# DPIF-Net: A Dual Path Network for Rural Road Extraction Based on the Fusion of Global and Local Information

Yuan Sun[1], Xingfa Gu[1], Xiang Zhou[1], Jian Yang[1], Wangyao Shen[2], Yuanlei Cheng[2], Jinming Zhang[2], Yunping Chen[2]

[1] Aerospace Information Research Institute, Chinese Academy of Sciences, Beijing, China

[2] School of Automation Engineering, University of Electronic Science and Technology of China, Chengdu, Sichuan, China

Corresponding Author:

Xingfa Gu[1]

No.9 Dengzhuang South Road, Haidian District, Beijing, 100094, China

Email address: xfgu@aircas.ac.cn

Yunping Chen[2]

No.2006, Xiyuan Ave., West High-tech Zone, Chengdu, 611731, China

Email address: chenyp@uestc.edu.cn

## Abstract

**Background**: Automatic extraction of roads from remote sensing images can facilitate many practical applications. However, thus far, thousands of kilometers or more of roads worldwide have not been recorded, especially low-grade roads in rural areas. Moreover, rural roads have different shapes and are influenced by complex environments and other interference factors, which has led to a scarcity of dedicated low level category road datasets.

**Methods**: To address these issues, based on convolutional neural networks (CNNs) and tranformers, this paper proposes the Dual Path Information Fusion Network (DPIF-Net). In addition, given the severe lack of low-grade road datasets, we constructed the GaoFen-2 (GF-2) rural road dataset to address this challenge, which spans three regions in China and covers an area of over 2300 kilometers, almost entirely composed of low-grade roads. To comprehensively test the low-grade road extraction performance and generalization ability of the model, comparative experiments are carried out on the DeepGlobe, and Massachusetts regular road datasets.

**Results**: The results show that DPIF-Net achieves the highest IoU and $F_1$ score on three datasets compared with methods such as U-Net, SegNet, DeepLabv3+, and D-LinkNet, with notable performance on the GF-2 dataset, reaching 0.6104 and 0.7608, respectively. Furthermore, multiple validation experiments demonstrate that DPIF-Net effectively preserves improved connectivity in low-grade road extraction with a modest parameter count of 63.9 MB. The constructed low-grade road dataset and proposed methods will facilitate further research on rural

40　roads, which holds promise for assisting governmental authorities in making informed decisions
41　and strategies to enhance rural road infrastructure.
42　**Keywords**: rural road extraction, remote sensing images, convolutional neural networks,
43　transformer
44

45　## Introduction

46　Roads are typical landscape features with complex topological relationships, and millions of
47　kilometers of roads in the world are still unrecorded, particularly low-grade roads in rural areas.
48　In *China (2003)*, low-grade roads are defined as those with an annual average daily traffic
49　volume of fewer than 6,000 cars. These roads are vital for promoting urban–rural economic
50　exchange and narrowing the gap between urban and rural areas. Therefore, it is imperative to
51　design an intelligent and automatic method for rural road extraction.
52　　　Although high-resolution remote sensing images have been studied for many years for road
53　extraction, accurately identifying rural roads in these images may face additional challenges in
54　reality (*Li et al. 2021*). High-resolution images provide rich discriminative features for road
55　identification but also contain many interference factors. For instance, shadows on roads are
56　produced due to occlusions from various vehicles, trees, and tall buildings under different
57　illumination conditions. Rural roads are often characterized by an absence of distinct geometric
58　features, and their connectivity can be impacted by nearby rivers, which in turn affects the
59　effectiveness of road extraction. In addition, rural roads made of dirt are more difficult to extract
60　than roads made of asphalt or cement.
61　　To solve these problems, this paper proposes an end-to-end network called the Dual Path
62　Information Fusion Network (DPIF-Net), which combines the strengths of convolutional neural
63　networks (CNNs) and transformers to further improve the accuracy of rural road extraction.
64　Furthermore, since there are few datasets related to rural roads, a dataset of rural roads is
65　specifically constructed. Finally, we present extensive experiments conducted on the DeepGlobe
66　and Massachusetts datasets as well as our dataset to test the model's generalization ability and
67　robustness.
68　　　The contributions of this paper are summarized as follows:
69　　　(1) The proposed DPIF-Net, which has a small number of network parameters and a simple
70　structure. It effectively combines the advantages of CNNs in spatial induction with the adaptive
71　weighting of input information in transformers to establish global dependencies. Moreover,
72　DPIF-Net can effectively extract both the local detailed features and global context features of
73　roads and fully integrate this information to produce more accurate road segmentation results.
74　　　(2) The constructed dataset of rural roads. Our dataset includes roads of different regions in
75　China, but most of them are various types of rural roads. This dataset is specifically constructed
76　for studying rural road extraction and our model's performance on rural roads.
77　　　The rest of this paper is organized as follows. Section 2 describes some related work on
78　deep learning for road extraction. Section 3 explains the road dataset used in the experiments and

79  describes the details of the method proposed in this paper. Section 4 presents the experimental
80  results and analysis. Sections 5 and 6 provide a discussion and conclusions.
81
82  **Related works**
83  At present, road extraction and monitoring operations are still performed manually or
84  semimanually, making them ineffective and costly (*Abdollahi et al. 2020*). Therefore, new robust
85  techniques, such as deep learning methods, are needed to accurately extract road networks of
86  various scales from remote sensing imagery (*Panboonyuen et al. 2017*), which has gradually
87  become a prominent direction of research. With the development of artificial intelligence, deep
88  convolutional neural networks (DCNNs) are gradually gaining dominance in the field of image
89  processing. In recent years, there has been an explosion of various papers on road segmentation
90  with DCNNs, and many excellent CNN models have emerged, such as U-Net (Ronneberger,
91  Fischer & Brox 2015), LinkNet (*Chaurasia & Culurciello 2017*), SegNet *(Badrinarayanan,*
92  *Kendall & Cipolla 2017*), D-LinkNet (*Zhou, Zhang & Wu 2018*), DeepLabv3+ (*Chen et al.*
93  *2018b*), and generative adversarial networks (GANs) (*Goodfellow et al. 2020*). These models
94  integrate features from multiple layers of a CNN to exploit the multiscale information at different
95  semantic levels (*Zhu et al. 2021*). Many road segmentation methods are based on the above
96  models.
97      *Zhang, Liu & Wang (2018)* integrated residual units into a U-Net-like network for road
98  extraction. Residual units can make it easier for a network to learn features and achieve better
99  results. *Moradi et al. (2019)* proposed a modified U-Net architecture combined with a feature
100  pyramid network and concatenated the feature maps from all levels of the U-Net decoder path as
101  input. Their method achieved good performance in medical image segmentation. *Chen et al.*
102  *(2021)* proposed a reconstruction bias U-Net for road extraction from remote sensing images.
103  This method obtains multiple levels of semantic information from different upsampling scales by
104  adding decoding branches. However, the extraction effect of the modified method is not good for
105  low-grade roads, such as rural roads. *Yang et al. (2019)* constructed a U-Net network consisting
106  entirely of Region CNN (RCNN) blocks, which preserve rich low-level spatial features. Inspired
107  by U-Net and atrous spatial pyramid pooling (ASPP) (*Chen et al. 2018a*), *He et al. (2019)*
108  integrated an ASPP module into U-Net to obtain multiscale road information. *Lu et al. (2019)*
109  proposed a deep learning framework based on U-Net, which can extract roads and road
110  centerlines, and integrate feature information from different scales to improve the robustness of
111  the model. *He et al. (2019)* added ASPP between the encoder and decoder in U-Net. At the same
112  time, a loss function that considered the digital number (DN) value, contrast, structure and other
113  factors of the image was proposed. *Lu et al. (2019)* replaced the first convolutional layer of each
114  group in U-Net with a multiscale module and constructed a pyramid-like structure to complete
115  the extraction of roads and road centerlines. To capture more information, a weighted loss for
116  roads and centerlines was built. Each loss component was weighted in accordance with the
117  relative proportions of background and target to solve the problem of target class imbalance.

118    Based on LinkNet, *Wang, Seo & Jeon (2022)* proposed an efficient nonlocal LinkNet with
119    nonlocal blocks (NLBs) that can grasp relations between global features. This enables each
120    spatial feature point to refer to all other contextual information and results in more accurate road
121    segmentation. *Zhu et al. (2021)* added an attentive GCA block between the encoder and decoder
122    to make the extracted road information more complete. They used FRN normalization to
123    improve the robustness of the model. *Xie et al. (2019)* replaced the D-LinkNet intermediate
124    structure with a global perception block for higher-order information. The design of the high-
125    order information global perception block was inspired by bilinear pooling. Experiments showed
126    that it achieved better performance than atrous convolution and could reduce the number of
127    parameters by 1/4. *Zhu et al. (2020)* proposed a model based on D-LinkNet and conditional
128    random fields (CRFs) to solve the edge smoothing problem in the process of building extraction.
129    *Tao et al. (2019)* proposed a network composed of a spatial information inference structure
130    (SIIS) for road extraction, and the overall framework was based on DeepLabv3+. The SIIS
131    consisted of two groups of RCNN units. A weighted loss function combining the mean squared
132    error (MSE) and intersection over union (IoU) was adopted. To solve the problem of imbalanced
133    samples, images with excessively small target proportions relative to the background were
134    removed. *Lourenco et al. (2023)* proposed combining DeepLabv3+ with an optimization strategy
135    to extract rural roads.
136    Many road extraction methods based on Generative Adversarial Network (GAN) have
137    achieved impressive results. *Zhang et al. (2019b)* proposed a GAN for road extraction that had
138    multiple discriminators. In the experiments, it was found that a combination of 4 discriminators
139    and 1 generator was best. At the same time, a road label generation method that needed less
140    manual intervention was proposed. *Shamsolmoali et al. (2021)* integrated feature pyramids into
141    GANs for road detection. *Zhang et al. (2019a)* explored different types of GANs. An end-to-end
142    model for road extraction based on GANs was proposed. The influence of convolution kernels of
143    different sizes was discussed, and it was concluded that large convolution kernels were not
144    needed to improve the receptive field for road extraction.
145    In addition to the above models, some scholars have used other road extraction methods and
146    have also achieved promising results. *Bastani et al. (2018)* presented a method to extract road
147    networks based on iterative graph construction. The final road map was generated by iteratively
148    adding new candidate road regions. A decision function was used to determine whether a
149    candidate area is a road by training a CNN. However, this method requires knowledge of the
150    initial points of the roads. *Shao et al. (2021)* proposed a two-task end-to-end CNN named the
151    Multitask Road-related Extraction Network (MRENet) for road surface extraction and road
152    centerline extraction. The network design of MRENet uses atrous convolutions and a pyramid
153    scene parsing pooling (PSP pooling) module to expand the network's receptive field, integrate
154    multilevel features, and obtain more abundant information. In addition, the authors used a
155    weighted binary cross-entropy function to alleviate the background imbalance problem. *Zhang &*
156    *Wang (2019)* introduced a network consisting of dense cavity convolution modules for road and
157    building extraction. *Batra et al. (2019)* proposed joint learning based on orientation and

158  segmentation maps to enhance the connectivity rate in road extraction. The CNN-based structure
159  achieved good road extraction results, but the accuracy was not high for complex road networks,
160  and the method was not effective for low-grade roads.
161      The transformer model has made a vital difference in the natural language processing (NLP)
162  field because of its attention mechanism (*Vaswani et al. 2017*). Inspired by the powerful
163  representation capabilities of transformers, researchers have extended transformers to computer
164  vision tasks (Han et al. 2020). Compared with other networks, transformer-based networks can
165  achieve comparable performance with less computation. *Dosovitskiy et al. (2020)* built a
166  framework consisting of a pure transformer for image classification tasks. The architecture was
167  trained using large-scale data to obtain pretrained models. When transferred to vision tasks, it
168  achieved a performance comparable to that of CNNs. *Xie et al. (2021)* combined a fully
169  convolutional network with an attention mechanism to learn information from long-range
170  contexts and achieved good results in image segmentation tasks. In addition, *Xie et al. (2021)*
171  built a semantic segmentation framework combining transformer and Multi-Layer Perceptron
172  (MLP). The framework was simple, efficient and powerful and consisted of a hierarchical
173  transformer encoder and a decoder composed of MLPs. It could output multiscale features and
174  did not require positional encoding, resulting in significantly improved performance and
175  efficiency compared with similar algorithms. Therefore, a structure based on the transformer
176  demonstrates a clear advantage in global feature extraction. In summary, current research is
177  predominantly based on CNN, which has shown good performance in typical road extraction.
178  However, the accuracy of this approach tends to diminish in complex road networks. Given the
179  complexity of shape for low-level roads, this paper aims to explore the potential of combining
180  CNN and transformer architectures specifically for low-level road contexts.
181

## Materials & Methods

### Dataset

184  Although the currently available public road datasets cover a wide range of road categories in
185  cities, suburbs and rural areas in many countries worldwide, they contain many normal roads and
186  few rural roads. Therefore, they are not suitable for analysis with a special focus on rural roads,
187  but they can be used as test data for model generalization performance.
188      In this study, a rural road dataset was constructed based on the GaoFen-2 (GF-2) satellite.
189  The GF-2 satellite carries a range of sensors, including a Panchromatic and Multispectral sensor
190  (PMS), a wide-field-of-view sensor (WFV), and a hyperspectral sensor (HSI), which provide
191  high-resolution imagery with spatial resolutions ranging from 0.8 m to 16 m. The images for the
192  PMS sensor (450 to 900 nm) at 0.8 m were utilized in this paper. The images include three
193  regions covering an area of over 2300 square kilometers : the junction between Jiancaoping
194  District and Gujiao City in Shanxi characterized by imagery measuring $36500 \times 34258$ pixels,
195  covering an approximate area of 752 square kilometers,  the junction between Anyang and
196  Shijiazhuang cities in Hebei characterized by imagery measuring $39695 \times 31311$ pixels, covering
197  an approximate area of 795 square kilometers, and the junction area of Guangzhou and Foshan in

**Deleted:** details

**Commented [sst1]:** Province?

199　Guangdong characterized by imagery measuring 36020 × 32431 pixels, covering an approximate
200　area of 747 square kilometers. As an example, the study area at the junction between Taiyuan
201　and Jiancaoping District and its mask samples, are illustrated in Figure 1. All GF-2 data in our
202　study is sourced from the China Centre for Resources Satellite Data and Application.
203
204　**Figure.1 Partial study area schematic diagram and masks.**
205
206　　　　Nevertheless, the considerable size of the satellite images poses challenges in terms of
207　efficient data loading, potentially leading to a substantial increase in training duration. Moreover,
208　the absence of masks within the original dataset introduces complexities in conducting effective
209　supervised training. To overcome these challenges, the images underwent cropping to achieve
210　dimensions of 512x512 pixels initially. Subsequently, manual annotation based on image texture
211　was executed to generate corresponding masks. In the end, we generated a dataset similar to the
212　examples shown in Figure 2. To address the potential sample imbalance, our dataset excluded the
213　images containing abundant higher level roads, consequently encompassing intricate details of
214　rural road attributes, such as tree coverage, agricultural field irrigation channels, and road
215　incompleteness. These adjustments were intended to enhance the model's performance in
216　extracting low-grade rural roads. After preprocessing the data, 5501 samples remained, with
217　5421 as training samples, 40 as validation samples, and 40 as test samples.
218
219　**Figure. 2. Overview of the data.** (All images and masks from the Massachusetts dataset)
220
221　　　　In addition, experiments were carried out on two public datasets, DeepGlobe *(Demir et al.
222　2018)* and Massachusetts *(Mnih 2013)*. However, it is imperative to emphasize that these two
223　datasets contain a limited quantity of low-grade roads in comparison to a substantial volume of
224　regular highways. They are specifically employed to enhance our model's extraction performance
225　and validate findings. The images in the DeepGlobe road dataset come from three countries,
226　namely, India, Thailand, and Indonesia, and include multiple imaged scenes covering an area of
227　over 2220 square kilometers, such as cities, villages, wilderness, suburbs, seashores, and tropical
228　rainforests. The ground resolution of the images is 0.5 m per pixel, and the image size is 1024 ×
229　1024 pixels. There are a total of 6226 images, of which 4976 are designated for training and
230　1250 are designated for testing. In this study, following *Zhu et al. (2021)*, the original images
231　were cropped to a resolution of 512 × 512 pixels with an overlap of 256 pixels. Finally, a total of
232　5000 images for training, 40 images for validation and 4500 images for testing were obtained.
233　　　　The Massachusetts road dataset consists of 1171 aerial images of the Massachusetts region,
234　which cover a wide variety of urban, suburban, and rural regions and an area of over 2600 square
235　kilometers. With a spatial resolution of 1 m per pixel, the images in this dataset have a size of
236　1500 × 1500 pixels and are composed of red, green, and blue channels. Similarly, the original
237　data were cropped into nonoverlapping images with a resolution of 512 × 512 pixels. Finally,

Deleted: For

Deleted: as

Deleted: .

241 3744 images were used for training, 30 images were used for validation, and 196 images were
242 used for testing. Moreover, all images with large white blank areas were removed manually.
243
244 **Details of the model structure**
245     We propose Dual Path Information Fusion Network (DPIF-Net) to improve the
246 performance of rural road extraction by exploring the potential of combining the capabilities of
247 transformers and CNNs for road segmentation. The schematic structure of DPIF-Net is displayed
248 in Figure 3. First, the top encoder branch uses a transformer to model global road information in
249 the input remote sensing image, while the other encoder branch uses convolution operations to
250 extract local details of roads and process spatial and channel information. Second, the feature
251 information of the two branches is effectively fused. Finally, each layer of the decoder fuses
252 high-level features from the previous layer with low-level features from the convolutional branch
253 and gradually upsamples the image to the original resolution to obtain a binary image containing
254 only roads.
255
256 **Figure. 3. Overview of the proposed model for rural road extraction.** (All satellite images
257 and masks from the Massachusetts dataset)
258
259 **Local detail information encoder based on a CNN**
260 In DPIF-Net, we propose a convolution module called the local detail feature extraction (LDFE)
261 block as shown in Figure 4. This block is composed of three parts to efficiently extract road
262 features while keeping the number of network parameters low.
263     In part A, a traditional 3x3 convolution is applied to the input feature map to extract
264 preliminary feature information without altering its resolution, resulting in a feature map that is
265 four times larger than the input in the channel dimension. Then, the obtained features are spliced
266 and fused using a skip connection.
267     In part B, each channel of the feature map is processed separately using a 3x3 convolution
268 to extract semantic feature information, and a leaky-ReLU (*Maas, Hannun & Ng 2013*)activation
269 layer is used to obtain nonlinear features, as expressed in formula (1).

$$Leak - ReLU(x) = \begin{cases} x & (x > 0) \\ \alpha x & (x <= 0) \end{cases} \qquad (1)$$

271 $\alpha$ is a small gradient value, which is set to 0.2 in this paper.
272     In part C, the feature information from part B is combined, cross-channel interaction is
273 achieved using a 1x1 convolution, and then another leaky-ReLU activation layer is applied.
274
275 **Figure. 4. Local detail feature extraction block.**
276
277 **Global information encoder based on a Transformer module**
278 The ASPP module was originally introduced to increase the receptive field without sacrificing
279 too much resolution, allowing for the preservation of image details as much as possible(Chen et

al. 2018a). To process multiscale road feature information extraction, we propose improvements to the ASPP module. Specifically, there are three different dilated convolution modules and a global adaptive pooling module, which can extract more feature information, as shown in Figure 5. Moreover, different dilation ratios are used to obtain abstract feature information at different scales, which provides a basis for the subsequent integration of features and accurate road segmentation. In this way, the modified ASPP module obtains more abstract semantic information and road topology information and achieves better robustness; it can accept input feature images of any size and finally produce output of a fixed size.

**Figure. 5. Modified ASPP module.**

In addition to the improvements to the CNN module proposed above, a Trans block is proposed to employ the Transformer architecture at the beginning of the model to capture global context information while avoiding interference with the CNN branch. Specifically, the Trans block can compensate for the shortcomings of the CNN in capturing global context information.

To process image blocks with a transformer, we convert them into one-dimensional data by using a fully connected (FC) layer. Due to the expansion into one-dimensional data, the position information is lost; therefore, the input feature data are converted from index numbers into a one-hot encoding matrix to maintain the position relationships of the sequence. Moreover, a random weight matrix is right-multiplied to complete the input position embedding. This method ensures that each token in the sequence has a unique position representation, which is crucial for the transformer to capture long-range dependencies.

**Figure. 6. The structure of the transformer block.**

The schematic structure of an individual transformer block is shown in Figure 6 (Dosovitskiy et al. 2020). The Transformer encoder mainly consists of alternating multihead self-attention (MSA) layers and multilayer perceptron (MLP) blocks. Before these two modules, layer normalization (LN) is applied for normalization, and a residual connection structure is used. The MLP block consists of two FC layers, in which the Gaussian error linear unit (GeLU) function is applied for nonlinear activation to obtain nonlinear features. The definition of the GeLU (*Hendrycks & Gimpel 2016*) activation function is shown in formula (2).

$$GeLU(x) = xP(X \le x) = x \times \phi(x) \quad X \sim N(0,1) \qquad (2)$$

$x$ is the input, and X is a random variable that follows a Gaussian distribution with mean 0 and variance 1. $P(X \le x) \phi(x)$ is the cumulative distribution of the Gaussian normal distribution of $x$, for which there is no analytical expression; instead, its approximate calculation method is shown in formula (3).

$$GeLU(x) = \frac{1}{2} x(1 + \tanh(\sqrt{\frac{2}{\pi}}(x + 0.044715x^3)))\qquad (3)$$

In our work, the number of attention heads is set to 16, assuming that the input image is $x$ ($x \in i^{H \times W \times C}$), where H , W and C are the height, width and channels, respectively, of $x$. The original image block feature $x_p$ ($x_p \in i^{N \times P^2 \times C}$) is obtained by expansion into one-dimensional data, $x_p \in i^{N \times P^2 \times C}$, where N is the number of small P×P patches into which the original input image is cropped. N and P satisfy the relationship shown in formula (4).

$$N = \frac{H \times W}{P^2} \tag{4}$$

Then, after dimension embedding and nonlinear mapping of the GeLU function, the i-th image block feature $z^i$ can be expressed as shown in formula (5).

$$z^i = GeLU(x_p^i E) \tag{5}$$

$E \in i^{P^2 \times C \times D}$. After the corresponding position encoding, the vector $z$ input into the transformer (TF) can be expressed as shown in formula (6).

$$z = \left[ GeLU(x_p^1 E), GeLU(x_p^2 E), \cdots, GeLU(x_p^N E) \right] + E_{pos} \tag{6}$$

$E_{pos} \in i^{N \times D}$. Finally, the input vector z is subjected to the calculations shown in formula (7) and formula (8).

$$z' = MSA(LN(z)) + z \tag{7}$$

$$z_1 = MLP(LN(z')) + z' \tag{8}$$

Multiple TF blocks are used to build the Trans block, as shown in Figure 7, and the generated output of global road information from the original remote sensing image, $z_{out}$, is expressed as shown in formula (9).

$$z_{out} = z_1 + z_2 + z_3 \tag{9}$$

**Figure. 7. Trans block structure.**

**Decoder based on context information fusion**

After the two different kinds of information discussed above are obtained, the final branch combines the two different kinds of road semantic information to fuse the global and local road features, as shown in part A of Figure 8. The high-level feature information is extracted using two 3×3 convolutional layers without changing the feature map resolution. Group normalization (GN) is used to normalize the features, and Leaky-ReLU is used as the nonlinear activation function to capture the nonlinear characteristics of the roads. This process fully integrates the semantic information from the two different sources, resulting in more comprehensive road semantic feature information that covers the feature information of various road categories.

**Figure. 8. Information fusion strategy.**

368 To solve the problem that the spatial location information of the roads is not obvious due to
369 multiple convolutions, this paper adopts the design concept of the U-Net network structure and
370 employs skip connections in the decoder to transfer the low-level feature information from the
371 output of the LDFE module to the corresponding decoder port for splicing, as depicted in part B
372 of Figure 8. Subsequently, the features are fused layer by layer via 3 × 3 convolutional layers,
373 and the feature map's resolution is incrementally restored via bilinear interpolation until it finally
374 reaches the original image resolution.
375     Finally, Wu & He (2018) proposed group normalization (GN) to improve the efficiency of
376 training rather than batch normalization (BN) BN works effectively for a relatively large batch
377 size. However, a small batch size leads to inaccurate estimation of the batch statistics, and
378 reducing the batch size for BN dramatically increases the model error. GN can achieve
379 approximately the same accuracy performance as BN for a moderate batch size and outperforms
380 other normalization variants because it can still achieve a small error rate even when the batch
381 size undergoes large fluctuations. The GN calculation is expressed as shown in formula (10):
382 (*Wu & He 2018*)

383
$$y_i = \frac{\gamma}{\sigma_i}\left(x_i - \mu_i\right) + \beta \tag{10}$$

384 $\gamma$ and $\beta$ are trainable scale and shift parameters, respectively, and $\mu_i$ and $\sigma_i$ in formula (10)
385 are the mean and standard deviation computed as shown in formula (11): (*Wu & He 2018*)

386
$$\mu_i = \frac{1}{m}\sum_{k \in S_i} x_k, \sigma_i = \sqrt{\frac{1}{m}\sum_{k \in S_i}\left(x_k - \mu_i\right)^2 + \varepsilon} \tag{11}$$

387 $\varepsilon$ is a small nonzero constant. $S_i$ is the set of pixels, and $m$ is the size of $S_i$.
388
389 **Experimental design description**
390 DPIF-Net was trained on an RTX 3090 GPU. At the beginning of training, we utilized common
391 data augmentation techniques, including translation, rotation, flipping, scaling, and random color
392 jitter as shown in Figure. 9, which contributed to improving the model's robustness and
393 performance. The model implementation was based on PyTorch, the optimizer for all structures
394 was Adam, the initial learning rate was set to 0.0002, and the batch size during training was set
395 to 2. The MSELoss was used to calculate the loss.
396
397 **Figure. 9. Data Augmentation Strategy**. (a) Original image, (b) HSV color jitter, (c)
398 translation, (d) flip, (e) random rotation, (f) translation and rotation. (All satellite images and
399 masks from the Massachusetts dataset)
400
401     To assess the performance of DPIF-Net in rural road extraction, we conducted three
402 primary comparative experiments across distinct datasets. Each of these comparative
403 experiments entailed a juxtaposition between mainstream models and ours. In addition, the
404 performance of DPIF-Net was compared with that of the U-Net, SegNet, D-LinkNet, and

409  DeepLabv3+ models on the above three road datasets. In these experiments, our goal was to
410  observe its capabilities for extracting roads in complex scenarios through evaluation metrics,
411  assessing both the completeness and accuracy of road extraction. In the discussion, we assessed
412  the parameter sizes among the models and the visual assessments of both completeness and
413  accuracy. Furthermore, the ablation experiments are conducted to observe the contributions of its
414  two branches, which are designed to identify which modules are most crucial for DPIF-Net's
415  performance.
416
417  **Experimental Evaluation Metrics**
418  To evaluate the effectiveness of our model, four common metrics are selected: intersection over
419  union (IoU), precision, recall, and $F_1$ score ($F_1$). High IoU indicates the model's accuracy in
420  predicting the location and shape of roads, which is particularly crucial for assessing the model's
421  ability to recognize roads (Lian et al. 2020). The IoU is calculated as shown in formula (12).

$$IoU = \frac{TP}{TP + FP + FN} \tag{12}$$

423      TP, FP, TN, and FN denote true positives, false positives, true negatives, and false
424  negatives, respectively.
425      In the context of rural roads, where non-road areas often constitute a significant proportion,
426  leading to an imbalance between positive and negative samples in the dataset, $F_1$ score becomes
427  particularly crucial for assessing performance under such conditions, which are calculated as
428  shown in formula (13).

$$F_1 = \frac{2 \times precision \times recall}{precision + recall} \tag{13}$$

430      The precision and recall are employed to evaluate the model's capability to correctly
431  identify roads, which are calculated as shown in formula (14) and formula (15), respectively.

$$precision = \frac{TP}{TP + FP} \tag{14}$$

$$recall = \frac{TP}{TP + FN} \tag{15}$$

434
435  **Figure. 10. Road extraction results for our dataset.**
436
437  # Results
438
439  **Road extraction experiment based on our road dataset**
440  The results on our dataset are shown in Figure 10. Through detailed comparisons, we find that in
441  Examples 1–2 and Examples 4–7, DPIF-Net generates fewer broken road segments
442  (misidentified connected vectors) than the other four models. At the same time, the occluded
443  parts caused by trees can be completely extracted to ensure the connectivity of the roads. In
444  Example 3, DPIF-Net extracts the least erroneous information, yielding results almost consistent

with the ground truth, whereas the other four models incorrectly extract road information in this example.

**Table 1**. **Comparison of different road extraction methods on the GF2 road dataset. The best results are highlighted in boldface.**

More detailed comparison results are presented in Table 1. DPIF-Net achieves the best performance in two metrics, namely, IoU and $F_1$ score, reaching 61.40% and 76.08%, respectively. In addition, DPIF-Net reaches 77.27% precision, which is 1.05% lower than the highest score and 0.6% lower than U-Net's score. DPIF-Net also reaches 74.94% recall, which is 0.08% lower than D-LinkNet's highest recall score of 75.02%, representing a small, almost negligible fluctuation.

Compared with U-Net, the IoU, recall and $F_1$ values of DPIF-Net are increased by 3.34%, 5.41%, and 2.61%, respectively. Compared with DeepLabv3+, the IoU, recall and F1 values are increased by 4.8%, 7.83% and 3.8%, respectively. Compared with D-LinkNet, the IoU, precision and F1 values increased by 0.35%, 0.64%, and 0.26%, respectively. Compared with SegNet, the IoU, precision, recall and F1 values are increased by 7.52%, 1.52%, 9.83%, and 6.05%, respectively. The largest increases in IoU, precision, recall, and F1 score reach 7.52%, 1.52%, 9.83%, and 6.05%, respectively. In the extraction of rural roads, the crucial aspects lie in correctly identifying roads and preserving their completeness, both reflected in the IoU and F1 metrics. DPIF-Net achieves the highest IoU and F1 on the lower level roads, highlighting its superiority in extracting rural roads, which underscores its proficiency in correctly extracting rural roads and recognizing road shapes.

**Road extraction experiment based on the DeepGlobe road dataset**

A second experiment was conducted on the DeepGlobe road dataset, and the results of the visual assessment comparison are presented in Figure 11. Figure 11 shows the detailed effects of rural road segment extraction for 8 examples. In the first example, compared with U-Net and SegNet, DPIF-Net extracts more complete road information. In contrast, D-LinkNet produces more mistakenly extracted road segments and more fractures than DPIF-Net. In the 2nd to 6th examples, the extraction results of U-Net and SegNet show more broken segments, while DPIF-Net achieves almost the same road integrity as D-LinkNet, whereas D-LinkNet misextracts more road segments than DPIF-Net. In the 7th and 8th examples, the integrity of the roads extracted by DPIF-Net is higher than that of the other three methods.

**Figure. 11. Road extraction results for the DeepGlobe road dataset.** (a) Ground truth. (b) U-Net. (c) SegNet. (d) DeepLabv3+. (e) D-LinkNet. (f) Ours.

**Table 2. Comparison of different road extraction methods on the DeepGlobe road dataset. The best results are highlighted in boldface.**

491   The detailed evaluation index comparisons are shown in Table 2. From the indicator data in

492 Table 2, it can be seen that the model with the worst comprehensive performance is

493 DeepLabv3+, which has the lowest score in all metrics. Its recall is lower by more than 30%, and

494 the other three metrics are lower by nearly 20%. DPIF-Net achieves the best results in terms of

495 IoU, precision and $F_1$ score, reaching values of 57%, 74.76% and 72.61%, respectively.

496 However, its recall is lower than those of D-LinkNet and SegNet at only 70.58%. D-LinkNet has

497 the highest recall score of 89.62%.

498   In summary, DPIF-Net improves the IoU by 0.6–20.08%, the precision by 0.08–20.48%,

499 and the $F_1$ score by 0.49–18.68% on the DeepGlobe road dataset. Its recall is weaker than those

500 of SegNet and D-LinkNet but 6.96% higher than that of U-Net and 16.99% higher than that of

501 DeepLabv3+.

502 **Road extraction experiment based on the Massachusetts road dataset**

503   To further test the generalization ability of the proposed DPIF-Net, a similar experiment

504 was carried out on the Massachusetts road dataset. The visual assessment of the compared

505 methods on this dataset is shown in Figure 12. All the models almost completely extract the road

506 information, but from Figure 12, it can be seen that the other models do not extract some road

507 details completely enough, resulting in various fractures. Comprehensive comparisons show that

508 the proposed model is better than the others in extracting many details.

509

510 **Figure. 12. Road extraction results on the Massachusetts road dataset.** (a) Ground truth. (b)

511 U-Net. (c) SegNet. (d) DeepLabv3+. (e) D-LinkNet. (f) Ours.

512

513 **Table 3. Comparison of different road extraction methods on the Massachusetts road**

514 **dataset. The best results are highlighted in boldface.**

515

516   Table 3 displays the detailed road extraction results on the Massachusetts road dataset,

517 revealing that DPIF-Net surpasses the other models in terms of IoU, precision, and $F_1$ score, with

518 values of 53.82%, 82.48%, and 70%, respectively. Although DeepLabv3+ achieves the highest

519 recall score of 63.92%, the IoU, $F_1$ score and precision of DeepLabv3+ are significantly lower

520 than those of all other models, with differences of 10–40%.

521 Compared to other models, DPIF-Net integrates global and local information more extensively,

522 enabling it to capture more features and thereby enhance its recognition capabilities. DPIF-Net

523 achieves superior performance in predicting road accuracy and completeness, demonstrated by

524 attaining maximum values in IoU, precision, and F1 on the DeepGlobe and Massachusetts

525 datasets, showcasing its advantages in overall road extraction.

526 **Discussion**

527 In this section, we provide a detailed discussion on several important aspects of DPIF-Net. First,

528 we elaborate on the input and output mechanisms within the network, highlighting the various

529 components and their respective roles. Second, we discuss the significance of the network

530 architecture for road extraction and how DPIF-Net utilizes the capabilities of both CNNs and

Deleted: .

Deleted: .

533  transformers for effective feature representation and information fusion. Moreover, we present a
534  comprehensive parameter comparison of DPIF-Net with other state-of-the-art models for road
535  extraction, demonstrating the effectiveness and efficiency of our proposed model. Finally, we
536  report an ablation study conducted to analyze the role of each branch of DPIF-Net, which
537  provides insights into the contribution of each component toward the final performance of the
538  model.
539
540  **Evaluation of the generalization performance and road representation ability of the**
541  **proposed model**
542       First, we evaluate the generalization performance and road representation ability of the
543  proposed model. To gain deeper insight into the inner workings of DPIF-Net, feature heatmaps
544  for five different stages in the decoder are displayed in Figure 13 to illustrate how the model
545  extracts roads. The results from the three different datasets indicate that DPIF-Net effectively
546  learns road features with clear boundaries. The feature maps show that the encoder learns various
547  levels of feature representations of the input image, and the decoder combines these
548  representations to generate accurate road masks. As the decoder performs upsampling four times,
549  the semantic information of the extracted roads becomes increasingly abstract, which highlights
550  the ability of DPIF-Net to learn high-level features. These results demonstrate that DPIF-Net not
551  only has excellent road extraction performance but also possesses good robustness and feature
552  representation capabilities; thus, it shows promising potential for various applications in remote
553  sensing image analysis.
554
555  **Figure. 13. Visualization of features at different levels.** (a) Ground truth. (b) First decoder
556  output. (c) Second decoder output. (d) Third decoder output. (e) Fourth decoder output. (f) Last
557  convolution output. (g) Final extraction result. The data sources are (1) – (2) our dataset, (3) – (4)
558  the Massachusetts road dataset, and (5) – (6) the DeepGlobe road dataset.
559
560  **Discussion on the quantity of model parameters**
561  Furthermore, a comparison of the parameters used in the experiments for each model provides
562  insight into their respective strengths and weaknesses. Table 4 presents a comprehensive
563  overview of the parameters for each model, indicating that DPIF-Net has the fewest parameters,
564  with a data volume of only 63.9 MB, which is similar to that of U-Net. While D-LinkNet
565  achieves better feature extraction in some cases, it also has a significantly larger number of
566  parameters due to its use of ResNet101 as the encoder and deeper network layers. On the other
567  hand, DeepLabv3+ also has a high number of parameters due to the use of Xception as the
568  encoding network, but it performs poorly in road extraction experiments. SegNet, with
569  approximately twice as many parameters as DPIF-Net, also yields inferior experimental results
570  for road segmentation. These comparisons highlight the trade-off between the number of
571  parameters and the effectiveness of a model for road extraction tasks. While having more
572  parameters may improve feature extraction, it also increases a model's complexity and

Deleted: .

574  computational cost. DPIF-Net, with its simple yet effective structure and relatively few
575  parameters, proves to be a promising model for road extraction from remote sensing images.
576
577  **Table 4. Comparison of the parameters of each model. The best results are highlighted in**
578  **boldface.**
579  **Ablation study**
580      To better understand the contribution of each branch in the encoder part of DPIF-Net to the
581  road extraction results, an ablation study was conducted using our dataset. A comparison of the
582  results reveals that the two branches make different levels of contributions to the road extraction
583  results, as presented in Table 5. Specifically, the branch that incorporates the CNN-based feature
584  extractor makes a stronger contribution to the final road extraction results than the branch that
585  employs the transformer-based feature extractor.
586      Notably, due to the strict parameter limitations of the final model presented in this paper,
587  only three transformer blocks were used in the transformer-based feature extractor. This resulted
588  in suboptimal performance in comparison to the CNN-based feature extractor. However, the
589  performance of the transformer-based feature extractor could be improved by stacking more
590  transformer blocks, albeit at the cost of dramatically increasing the number of parameters.
591      While DPIF-Net has demonstrated impressive performance by effectively combining a
592  CNN and a transformer, there are several potential areas for improvement in future research on
593  road extraction using DPIF-Net. First, more advanced architectures for the transformer block
594  could be explored to further enhance the model's performance without significantly increasing
595  the number of parameters. Second, modified attention mechanisms or other forms of spatial
596  information modeling might improve the model's ability to capture fine details and connectivity
597  information in the road network. Third, methods to improve the accuracy of road boundary
598  delineation and reduce false positive rates could further increase the model's utility in practical
599  applications. Overall, these potential areas for improvement could help advance the state of the
600  art in road extraction using deep learning models.
601
602  **Table 5. Comparison of the contribution of encoder branches to road extraction. The best**
603  **results are highlighted in boldface.**
604
605  ## Conclusions
606  In this study, we successfully constructed a dedicated lower level category roads and developed
607  DPIF-Net. By effectively harnessing the strengths of both transformers and CNNs, our model
608  has demonstrated excellent performance in road extraction tasks with in comparison with
609  advanced models of the same period. It not only enhances extraction accuracy but also achieves
610  higher levels of road connectivity. This achievement holds significant implications not only for
611  the field of road extraction but also for increasing attention to rural road issues in research and
612  government decision-making. Despite constraints on research time and workload, DPIF-Net can
613  leverage the rapid advancements in deep learning to enhance its two branches by incorporating

614 state-of-the-art transformer or CNN modules, thus further improving the model's performance.
615 Additionally, examining the model's generalization capabilities across different geographical
616 regions and complex weather conditions to validate its practicality is crucial. These efforts will
617 contribute to the advancement of road extraction technology in remote sensing images, providing
618 strong support for rural development and infrastructure construction.
619

627 **References**
628 Abdollahi A, Pradhan B, Shukla N, Chakraborty S, and Alamri A. 2020. Deep learning
629 approaches applied to remote sensing datasets for road extraction: A state-of-the-art
630 review. *Remote Sensing* 12:1444.
631 Badrinarayanan V, Kendall A, and Cipolla R. 2017. Segnet: A deep convolutional encoder-
632 decoder architecture for image segmentation. *IEEE transactions on pattern analysis and
633 machine intelligence* 39:2481-2495.
634 Bastani F, He S, Abbar S, Alizadeh M, Balakrishnan H, Chawla S, Madden S, DeWitt D, and
635 Ieee. 2018. RoadTracer: Automatic Extraction of Road Networks from Aerial Images.
636 31st IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Salt
637 Lake City, UT. p 4720-4728.
638 Batra A, Singh S, Pang G, Basu S, Jawahar CV, Paluri M, and Soc IC. 2019. Improved Road
639 Connectivity by Joint Learning of Orientation and Segmentation. 32nd IEEE/CVF
640 Conference on Computer Vision and Pattern Recognition (CVPR). Long Beach, CA. p
641 10377-10385.
642 Chaurasia A, and Culurciello E. 2017. Linknet: Exploiting encoder representations for efficient
643 semantic segmentation. 2017 IEEE Visual Communications and Image Processing
644 (VCIP): IEEE. p 1-4.
645 Chen L-C, Papandreou G, Kokkinos I, Murphy K, and Yuille AL. 2018a. DeepLab: Semantic
646 Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully
647 Connected CRFs. *Ieee Transactions on Pattern Analysis and Machine Intelligence*
648 40:834-848. 10.1109/tpami.2017.2699184
649 Chen L-C, Zhu Y, Papandreou G, Schroff F, and Adam H. 2018b. Encoder-decoder with atrous
650 separable convolution for semantic image segmentation. Proceedings of the European
651 conference on computer vision (ECCV). p 801-818.
652 Chen Z, Wang C, Li J, Xie N, Han Y, and Du J. 2021. Reconstruction Bias U-Net for Road
653 Extraction From Optical Remote Sensing Images. *Ieee Journal of Selected Topics in
654 Applied Earth Observations and Remote Sensing* 14:2284-2294.
655 10.1109/jstars.2021.3053603
656 China P. 2003. Ministry of Communications,"JTG B01-2003 Technical Standard of Highway
657 Engineering". Beijing: China Communications Press.
658 Demir I, Koperski K, Lindenbaum D, Pang G, Huang J, Basu S, Hughes F, Tuia D, and Raskar
659 R. 2018. Deepglobe 2018: A challenge to parse the earth through satellite images.

660          Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition
661          Workshops. p 172-181.
662 Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, Dehghani M,
663          Minderer M, Heigold G, and Gelly S. 2020. An image is worth 16x16 words:
664          Transformers for image recognition at scale. *arXiv preprint arXiv:201011929*.
665 Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, and
666          Bengio Y. 2020. Generative Adversarial Networks. *Communications of the Acm* 63:139-
667          144. 10.1145/3422622
668 Han K, Wang Y, Chen H, Chen X, Guo J, Liu Z, Tang Y, Xiao A, Xu C, and Xu Y. 2020. A
669          survey on visual transformer. *arXiv preprint arXiv:201212556* 2.
670 He H, Yang D, Wang S, Wang S, and Li Y. 2019. Road extraction by using atrous spatial
671          pyramid pooling integrated encoder-decoder network and structural similarity loss.
672          *Remote Sensing* 11:1015.
673 Hendrycks D, and Gimpel K. 2016. Gaussian error linear units (gelus). *arXiv preprint*
674          *arXiv:160608415*.
675 Li P, He X, Qiao M, Miao D, Cheng X, Song D, Chen M, Li J, Zhou T, and Guo X. 2021.
676          Exploring multiple crowdsourced data to learn deep convolutional neural networks for
677          road extraction. *International Journal of Applied Earth Observation and Geoinformation*
678          104:102544.
679 Lian R, Wang W, Mustafa N, and Huang L. 2020. Road extraction methods in high-resolution
680          remote sensing images: A comprehensive review. *IEEE Journal of Selected Topics in*
681          *Applied Earth Observations and Remote Sensing* 13:5489-5507.
682 Lourenco M, Estima D, Oliveira H, Oliveira L, and Mora A. 2023. Automatic Rural Road
683          Centerline Detection and Extraction from Aerial Images for a Forest Fire Decision
684          Support System. *Remote Sensing* 15. 10.3390/rs15010271
685 Lu X, Zhong Y, Zheng Z, Liu Y, Zhao J, Ma A, and Yang J. 2019. Multi-scale and multi-task
686          deep learning framework for automatic road extraction. *IEEE Transactions on*
687          *Geoscience and Remote Sensing* 57:9362-9377.
688 Maas AL, Hannun AY, and Ng AY. 2013. Rectifier nonlinearities improve neural network
689          acoustic models. Proc icml: Atlanta, GA. p 3.
690 Mnih V. 2013. Machine learning for aerial image labeling. University of Toronto.
691 Moradi S, Oghli MG, Alizadehasl A, Shiri I, Oveisi N, Oveisi M, Maleki M, and Dhooge J. 2019.
692          MFP-Unet: A novel deep learning based approach for left ventricle segmentation in
693          echocardiography. *Physica Medica-European Journal of Medical Physics* 67:58-69.
694          10.1016/j.ejmp.2019.10.001
695 Panboonyuen T, Vateekul P, Jitkajornwanich K, and Lawawirojwong S. 2017. An enhanced
696          deep convolutional encoder-decoder network for road segmentation on aerial imagery.
697          International conference on computing and information technology: Springer. p 191-201.
698 Ronneberger O, Fischer P, and Brox T. 2015. U-net: Convolutional networks for biomedical
699          image segmentation. International Conference on Medical image computing and
700          computer-assisted intervention: Springer. p 234-241.
701 Shamsolmoali P, Zareapoor M, Zhou H, Wang R, and Yang J. 2021. Road Segmentation for
702          Remote Sensing Images Using Adversarial Spatial Pyramid Networks. *IEEE*
703          *Transactions on Geoscience and Remote Sensing* 59:4673-4688.
704          10.1109/tgrs.2020.3016086
705 Shao Z, Zhou Z, Huang X, and Zhang Y. 2021. MRENet: Simultaneous Extraction of Road
706          Surface and Road Centerline in Complex Urban Scenes from Very High-Resolution
707          Images. *Remote Sensing* 13. 10.3390/rs13020239
708 Tao C, Qi J, Li Y, Wang H, and Li H. 2019. Spatial information inference net: Road extraction
709          using road-specific contextual information. *ISPRS Journal of Photogrammetry and*
710          *Remote Sensing* 158:155-166. 10.1016/j.isprsjprs.2019.10.001

711    Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser L, and Polosukhin I.
712        2017. Attention Is All You Need. 31st Annual Conference on Neural Information
713        Processing Systems (NIPS). Long Beach, CA.
714    Wang Y, Seo J, and Jeon T. 2022. NL-LinkNet: Toward Lighter But More Accurate Road
715        Extraction With Nonlocal Operations. *IEEE Geoscience and Remote Sensing Letters* 19.
716        10.1109/lgrs.2021.3050477
717    Wu Y, and He K. 2018. Group normalization. Proceedings of the European conference on
718        computer vision (ECCV). p 3-19.
719    Xie E, Wang W, Yu Z, Anandkumar A, Alvarez JM, and Luo P. 2021. SegFormer: Simple and
720        efficient design for semantic segmentation with transformers. *Advances in Neural*
721        *Information Processing Systems* 34:12077-12090.
722    Xie Y, Miao F, Zhou K, and Peng J. 2019. HsgNet: A Road Extraction Network Based on Global
723        Perception of High-Order Spatial Information. *Isprs International Journal of Geo-*
724        *Information* 8. 10.3390/ijgi8120571
725    Yang X, Li X, Ye Y, Lau RYK, Zhang X, and Huang X. 2019. Road Detection and Centerline
726        Extraction Via Deep Recurrent Convolutional Neural Network U-Net. *IEEE Transactions*
727        *on Geoscience and Remote Sensing* 57:7209-7220. 10.1109/tgrs.2019.2912301
728    Zhang X, Han X, Li C, Tang X, Zhou H, and Jiao L. 2019a. Aerial Image Road Extraction Based
729        on an Improved Generative Adversarial Network. *Remote Sensing* 11.
730        10.3390/rs11080930
731    Zhang Y, Xiong Z, Zang Y, Wang C, Li J, and Li X. 2019b. Topology-Aware Road Network
732        Extraction via Multi-Supervised Generative Adversarial Networks. *Remote Sensing* 11.
733        10.3390/rs11091017
734    Zhang Z, Liu Q, and Wang Y. 2018. Road Extraction by Deep Residual U-Net. *IEEE*
735        *Geoscience and Remote Sensing Letters* 15:749-753. 10.1109/lgrs.2018.2802944
736    Zhang Z, and Wang Y. 2019. JointNet: A Common Neural Network for Road and Building
737        Extraction. *Remote Sensing* 11. 10.3390/rs11060696
738    Zhou L, Zhang C, and Wu M. 2018. D-linknet: Linknet with pretrained encoder and dilated
739        convolution for high resolution satellite imagery road extraction. Proceedings of the IEEE
740        Conference on Computer Vision and Pattern Recognition Workshops. p 182-186.
741    Zhu Q, Zhang Y, Wang L, Zhong Y, Guan Q, Lu X, Zhang L, and Li D. 2021. A Global Context-
742        aware and Batch-independent Network for road extraction from VHR satellite imagery.
743        *ISPRS Journal of Photogrammetry and Remote Sensing* 175:353-365.
744    Zhu QQ, Li Z, Zhang YN, and Guan QF. 2020. Building Extraction from High Spatial Resolution
745        Remote Sensing Images via Multiscale-Aware and Segmentation-Prior Conditional
746        Random Fields. *Remote Sensing* 12. 10.3390/rs12233983
747