# Foreground-Aware Relation Network for Geospatial Object Segmentation in High Spatial Resolution Remote Sensing Imagery

Zhuo Zheng      Yanfei Zhong*      Junjue Wang      Ailong Ma
Wuhan University, Wuhan, China
{zhengzhuo, zhongyanfei, kingdrone, maailong007}@whu.edu.cn

## Abstract

*Geospatial object segmentation, as a particular semantic segmentation task, always faces with larger-scale variation, larger intra-class variance of background, and foreground-background imbalance in the high spatial resolution (HSR) remote sensing imagery. However, general semantic segmentation methods mainly focus on scale variation in the natural scene, with inadequate consideration of the other two problems that usually happen in the large area earth observation scene. In this paper, we argue that the problems lie on the lack of foreground modeling and propose a foreground-aware relation network (FarSeg) from the perspectives of relation-based and optimization-based foreground modeling, to alleviate the above two problems. From perspective of relation, FarSeg enhances the discrimination of foreground features via foreground-correlated contexts associated by learning foreground-scene relation. Meanwhile, from perspective of optimization, a foreground-aware optimization is proposed to focus on foreground examples and hard examples of background during training for a balanced optimization. The experimental results obtained using a large scale dataset suggest that the proposed method is superior to the state-of-the-art general semantic segmentation methods and achieves a better trade-off between speed and accuracy. Code has been made available at: https://github.com/Z-Zheng/FarSeg.*

## 1. Introduction

High spatial resolution earth observation technique has provided a large number of high spatial resolution (HSR) remote sensing images that can finely describe various geospatial objects, such as ship, vehicle and airplane, etc. Automatically extracting objects of interest from HSR re-
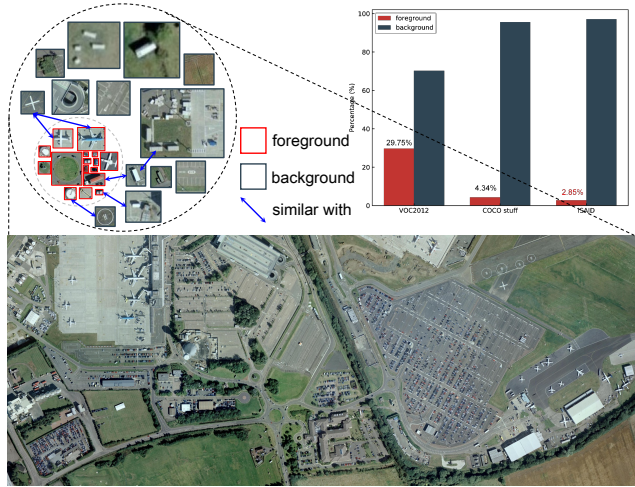
Figure 1. The main challenges of object segmentation in the HSR remote sensing imagery. (1) larger-scale variation. (2) foreground-background imbalance. (3) intra-class variance of background.

mote sensing imagery is very helpful for urban management, planing and monitoring [39, 40, 25, 26]. Geospatial object segmentation, as a significant role in object extraction, can provide semantic and location information for the objects of interest, which belongs to a particular semantic segmentation task with the goal to divide image pixels into two subsets of the foreground objects and the background area. And meanwhile, it needs to further assign a unique semantic label to each pixel in the foreground object area.

Compared with natural scene, geospatial object segmentation is more challenging in the HSR remote sensing images. There are three reasons at least:

*1) The object always has larger-scale variation in the HSR remote sensing images [14, 42].* This causes the multi-scale problem, which makes it difficult to locate and recognize the object.

*2) The background is much more complex in the HSR remote sensing images [36, 13]*, which causes serious false alarms due to larger intra-class variance.

*3) The foreground ratio is much less than it in the natural images,* as Fig. 1 shows, which causes foreground-background imbalance problem.

For natural images, the object segmentation task is directly seen as a semantic segmentation task in the computer vision field, the performance of which is mainly limited by the multi-scale problem. Therefore, current state-of-the-art general semantic segmentation methods focus on scale-aware [7] and multi-scale [5, 6, 8, 44] modeling. However, for the HSR remote sensing images, false alarms problem and foreground-background imbalance problem are ignored in these general semantic segmentation methods. We argue that this is because these methods are lack of explicit modeling for the foreground. This seriously limits the further improvement of object segmentation in the HSR remote sensing images.

In this paper, a foreground-aware relation network (FarSeg) is proposed to tackle aforementioned two problems by exploiting explicitly foreground modeling for more robust object segmentation in the HSR remote sensing imagery. We explore two perspectives of explicitly foreground modeling: relation-based and optimization-based foreground modeling, and we further propose two modules in the FarSeg: foreground-scene relation module and foreground-aware optimization. The foreground-scene relation module learns the symbiotic relation between scene and foreground to associate foreground-correlated contexts to enhance the foreground features, thus reducing false alarms. The foreground-aware optimization focus the model on the foreground by suppressing numerous easy examples in the background to alleviate the foreground-background imbalance problem.

The main contributions of our study are summarized as follows:

1. A foreground-aware relation network (FarSeg) is proposed for geospatial object segmentation in HSR remote sensing imagery.

2. To inherit multi-scale context modeling and learn geospatial scene representation, FarSeg builds a foreground branch based on the feature pyramid network (FPN) and a scene embedding branch upon a shared backbone network, namely multi-branch encoder.

3. To suppress false alarms, F-S relation module leverages the symbiotic relation between geospatial scene and geospatial objects, to associate foreground-correlated contexts and enhance the discrimination of foreground features. And meanwhile, the background without any contribution is suppressed by this symbiotic relation, thus suppressing false alarms.

4. To alleviate foreground-background imbalance, F-A optimization is proposed to focus the network on hard examples progressively, thus down-weighting gradient contribution of numerous easy examples in the background, for the foreground-background balanced training.

## 2. Related Work

**General Semantic Segmentation** Traditional methods first extract features for each pixel by the handcrafted feature descriptor. The further promotion of these traditional methods mainly depends on the improvement of handcrafted feature descriptors. However, designing a feature descriptor is time-consuming and the handcrafted feature is not robust due to limitation of prior knowledge of the expert.

The success of deep learning-based methods lies in solving this problem by learning feature representation from data directly [17]. Convolutional neural network (CNN), as structured feature representation framework in deep learning, has been explored for semantic segmentation via patch-wise classification [11, 17, 19, 18, 37]. However, patch-wise fashion limits the spatial context modeling and brings redundant computation on overlapped areas between patches. To solve this problem, fully convolutional network (FCN) [33] was proposed, which directly outputs the pixel-wise prediction from the input with arbitrary size via the in-network upsampling layer. FCN was the first pixels-to-pixels semantic segmentation method and was end-to-end trained.

To further exploit spatial context for semantic segmentation, deeplab v1 [4] utilized atrous convolution to enlarge receptive field of the CNN for wider spatial context modeling. And a dense conditional random field (CRF) was used as a postprocess to smooth the prediction.

To learn multi-scale feature representation, atrous spatial pyramidal pooling (ASPP) [5] and pyramid pooling module (PPM) [48] were proposed. ASPP utilized multiple atrous convolutions with different atrous rate to extract features with the different receptive field, while PPM generated pyramidal feature maps via pyramid pooling [20]. The image-level features and batch normalization were embedded into ASPP for further improvement of accuracy in deeplab v3 [6]. DenseASPP [44] further enhanced multi-scale feature representation via densely connected ASPP to make the multi-scale features covering larger and denser scale range. However, these methods failed to extract fine details of the object, such as the edge.

U-Net [38] and SegNet [1] utilized a new "encoder-decoder" network architecture, which reused the shallow features with high spatial resolution to enhance the deep features with strong semantics on spatial detail. RefineNet [30] proposed a multi-path refinement network to progressively recover the spatial detail of deep features for better accuracy and visual performance. Deeplab v3+ also
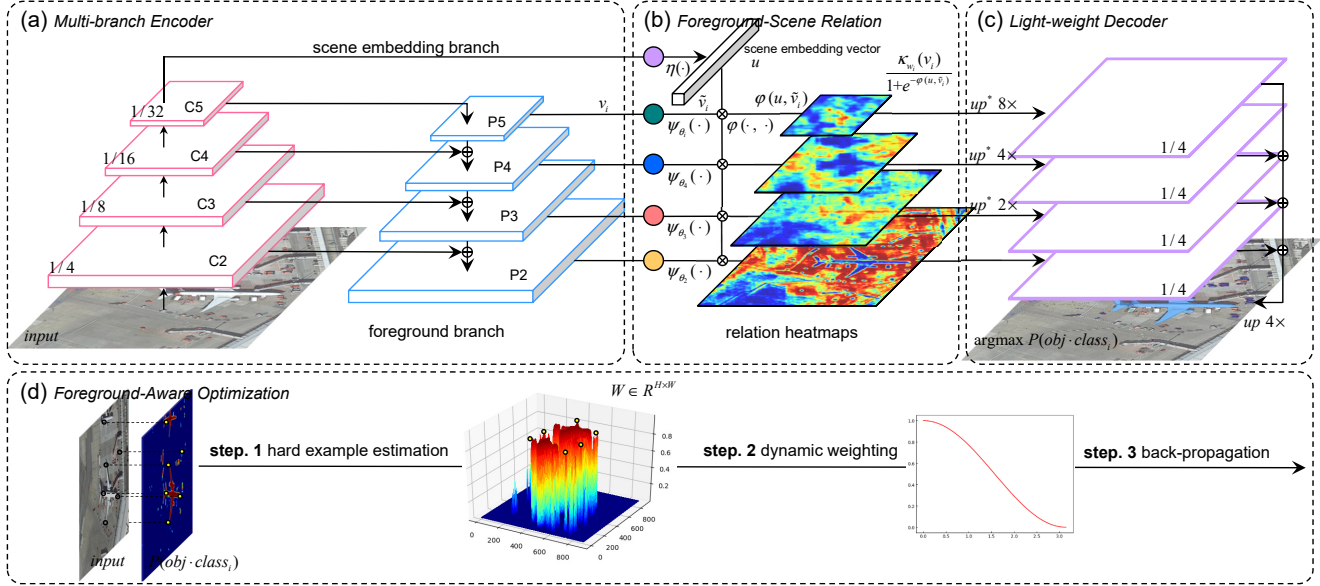
Figure 2. Overview of FarSeg. (a) Multi-branch Encoder for multi-scale object segmentation. (b) Foreground-scene relation module. (c) Light-weight decoder. (d) Foreground-aware optimization. The yellow dots indicate the relative positions of hard example in the raw image, probability map and estimation surface for a simple demonstration.

adopted "encoder-decoder" framework to further improve performance via a more powerful backbone Xception [10] and a light-weight decoder to recover the spatial resolution of features with a small overhead.

These general semantic segmentation methods mainly focus on multi-scale context modeling, ignoring the special issues in the HSR remote sensing imagery, such as false alarms and foreground-background imbalance. This causes that these methods are lack of explicit modeling for the foreground. Therefore, a foreground-aware method is needed for object segmentation in the HSR remote sensing imagery.

**Semantic Segmentation in Remote Sensing Community** There are a lot of applications using semantic segmentation technique in the remote sensing community, such as land use and land cover (LULC) classification [46, 23, 47], building extraction [45, 24, 43, 15], road extraction [29, 9, 2, 34, 3], vehicle detection [35], etc. The main methodologies follow general semantic segmentation, but for special application scenario (e.g. road or building), there were many improved techniques [2, 15, 3] for its application scenario.

However, these methods mainly focus on the improvement under the special application scenario, ignoring the consideration of common issues for object segmentation in the HSR remote sensing imagery, such as false alarms problem and foreground-background imbalance problem, especially for large scale HSR remote sensing imagery. Hence, we propose a foreground-aware relation network (FarSeg) to tackle these problems.

## 3. Foreground-Aware Relation Network

To explicit model the foreground for object segmentation in the HSR remote sensing imagery, we propose a foreground-aware relation network (FarSeg), as shown in Fig. 2. The proposed FarSeg consists of a variant of feature pyramid network (FPN), foreground-scene (F-S) relation module, light-weight decoder and foreground-aware (F-A) optimization. FPN is responsible for multi-scale object segmentation. In the F-S relation module, we first formulate false alarms problem as a problem of lacking discriminative information in the foreground, and then introduce the latent scene semantics and F-S relation to improve the discrimination of foreground features. The light-weight decoder is simply designed to recover the spatial resolution of semantic features. To make the network focus on foreground during training, the F-A optimization is proposed to alleviate foreground-background imbalance problem.

### 3.1. Multi-Branch Encoder

Multi-branch encoder is made up of a foreground branch and a scene embedding branch. As shown in Fig. 2 (a), these branches are built upon a backbone network. In the proposed method, ResNets [21] are chosen as the backbone network for basic feature extraction. $\{C_i | i = 2, 3, 4, 5\}$ denotes the set of feature maps extracted from ResNets, where the feature map $C_i$ has a output stride of $2^i$ pixels with respect to the input image. Similar to the original FPN, the top-down pathway and lateral connections are used to generated pyramidal feature maps $\{P_i | i = 2, 3, 4, 5\}$ with a

same number of channels $d$. We formulate this procedure as follows:

$$P_i = \zeta(C_i) + \Gamma(P_{i+1}), i = 2, 3, 4, 5 \qquad (1)$$

where $\zeta$ denotes the lateral connection implemented by a learnable $1\times1$ convolutional layer and $\Gamma$ denotes a nearest neighbor upsampling with a scale factor of 2. By this top-down pathway and lateral connections, the feature maps can be enhanced with high spatial detail from shallow layers and strong semantics from deep layers, which is helpful to recover the detail of objects and multi-scale context modeling. Apart from the pyramidal feature maps $v_i$, a extra branch is attached on $C_5$ to generate a geospatial scene feature $C_6$ via global context aggregation. For simplicity, we use global average pooling as the aggregation function. $C_6$ is used to model the relation between geospatial scene and foreground, which is illustrated in the Section 3.2.

### 3.2. Foreground-Scene Relation Module

The background is much more complex in the HSR remote sensing imagery. It means that there is larger intra-class variance in the background, which causes the false alarms problem. To alleviate this problem, foreground-scene (F-S) relation module is proposed to improve the discrimination of foreground features by associating geospatial scene-relevant context. The main idea is shown in Fig. 3. F-S relation module first explicitly models the relation between foreground and geospatial scene and use latent geospatial scene to associate the foreground and relevant context. And then the relation is used to enhance the input feature maps to increase the disparity between foreground features and background features, thereby improving the discrimination of foreground features.

As shown in Fig. 2 (b), for the pyramidal feature map $v_i$, F-S relation module will produce a new feature map $z_i$. The feature map $z_i$ is obtained by re-encoding $v_i$ and then re-weighting it using the relation map $r_i$. The relation map $r_i$ is the similarity matrix between geospatial scene representation and foreground representation. To align these two feature representations into a shared manifold $R^{d_u}$, there are two projection functions needed to learn for geospatial scene and foreground, respectively. $\tilde{v}_i$ is the feature map $v_i$ transformed by the scale-aware projection function $\psi_{\theta_i}(\cdot) : R^{d \times H \times W} \mapsto R^{d_u \times H \times W}$, as shown in Eqn. 2.

$$\tilde{v}_i = \psi_{\theta_i}(v_i) \qquad (2)$$

where $\theta_i$ denotes the learnable parameters of $\psi_{\theta_i}(\cdot)$. We adopt a simple form of $\psi_{\theta_i}(\cdot)$ which is just implemented by $1\times1$ convolutional layer followed by batch normalization and ReLU in order.

To compute the relation map $r_i$, a 1-D scene embedding vector $u \in R^{d_u}$ is needed to interact with the foreground
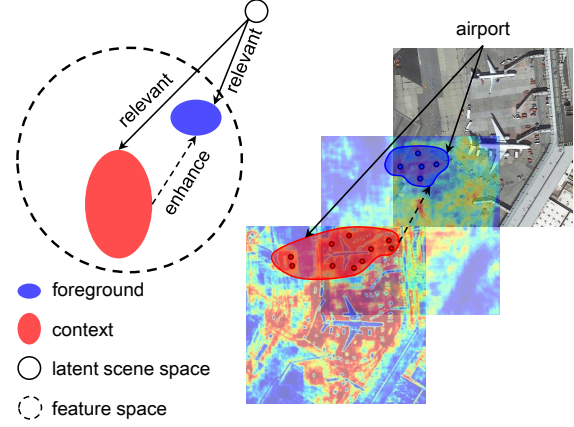


Figure 3. Concept of F-S relation. The foreground features are associated with relevant context features by their collaborative latent geospatial scene space. Meanwhile, the relevant context features are utilized to enhance the discrimination of the foreground features.
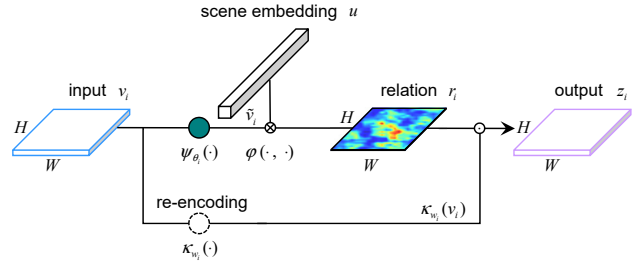


Figure 4. The computation detail of relation modeling for the pyramid level $i$ in the F-S relation module. The input and output have the same spatial size.

feature maps $\tilde{v}_i$ in the shared manifold. The scene embedding vector $u$ is computed by applying $\eta(\cdot)$ on $C_6$, as shown Eqn. 3.

$$u = \eta(C_6) \qquad (3)$$

where $\eta$ denotes a projection function for geospatial scene representation and it is implemented by a learnable $1\times1$ convolutional layer with output channels of $d_u$. The scene embedding vector $u$ is shared for each pyramid because the latent geospatial scene semantics is scale-invariant cross all pyramids. Hence, the relation map $r_i$ can be naturally obtained by Eqn. 4.

$$r_i = \varphi(u, \tilde{v}_i) = u \odot \tilde{v}_i \qquad (4)$$

where $\varphi$ denotes the similar estimation function and it is implemented by pointwise inner product for simplicity and efficient computational complexity.

For each pyramid level, the process detail of the relation modeling is illustrated in Fig. 4 and relation enhanced fore-

ground feature maps $z_i$ is computed as follows:

$$z_i = \frac{1}{1 + \exp(-r_i)} \cdot \kappa_{w_i}(v_i) \qquad (5)$$

where $\kappa_{w_i}(\cdot)$ is the encoder with learnable parameters $w_i$ for input feature maps $v_i$. The encoder is designed to introduce a extra non-linear unit to avoid feature degradation since the weighting operation is a linear function. Therefore, we adopt a simple form of this encoder, which implemented by a $1 \times 1$ convolutional layer followed by batch normalization and ReLU for high efficiency of parameters and computation. The item including $r_i$ of Eqn. 5 is used to weight the re-encoded feature maps, which is the normalized relation map using the sigmoid gate function based on a simple self-gating mechanism [22].

### 3.3. Light-weight decoder

The light-weight decoder is designed to recover the spatial resolution of relation enhanced semantic feature maps from F-S relation module in a light-weight fashion. The detailed architecture of the light-weight decoder is illustrated in Fig. 5.
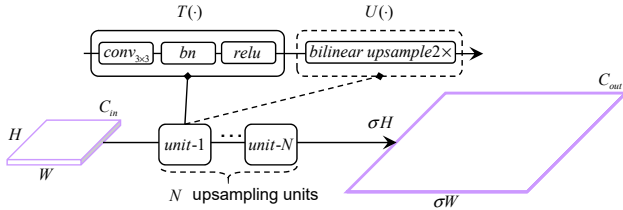


Figure 5. Abstract architecture of the light-weight decoder for each pyramid level.

Given the pyramid feature maps $z_i \in R^{C_{in} \times H \times W}$ from F-S relation module, the upsampled feature maps $z_i' \in R^{C_{out} \times \sigma H \times \sigma W}$ is computed via the light-weight decoder. The light-weight decoder is stacked by many upsampling units. The upsampling unit is made up of a channel transformation $T(\cdot)$ and an optional $2 \times$ upsampling operation $U(\cdot)$, which only includes $T(\cdot)$ if the scale factor $\sigma = 1$. Hence, the light-weight decoder for pyramid level $i$ can be simply formulated as:

$$z_i' = \begin{cases} \underbrace{U \circ T(z_i)}_{N}, & N > 0, \\ T(z_i), & N = 0. \end{cases} \qquad (6)$$

where $N$ denotes the number of upsampling units and $N = i - 2$.

$T(\cdot)$ is implemented by a $3 \times 3$ convolutional layer followed by batch normalization and ReLU. $U(\cdot)$ is the bilinear upsampling with a scale factor of 2. The total upsampling scale $\sigma$ is equal to $2^N$ because the output stride is 4.

To aggregate upsampled feature maps from each pyramid, the point-wise mean operation followed by $1 \times 1$ convolutional layer is adopted for computation and parameter efficiency. And a $4 \times$ bilinear upsampling is used to produce the final class probability map of the same size as the input image.

### 3.4. Foreground-Aware Optimization

The foreground-background imbalance problem usually causes the fact that background examples dominate the gradients during training. However, only the hard part of the background examples is valuable for optimization in the late period of training, where the hard examples are much less than easy examples in the background. Motivated by this, the foreground-aware optimization is proposed to make the network focus on foreground and hard examples in the background for a balanced optimization. The foreground-aware optimization includes three steps: hard example estimation, dynamic weighting and back-propagation, as shown in Fig. 2 (d).

**hard example estimation.** This step is used to obtain the weights reflecting the hard degree of examples to adjust the distribution of pixel-wise loss. That the example is harder means that its weight is larger. Motivated by focal loss [31], we adopt $(1 - p)^\gamma$ as weight to estimate hard examples, where $p \in [0, 1]$ is the predicted probability by the network and $\gamma$ is the focusing factor. This formulation was used in object detection, but for the pixel-level task with foreground-background imbalance, we only expect to adjust the loss distribution without change of sum for avoiding gradient vanishing. Therefore, we generalize it for object segmentation in the HSR remote sensing imagery by introducing a normalization constant $Z$ that guarantees $\sum l(p_i, y_i) = \frac{1}{Z} \sum (1 - p_i)^\gamma l(p_i, y_i)$, where $l(p_i, y_i)$ denotes the cross entropy loss of $i$-th pixel computed by predicted probability $p_i$ and its ground truth $y_i$. Hence, for the loss of each pixel, it has a weight $\frac{1}{Z}(1 - p_i)^\gamma$.

**dynamic weighting.** The hard example estimation relies on the discrimination of the model. However, the discrimination is unconfident in the initial period of training, which makes the hard example estimation unconfident. If this unconfident hard example weights are used, the model training will be unstable, influencing the converged performance. To solve this problem, we propose a dynamic weighting strategy based on an annealing function. We design three annealing functions as the candidates, as Table 1 lists. Given the cross entropy loss $l(p_i, y_i)$, the dynamic weighted loss is formulate as:

$$l'(p_i, y_i) = [\frac{1}{Z}(1 - p_i)^\gamma + \zeta(t)(1 - \frac{1}{Z}(1 - p_i)^\gamma)] \cdot l(p_i, y_i) \qquad (7)$$

where $\zeta(\cdot)$ denotes an annealing function with respect to current training step $t$ and $\zeta(t) \in [0, 1]$ is a monotonically

Table 1. Candidates of annealing functions.

| Annealing function | Formula | Hyperparameter |
|---|---|---|
| Linear | $\zeta(t) = 1 - \frac{t}{annealing\_step}$ | $annealing\_step$ |
| Poly | $\zeta(t) = (1 - \frac{t}{annealing\_step})^{decay\_factor}$ | $annealing\_step, decay\_factor$ |
| Cosine | $\zeta(t) = 0.5 * (1 + \cos(\frac{t}{annealing\_step}\pi))$ | $annealing\_step$ |

decreasing function. By this way, the focus of loss distribution can progressively move on hard examples with the increase of the confidence of hard example estimation.

# 4. Experiments

## 4.1. Experimental setting

**Dataset.** iSAID [41] dataset consists of 2,806 HSR remote sensing images. These images were collected from multiple sensors and platforms with multiple resolutions. The original image sizes range from $\sim 800 \times 800$ pixels to $\sim 4000 \times 13000$ pixels. The iSAID dataset provides 655,451 instances annotations over 15 categories[1] of the object, which is the largest dataset for instance segmentation in the HSR remote sensing imagery. The predefined training set contains 1,411 images, while validation (*val*) set contains 458 images and test set has 937 images. In this work, we only use semantic mask annotations for object segmentation. And we use the predefined training set to train models and evaluate on the validation set. Because the test set is unavailable.

**Implementation detail.** The backbone used in FarSeg was ResNet-50 for all the experiments, which was pretrained on ImageNet [12]. The channels $d$ in FPN was set to 256 and the dimension of shared manifold $d_u$ in F-S relation module was set to 256 if not specified. The default focusing factor $\gamma$ in F-A optimization was 2. For hyperparameters introduced by F-A optimization, $annealing\_step$ was set to 10k and $decay\_factor$ was set to 0.9 for the poly annealing function. For all the experiments, these models were trained for 60k iterations with a "poly" learning rate policy, where the initial learning rate was set to 0.007 and multiplied by $(1 - \frac{step}{max\_step})^{power}$ with $power = 0.9$. We used synchronized SGD over 2 GPUs with a total of 8 images per mini-batch (4 images per GPU), weight decay of 0.0001 and momentum of 0.9. The synchronized batch normalization was used for cross-gpu communication of statistic in the batch normalization layer. For data augmentation, horizontal and vertical flip, rotation of $90 \cdot k$ ($k = 1, 2, 3$) degree were adopted during training. For extra data preprocessing, we crop the image into a fixed size of (896, 896) using a sliding window striding 512 pixels.

[1]The categories are defined as: ship (Ship), storage tank (ST), baseball diamond (BD), tennis court (TC), basketball court (BC), ground field track (GTF), bridge (Bridge), large vehicle (LV), small vehicle (SV), helicopter (HC), swimming pool (SP), roundabout (RA), soccerball field (SBF), plane (Plane), harbor (Harbor).

**Evaluation metric.** Following the common practice [16, 32], we used the mean intersection over union (mIoU) as the main metric for object segmentation to evaluate the proposed method.

## 4.2. Comparison to General methods

To evaluate the FarSeg, we conduct comprehensive experiments on a larger scale HSR remote sensing images dataset. We compared FarSeg with several CNN-based methods from classical to state-of-the-art, including U-Net [38], FCN-8s [33], DenseASPP [44], Deeplab v3 [6], Semantic FPN [27], Deeplab v3+ [8], RefineNet [30], PSPNet [48]. The quantitative results listed in Table 2 suggest that FarSeg outperforms other methods in HSR scenario.

Fig. 6 shows the trade-off between speed and accuracy. It indicates that FarSeg achieves a better trade-off between speed and accuracy, which benefits from the light-weight and effective module design.



Figure 6. Speed (FPS) versus accuracy (mIoU) on iSAID *val* set. The radius of circles represents the number of parameters.

## 4.3. Ablation Study

In this section, we conduct comprehensive experiments to analyze the proposed modules and many important hyper-parameters in FarSeg. The baseline is composed of a FPN and a light-decoder, optimizing cross entropy loss. The mIoU is evaluated on iSAID *val* set with the same experimental settings if not specified.

### 4.3.1 Foreground-Scene Relation Module

**The effect of F-S relation module.** Table 3 (b) and (c) show the ablation results of adding F-S relation based on baseline method (Table 3 (a)). F-S relation modules (w/o and w/ scale-aware projection) brings 1.11% and 1.18%

Table 2. Object segmentation mIoU (%) on iSAID *val* set. The bold values in each column means the best entries.

| Method | backbone | mIoU (%) | IoU per category (%) | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Ship | ST | BD | TC | BC | GTF | Bridge | LV | SV | HC | SP | RA | SBF | Plane | Harbor |
| U-Net [38] | - | 37.39 | 49.0 | 0 | 6.51 | 78.60 | 22.89 | 5.52 | 7.48 | 49.89 | 35.62 | 0 | 38.03 | 46.49 | 9.67 | 74.74 | 45.64 |
| FCN-8s [33] | VGG-16 | 41.66 | 51.74 | 22.91 | 26.44 | 74.81 | 30.24 | 27.85 | 8.17 | 49.35 | 37.05 | 0 | 30.74 | 51.91 | 52.07 | 62.90 | 42.02 |
| DenseASPP [44] | DenseNet-121 | 56.81 | 61.15 | 50.05 | 67.54 | 86.09 | 56.56 | 52.28 | 29.61 | 57.10 | 38.44 | 0 | 43.26 | 64.80 | 74.10 | 78.12 | 51.09 |
| Deeplab v3 [6] | ResNet-50 | 59.05 | 59.74 | 50.49 | 76.98 | 84.21 | 57.92 | 59.57 | 32.88 | 54.80 | 33.75 | 31.29 | 44.74 | 66.03 | 72.13 | 75.84 | 45.68 |
| Semantic FPN [27] | ResNet-50 | 59.31 | 63.68 | 59.49 | 71.75 | **86.61** | 57.78 | 51.64 | 33.99 | 59.15 | 45.14 | 0 | 46.42 | 68.71 | 73.58 | 80.83 | 51.27 |
| Deeplab v3+ [8] | ResNet-50 | 59.33 | 59.02 | 55.15 | 75.94 | 84.18 | 58.52 | 59.24 | 32.11 | 54.54 | 33.79 | 31.14 | 44.24 | 67.51 | 73.78 | 75.70 | 45.76 |
| RefineNet [30] | ResNet-50 | 60.20 | 63.80 | 58.56 | 72.31 | 85.28 | 61.09 | 52.78 | 32.63 | 58.23 | 42.36 | 22.98 | 43.40 | 65.63 | **74.42** | 79.89 | 51.10 |
| PSPNet [48] | ResNet-50 | 60.25 | 65.2 | 52.1 | 75.7 | 85.57 | 61.12 | **60.15** | 32.46 | 58.03 | 42.96 | 10.89 | 46.78 | 68.6 | 71.9 | 79.5 | **54.26** |
| FarSeg | ResNet-50 | **63.71** | **65.38** | **61.80** | **77.73** | 86.35 | **62.08** | 56.70 | **36.70** | **60.59** | **46.34** | **35.82** | **51.21** | **71.35** | 72.53 | **82.03** | 53.91 |

Table 3. Object segmentation mIoU (%) on iSAID *val* set. Starting from Baseline, the proposed modules are gradually added in the proposed FarSeg for the module analysis.

| Method | F-S Relation | Scale-aware Proj. | F-A Opt. | mIoU(%) | Δ#params(M) |
|---|---|---|---|---|---|
| (a) Baseline | - | - | - | 59.31 | 0 |
| (b) Baseline w/ F-S Relation | ✓ | | | 60.42 | 1.12 |
| (c) Baseline w/ F-S Relation and Scale-aware Proj. | ✓ | ✓ | | 60.49 | 2.89 |
| (d) Baseline w/ F-A Opt. | | | ✓ | 61.51 | 0 |
| (e) Baseline w/ F-S Relation and F-A Opt. | ✓ | | ✓ | 63.21 | 1.12 |
| (f) FarSeg | ✓ | ✓ | ✓ | 63.71 | 2.89 |

Table 4. Foreground-aware optimization module analysis.

| Method | | Normalization | Annealing function | mIoU(%) |
|---|---|---|---|---|
| (a) FarSeg w/o F-A **Opt.** | | - | - | 60.49 |
| (b) Loss weighted with $(1-p)^\gamma$ | | | | 56.44 |
| (c) + **Norm.** | | ✓ | | 62.98 |
| (d) + **Norm.** | + Linear Annealing | ✓ | Linear | 63.18 |
| (e) + **Norm.** | + Poly Annealing | ✓ | Poly | 63.52 |
| (f) + **Norm.** | + Cosine Annealing | ✓ | Cosine | 63.71 |

performance gains in mIoU, respectively. Δ#params denotes the extra parameters introduced by the corresponding module. It indicates that F-S relation modules are parameter efficient with only 2.89 M and 1.12 M, where relative increments of parameters are ∼ 10% and ∼ 4%, respectively. This suggests that the performance gain not only comes from the gain of parameters, but also results from F-S relation design of using geospatial scene feature associates the relevant context features to enhance the foreground features.

**Scale-aware projection for scene embedding.** The projection function $\eta$ is used for geospatial scene representation in F-S relation module. We explore whether the scale-aware projection function $\eta$ is needed for each pyramid level. The results of Table 3 (b)/(c) and (e)/(f) suggest that scale-aware projection function performs better. With F-A optimization, the gain in mIoU from scale-aware projection is larger. It indicates that geospatial scene representation is related to scale and foregrounds.

**Visual interpretation for F-S relation module.** F-S relation module has good visual interpretability, combining with geoscience knowledge. Fig. 7 shows the visualization of F-S relations in the different pyramid levels. Each pixel represents the relation intensity between the latent geospatial scene and the pixel-self. There are three classical scenarios: airport, harbor, and parking-lot. We can find that different scenarios focus on different objects that are discriminative to this scenario. For example, the harbor mainly focuses on ship and water, while the airport focuses on the airplanes and their contexts. Meanwhile, these relation maps illustrate again that the geospatial scene is related to scale, foreground, and foreground-relative contexts. Because we can find that small objects are hot in the relation map with high spatial resolution ($OS = 4$), such as small vehicle and ship. The large objects are hot in the relation map with lower spatial resolution. However, contexts are not spatial resolution-specific in the relation map. It reveals that the geospatial scene is related to scale-specific foregrounds and scale-agnostic contexts.

### 4.3.2 Foreground-Aware Optimization

**The effect of F-A optimization.** Table 3 (d) and (f) show the ablation results of adding F-A optimization based on baseline method (Table 3 (a)) and baseline method with F-S relation and scale-aware projection (Table 3 (c)), respectively. F-A optimization boosts the performance with 2.2% and 3.24% in mIoU without any extra computation and memory footprint. It indicates that F-A optimization can significantly alleviate foreground-background imbalance problem for object segmentation in the HSR remote sensing imagery. Meanwhile, it suggests that F-A optimization is compatible with the F-S relation module well.

**Normalization.** Normalization is designed to only adjust the loss distribution without change of sum for avoiding gradient vanishing. Table 4 (c) shows the result of adding normalization on the naive softmax focal loss (Table 4 (b)).

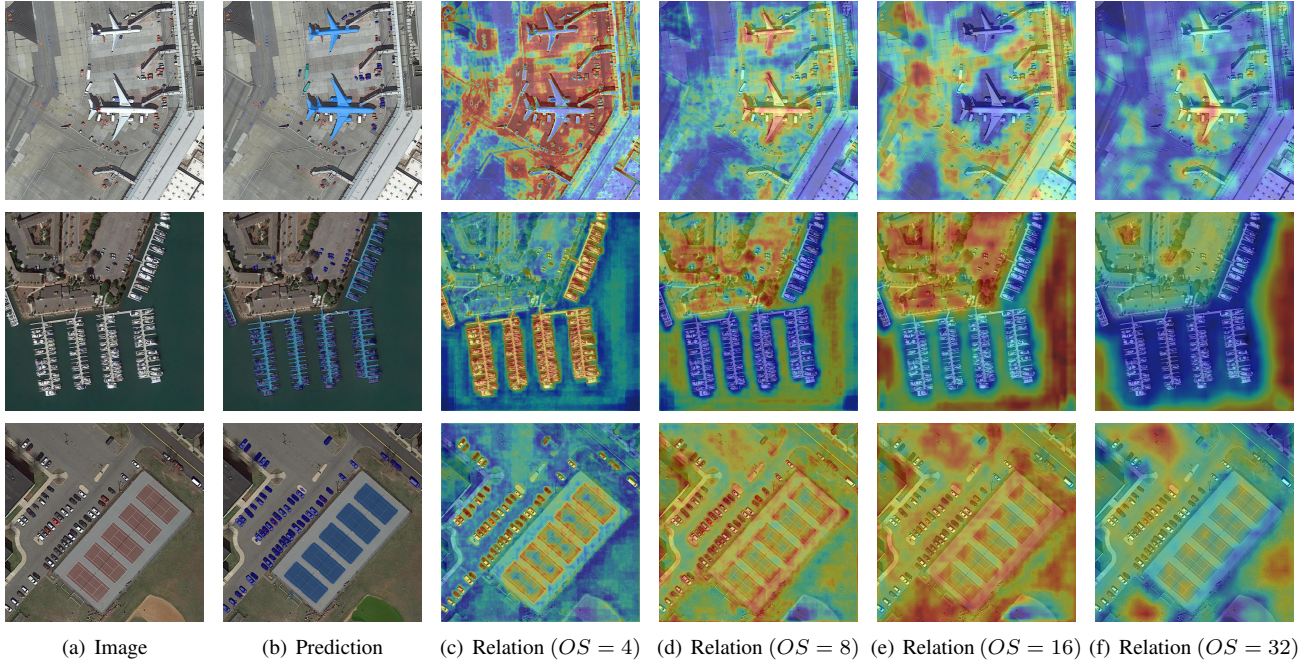| (a) Image | (b) Prediction | (c) Relation ($OS = 4$) | (d) Relation ($OS = 8$) | (e) Relation ($OS = 16$) | (f) Relation ($OS = 32$) |

Figure 7. Visualization of F-S relation heatmap in the different pyramid levels. (a) original images. (b) object segmentation results. (c)-(f) images with F-S relation heatmaps in the different pyramid level. $OS$ denotes "output stride" defined in FPN. For convenient visualization, we resize these relation maps to corresponding image sizes. Legend: Scene 1 (plane, large vehicle, small vehicle), Scene 2 (small vehicle, harbor, ship), Scene 3 (small vehicle, tennis court, baseball diamond), in a row order.

Due to instability the naive softmax focal loss for object segmentation, the mIoU drops 4.05%. However, when adding the normalization, the performance obtains significant improvement with 2.49% in mIoU. Compared with naive softmax focal loss, it gains 6.54% in mIoU. It suggests that the tuning the loss distribution without change of sum is the key to alleviate foreground-background imbalance problem.

**Annealing function.** The annealing function is used in dynamic weighting stage of F-A optimization. It aims to alleviate the training instability due to the wrong hard example estimation in the period of early training. Table 4 (d), (e), and (f) show the results of applying three proposed annealing functions. We can find that annealing-based dynamic weighting boosts the performance via reducing the wrong hard example estimation in the period of early training. Intuitively, the cosine annealing function obtains the most significant gains of 0.63% in mIoU. Because the cosine annealing function has a slow descent rate at the start and end of training, which can stably adjust the loss distribution for healthy convergence, compared with linear annealing function and polynomial annealing function.

**The choice of the focusing factor** $\gamma$**.** The focusing factor $\gamma$ is introduced to adjust the weight of hard examples.

Table 5. mIoU (%) on iSAID $val$ set using varying $\gamma$ for F-A optimization.

| $\gamma$ | 0 | 0.3 | 0.5 | 1 | 2 | 5 |
|---|---|---|---|---|---|---|
| mIoU (%) | 60.42 | 61.35 | 62.48 | 62.99 | **63.71** | 62.61 |

Larger $\gamma$, larger weight on hard examples. Following [31], we use varying $\gamma$ to conduct experiments. The results are presented in Table 5. As $\gamma$ increase, the performance obtains continually improvement. With $\gamma = 2$, F-A optimization yields 3.29% in mIoU improvement over the baseline, achieving the best result of 63.71% in mIoU. However, with $\gamma = 5$, the performance drops. The possible reason is that noise labels are wrongly seen as hard examples, as mentioned in [28].

## 5. Conclusion

In this work, we argue that false alarm and foreground-background imbalance problems are the bottlenecks of object segmentation in the HSR remote sensing imagery, while general semantic segmentation methods ignore it. To alleviate these two problems, we propose foreground-aware relation network (FarSeg), which learns foreground-scene relation to enhance the foreground features for less false alarms and trains the network using a foreground-aware optimization in foreground-background balanced fashion. The com-

prehensive experimental results show the effectiveness of FarSeg and a better trade-off between speed and accuracy.

## References

[1] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017.

[2] F. Bastani, S. He, S. Abbar, M. Alizadeh, H. Balakrishnan, S. Chawla, S. Madden, and D. DeWitt. Roadtracer: Automatic extraction of road networks from aerial images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4720–4728, 2018.

[3] A. Batra, S. Singh, G. Pang, S. Basu, C. Jawahar, and M. Paluri. Improved road connectivity by joint learning of orientation and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10385–10393, 2019.

[4] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014.

[5] 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 transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.

[6] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.

[7] L.-C. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille. Attention to scale: Scale-aware semantic image segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3640–3649, 2016.

[8] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018.

[9] G. Cheng, Y. Wang, S. Xu, H. Wang, S. Xiang, and C. Pan. Automatic road detection and centerline extraction via cascaded end-to-end convolutional neural network. *IEEE Transactions on Geoscience and Remote Sensing*, 55(6):3322–3337, 2017.

[10] F. Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017.

[11] D. Ciresan, A. Giusti, L. M. Gambardella, and J. Schmidhuber. Deep neural networks segment neuronal membranes in electron microscopy images. In *Advances in neural information processing systems*, pages 2843–2851, 2012.

[12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[13] Z. Deng, H. Sun, S. Zhou, and J. Zhao. Learning deep ship detector in sar images from scratch. *IEEE Transactions on Geoscience and Remote Sensing*, 57(6):4021–4039, 2019.

[14] Z. Deng, H. Sun, S. Zhou, J. Zhao, L. Lei, and H. Zou. Multi-scale object detection in remote sensing imagery with convolutional neural networks. *ISPRS journal of photogrammetry and remote sensing*, 145:3–22, 2018.

[15] M. Dickenson and L. Gueguen. Rotated rectangles for symbolized building footprint extraction. In *CVPR Workshops*, pages 225–228, 2018.

[16] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.

[17] C. Farabet, C. Couprie, L. Najman, and Y. LeCun. Learning hierarchical features for scene labeling. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1915–1929, 2012.

[18] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik. Learning rich features from rgb-d images for object detection and segmentation. In *European conference on computer vision*, pages 345–360. Springer, 2014.

[19] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Simultaneous detection and segmentation. In *European Conference on Computer Vision*, pages 297–312. Springer, 2014.

[20] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 37(9):1904–1916, 2015.

[21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[22] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.

[23] B. Huang, B. Zhao, and Y. Song. Urban land-use mapping using a deep convolutional neural network with high spatial resolution multispectral remote sensing imagery. *Remote Sensing of Environment*, 214:73–86, 2018.

[24] S. Ji, S. Wei, and M. Lu. Fully convolutional networks for multisource building extraction from an open aerial and satellite imagery data set. *IEEE Transactions on Geoscience and Remote Sensing*, 57(1):574–586, 2018.

[25] M. Kampffmeyer, A.-B. Salberg, and R. Jenssen. Semantic segmentation of small objects and modeling of uncertainty in urban remote sensing images using deep convolutional neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2016.

[26] R. Kemker, C. Salvaggio, and C. Kanan. Algorithms for semantic segmentation of multispectral remote sensing imagery using deep learning. *ISPRS journal of photogrammetry and remote sensing*, 145:60–77, 2018.

[27] A. Kirillov, R. Girshick, K. He, and P. Dollár. Panoptic feature pyramid networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6399–6408, 2019.

[28] B. Li, Y. Liu, and X. Wang. Gradient harmonized single-stage detector. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8577–8584, 2019.

[29] J. Liang, N. Homayounfar, W.-C. Ma, S. Wang, and R. Urtasun. Convolutional recurrent network for road boundary extraction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9512–9521, 2019.

[30] G. Lin, A. Milan, C. Shen, and I. Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1925–1934, 2017.

[31] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[32] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[33] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.

[34] X. Lu, Y. Zhong, Z. Zheng, Y. Liu, J. Zhao, A. Ma, and J. Yang. Multi-scale and multi-task deep learning framework for automatic road extraction. *IEEE Transactions on Geoscience and Remote Sensing*, 2019.

[35] L. Mou and X. X. Zhu. Vehicle instance segmentation from aerial image and video using a multitask learning residual fully convolutional network. *IEEE Transactions on Geoscience and Remote Sensing*, 56(11):6699–6711, 2018.

[36] J. Pang, C. Li, J. Shi, Z. Xu, and H. Feng. $\mathcal{R}^2$ -cnn: Fast tiny object detection in large-scale remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 57(8):5512–5524, Aug 2019.

[37] P. H. Pinheiro and R. Collobert. Recurrent convolutional neural networks for scene labeling. In *31st International Conference on Machine Learning (ICML)*, number CONF, 2014.

[38] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[39] F. Rottensteiner, G. Sohn, J. Jung, M. Gerke, C. Baillard, S. Benitez, and U. Breitkopf. The isprs benchmark on urban object classification and 3d building reconstruction. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences I-3 (2012), Nr. 1*, 1(1):293–298, 2012.

[40] M. Volpi and V. Ferrari. Semantic segmentation of urban scenes by learning local class interactions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–9, 2015.

[41] S. Waqas Zamir, A. Arora, A. Gupta, S. Khan, G. Sun, F. Shahbaz Khan, F. Zhu, L. Shao, G.-S. Xia, and X. Bai. isaid: A large-scale dataset for instance segmentation in aerial images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 28–37, 2019.

[42] G.-S. Xia, X. Bai, J. Ding, Z. Zhu, S. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang. Dota: A large-scale dataset for object detection in aerial images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3974–3983, 2018.

[43] Y. Xu, L. Wu, Z. Xie, and Z. Chen. Building extraction in very high resolution remote sensing imagery using deep learning and guided filters. *Remote Sensing*, 10(1):144, 2018.

[44] M. Yang, K. Yu, C. Zhang, Z. Li, and K. Yang. Denseaspp for semantic segmentation in street scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3684–3692, 2018.

[45] J. Yuan. Learning building extraction in aerial scenes with convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 40(11):2793–2798, 2017.

[46] C. Zhang, I. Sargent, X. Pan, H. Li, A. Gardiner, J. Hare, and P. M. Atkinson. An object-based convolutional neural network (ocnn) for urban land use classification. *Remote sensing of environment*, 216:57–70, 2018.

[47] C. Zhang, I. Sargent, X. Pan, H. Li, A. Gardiner, J. Hare, and P. M. Atkinson. Joint deep learning for land cover and land use classification. *Remote sensing of environment*, 221:173–187, 2019.

[48] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017.