified that the improved part in this study is effective. It can be observed from the above that the proposed method can effectively extract roads from high-resolution RS images and has advantages over the comparison method. The time-consuming problem is solved, which shows that the method in this study has a certain applicability. However, it can be seen from the verification of different data sets in this paper that the accuracy of the experimental results on the Ottawa road data set is better than that on the Massachusetts road data set. It shows that the method in this paper is more suitable for the research of road extraction from high-resolution remote sensing images. Although the method proposed in this paper has improved the number of parameters in the model, the long training time is still a problem faced by FCN networks. In future studies, lightweight networks can be used for road extraction. High-resolution RS images require a large amount of storage space, which is required in large-scale experiments. In future research, a suitable resolution RS image will be selected for large-scale road extraction research. The method in this study was only researched for the extraction of roads. However, it is not known whether it can effectively extract other ground features. In future research, this method will be applied to more target-extraction tasks.

ACKNOWLEDGMENT

The authors would like to thank Wei Yue and Xiangyuan Ding from the Chinese Academy of Forestry for their help in data processing. The authors would also like to thank the editors and referees for their constructive criticism of this study.

REFERENCES

- [1] G. Cheng, F. Zhu, S. Xiang, and C. Pan, "Road centerline extraction via semisupervised segmentation and multidirection nonmaximum suppression," *IEEE Geosci. Remote Sens. Lett.*, vol. 13, no. 4, pp. 545–549, Apr. 2016.
- [2] Y. H. Zhang, J. He, X. Kan, and J. Y. Li, "Summary of road extraction methods for RS images," *Comput. Eng. Appl.*, vol. 54, no. 13, pp. 1–10, 2018.
- [3] Y. Liu, J. Kang, H. Y. Guan, and H. Y. Wang, "High-resolution RS image road extraction model based on dual-attention residual network," *J. Geo-Inf. Sci.*, vol. 25, no. 2, pp. 396–408, 2023.
- [4] X. Lin, J. Zhang, Z. Liu, and J. Shen, "Semi-automatic road tracking by template matching and distance transform," in *Proc. Joint Urban Remote Sens. Event*, May 2009, pp. 1–7.
- [5] S. Udomhunsakul, "Semi-automatic road extraction from aerial images," *Proc. SPIE*, vol. 5239, pp. 26–32, Mar. 2004.
- [6] H. Mayer, I. Laptev, A. Baumgartner, and C. Steger, "Automatic road extraction based on multi-scale modeling, context, and snakes," *Int. Arch. Photogramm. Remote Sens.*, vol. 32, no. 3, pp. 106–113, 1997.
- [7] A. Baumgartner, C. Steger, H. Mayer, and H. Ebner, "Automatic road extraction based on multi-scale, grouping, and context," *Photogramm. Eng. Remote Sens.*, vol. 65, no. 7, pp. 777–785, 1999.
- [8] K. Treash and K. Amaratunga, "Automatic road detection in grayscale aerial images," *J. Comput. Civil Eng.*, vol. 14, no. 1, pp. 60–69, Jan. 2000.
- [9] R. Gaetano, J. Zerubia, G. Scarpa, and G. Poggi, "Morphological road segmentation in urban areas from high resolution satellite images," in *Proc. 17th Int. Conf. Digit. Signal Process. (DSP)*, Jul. 2011, pp. 1–8.
- [10] L. F. Wang and C. M. Yan, "A review of road scene semantic segmentation?" *Prog. Laser Optoelectron.*, vol. 58, no. 12, 2021, Art. no. 1200002.
- [11] R. Alshehhi and P. R. Marpu, "Hierarchical graph-based segmentation for extracting road networks from high-resolution satellite images," *ISPRS J. Photogramm. Remote Sens.*, vol. 126, pp. 245–260, Apr. 2017.
- [12] M. Song and D. Civco, "Road extraction using SVM and image segmentation," *Photogrammetric Eng. Remote Sens.*, vol. 70, no. 12, pp. 1365–1371, Dec. 2004.
- [13] L. Xiao, B. Dai, D. Liu, T. Hu, and T. Wu, "CRF based road detection with multi-sensor fusion," in *Proc. IEEE Intell. Vehicles Symp. (IV)*, Jun. 2015, pp. 192–198.
- [14] J. Dai, Y. Wang, Y. Du, T. Zhu, S. Xie, C. Li, and X. Fang, "Development and prospect of road extraction method for optical remote sensing image," *Nat. Remote Sens. Bull.*, vol. 24, no. 7, pp. 804–823, 2020.
- [15] M. Guo, H. Liu, Y. Xu, and Y. Huang, "Building extraction based on U-Net with an attention block and multiple losses," *Remote Sens.*, vol. 12, no. 9, p. 1400, Apr. 2020.
- [16] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 3431–3440.
- [17] X. W. Kong, C. Y. Wang, S. C. Zhang, J. H. Li, and Y. Sui, "Application of improved U-Net network in road extraction from RS image remote sensing information," vol. 37, no. 2, pp. 97–104, 2022.
- [18] J. Li, Y. Liu, Y. Zhang, and Y. Zhang, "Cascaded attention DenseUNet (CADUNet) for road extraction from very-high-resolution images," *ISPRS Int. J. Geo-Inf.*, vol. 10, no. 5, p. 329, May 2021.
- [19] L. A. R. Umama, J. E. M. Baquero, and R. J. Moreno, "Semantic segmentation for applications in autonomous vehicles," *Int. J. Appl. Eng. Res.*, vol. 13, no. 5, pp. 1–12, 2018.
- [20] L. C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected CRFs," *Comput. Sci.*, vol. 2014, no. 4, pp. 357–361, 2015.
- [21] L. C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking Atrous convolution for semantic image segmentation," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2017, pp. 1–14.

- [22] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, Atrous convolution, and fully connected CRFs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, Apr. 2018.
- [23] L. C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with Atrous separable convolution for semantic image SEG mentation," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 801–818.
- [24] L. L. Han, Z. Yang, H. L. Li, Z. X. Liu, and B. W. Huang, "Road extraction of high resolution RS imagery based on DeepLab V3," *Remote Sens. Inf.*, vol. 36, no. 1, pp. 22–28, 2021.
- [25] R. Liu and D. He, "Semantic segmentation based on Deeplabv3+ and attention mechanism," in *Proc. IEEE 4th Adv. Inf. Manage., Communicates, Electron. Autom. Control Conf. (IMCEC)*, vol. 4, Jun. 2021, pp. 255–259.
- [26] D. Y. Liu, H. X. Jiang, N. N. Rao, C. S. Luo, W. J. Du, Z. W. Li, and T. Gan, "Computer aided annotation of early esophageal cancer in gastroscopic images based on Deeplabv3+ network," in *Proc. 4th Int. Conf. Biomed. Signal Image Process.*, 2019.
- [27] H. Su, Y. Peng, C. Xu, A. Feng, and T. Liu, "Using improved DeepLabv3+ network integrated with normalized difference water index to extract water bodies in Sentinel-2A urban remote sensing images," *J. Appl. Remote Sens.*, vol. 15, no. 1, Mar. 2021.
- [28] W. Yang, J. L. Zhang, Z. Y. Xu, and K. Hu, "Real-time DeepLabv3+ for pedestrian segmentation," *J. Opt. Technol.*, vol. 86, no. 9, p. 570, Sep. 2019.
- [29] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient convolutional neural networks for mobile vision applications," 2017, *arXiv:1704.04861*.
- [30] S. Huang, Y. He, and X.-A. Chen, "M-YOLO: A nighttime vehicle detection method combining MobileNet V2 and YOLO v3," *J. Phys., Conf. Ser.*, vol. 1883, no. 1, Apr. 2021, Art. no. 012094.
- [31] J. Li, H. Zhao, S. P. Zhu, H. Huang, Y. Miao, and Z. Jiang, "An improved lightweight network architecture for identifying tobacco leaf maturity based on deep learning," *J. Intell. Fuzzy Syst.*, vol. 41, no. 2, pp. 4149–4158, Sep. 2021.
- [32] P. N. Huu, H. N. T. Thu, and Q. T. Minh, "Proposing a recognition system of gestures using MobileNetV2 combining single shot detector network for smart-home applications," *J. Electr. Comput. Eng.*, vol. 2021, pp. 1–18, Feb. 2021.
- [33] S.-Y. Lu, S.-H. Wang, and Y.-D. Zhang, "A classification method for brain MRI via MobileNet and feedforward network with random weights," *Pattern Recognit. Lett.*, vol. 140, pp. 252–260, Dec. 2020.
- [34] Y. Zhu, "ASPP-DF-PVNet: Atrous spatial pyramid pooling and distance-filtered pvnet for occlusion resistant 6D estimation," *Signal Process., Image Commun.*, vol. 95, no. 1, 2021.
- [35] Y. Zhu, L. Wan, W. Xu, and S. Wang, "ASPP-DF-PVNet: Atrous spatial pyramid pooling and distance-filtered PVNet for occlusion resistant 6D object pose estimation," *Signal Process., Image Commun.*, vol. 95, Jul. 2021, Art. no. 116268.
- [36] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, "Squeeze-and-excitation networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 8, pp. 2011–2023, Aug. 2020.
- [37] J. Li, Y. Tian, and T. Lee, "Convolution-based channel-frequency attention for text-independent speaker verification," 2022, *arXiv:2210.17310*.
- [38] J.-H. Park, S.-M. Seo, and J. H. Yoo, "Channel attention module in convolutional neural network and its application to SAR target recognition under limited angular diversity condition," *J. Korea Inst. Mil. Sci. Technol.*, vol. 24, no. 2, pp. 175–186, Apr. 2021.
- [39] X. Z. Zhu, D. Z. Cheng, Z. Zhang, S. Lin, and J. Dai, "An empirical study of spatial attention mechanisms in deep networks," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Seoul, South Korea, Oct. 2020, pp. 6687–6696.
- [40] Z. Li, Y. Huang, M. Cai, and Y. Sato, "Manipulation-skill assessment from videos with spatial attention network," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. Workshop (ICCVW)*, Oct. 2019, pp. 4385–4395.
- [41] Z. Cui, Q. Li, Z. Cao, and N. Liu, "Dense attention pyramid networks for multi-scale ship detection in SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 11, pp. 8983–8997, Nov. 2019.
- [42] H. Zhang, A. C. Berg, M. Maire, and J. Malik, "SVM-KNN: Discriminative nearest neighbor classification for visual category recognition," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jul. 2006, pp. 1–12.
- [43] A. Liaw and M. Wiener, "Classification and regression by random forest," *R News*, vol. 23, no. 23, 2002.
- [44] S. Heremans and J. Van Orshoven, "Machine learning methods for sub-pixel land-cover classification in the spatially heterogeneous region of Flanders (Belgium): A multi-criteria comparison," *Int. J. Remote Sens.*, vol. 36, no. 11, pp. 2934–2962, Jun. 2015.



RUI WANG received the degree from Wuhan University. He is currently with the China Transport Telecommunications and Information Center, Beijing, China. His research interests include the application of traffic data in the insurance industry, spatial big data management, and the study of image recognition algorithms based on machine learning.



MINGXIANG CAI received the degree from Wuhan University. He is currently with the China Transport Telecommunications Information Center, Beijing, China. His research interests include spatio-temporal information mining and the application of fundamental models in traffic big data.



ZIXUAN XIA is currently pursuing the degree in surveying and mapping engineering with the Heilongjiang University of Technology. His current research interests include flight calibration and validation of forestry applications of high-definition aerial systems and the study of vegetation cover differentiation based on the image dichotomous model.

...