# Remote Sensing Object Detection Meets Deep Learning: A Meta-review of Challenges and Advances

Xiangrong Zhang, *Senior Member, IEEE*, Tianyang Zhang, Guanchun Wang, Peng Zhu,
Xu Tang, *Senior Member, IEEE*, Xiuping Jia, *Fellow, IEEE*, and Licheng Jiao, *Fellow, IEEE*

*Abstract*—**Remote sensing object detection (RSOD), one of the most fundamental and challenging tasks in the remote sensing field, has received longstanding attention. In recent years, deep learning techniques have demonstrated robust feature representation capabilities and led to a big leap in the development of RSOD techniques. In this era of rapid technical evolution, this review aims to present a comprehensive review of the recent achievements in deep learning based RSOD methods. More than 300 papers are covered in this review. We identify five main challenges in RSOD, including multi-scale object detection, rotated object detection, weak object detection, tiny object detection, and object detection with limited supervision, and systematically review the corresponding methods developed in a hierarchical division manner. We also review the widely used benchmark datasets and evaluation metrics within the field of RSOD, as well as the application scenarios for RSOD. Future research directions are provided for further promoting the research in RSOD.**

*Index Terms*—**Object detection, Remote sensing images, Deep learning, Technical evolution**

## I. INTRODUCTION

With the rapid advances in earth observation technology, remote sensing satellites (e.g., Google Earth [1], WordWide-3 [2], and Gaofen series satellites [3]–[5]) have made significant improvements in spatial, temporal, and spectral resolutions and a massive number of remote sensing images (RSIs) are now accessible. Benefiting from the dramatic increase in available RSIs, human beings have entered an era of remote sensing big data, and the automatic interpretation of RSIs has become an active yield challenging topic [6]–[8].

RSOD aims to determine whether or not objects of interest exist in a given RSI and to return the category and position of each predicted object. The term 'object' in this survey refers to man-made or highly structured objects (such as airplanes, vehicles, and ships) rather than unstructured scene objects (e.g., land, sky, and grass). As the cornerstone in the automatic interpretation of RSIs, RSOD has received significant attention.

In general, RSIs are taken at an overhead viewpoint with different ground sampling distances (GSDs) and cover widespread regions of the Earth's surface. As a result, the geospatial objects exhibit more significant diversity in scale,

Xiangrong Zhang, Tianyang Zhang, Guanchun Wang, Peng Zhu, Xu Tang, and Licheng Jiao are with the School of Artificial Intelligence, Xidian University, Xi'an 710071, China (e-mail: xrzhang@mail.xidian.edu.cn).

Xiuping Jia is with the School of Engineering and Information Technology, University of New South Wales, Canberra, ACT 2612, Australia.

angle, and appearance. Based on the characteristics of geospatial objects in RSIs, we summarize the major challenges of RSOD in the following five aspects:

(1) **Huge Scale Variations.** On the one hand, there are generally massive scale variations across different categories of objects, as illustrated in Fig. 1(b): a vehicle may be as small as 10 pixel area, while an airplane can be 20 times larger than the vehicle. On the other hand, the intra-category objects also show a wide range of scales. Therefore, the detection models require to handle both large-scale and small-scale objects.

(2) **Arbitrary Orientations.** The unique overhead viewpoint leads to the geospatial objects often distributed with arbitrary orientations, as shown in Fig. 1(c). This rotated object detection task exacerbates the challenge of RSOD, making it important for the detector to be perceptive of orientation.

(3) **Weak Feature Responses.** Generally, RSIs contain complex context and massive background noises. As depicted in Fig. 1(a), some vehicles are obscured by shadows, and the surrounding background noises tend to have a similar appearance to vehicles. This intricate interference may overwhelm the objects of interest and deteriorate their feature representation, which results in the objects of interest being presented as weak feature responses [9].

(4) **Tiny Objects.** As shown in Fig. 1(d), tiny objects tend to exhibit extremely small scales and limited appearance information, resulting in a poor-quality feature representation. In addition, the current prevailing detection paradigms inevitably weaken or even discard the representation of tiny objects [10]. These problems in tiny object detection bring new difficulties to existing detection methods.

(5) **Expensive Annotation.** The complex characteristics of geospatial objects in terms of scale and angle, as well as the expert knowledge required for fine-grained annotations [11], make the accurate box-level annotations of RSIs a time-consuming and labor-intensive task. However, the current deep learning based detectors rely heavily on abundant well-labeled data to reach performance saturation. Therefore, the efficient RSOD methods in a lack of sufficient supervised information scenario remain challenging.

To tackle these challenges, numerous RSOD methods have emerged in the past two decades. At the early stage, researchers adopted template matching [12]–[14] and prior knowledge [15]–[17] for object detection in remote sensing scenes. These early methods rely more on hand-crafted templates or prior knowledge, leading to unstable results.
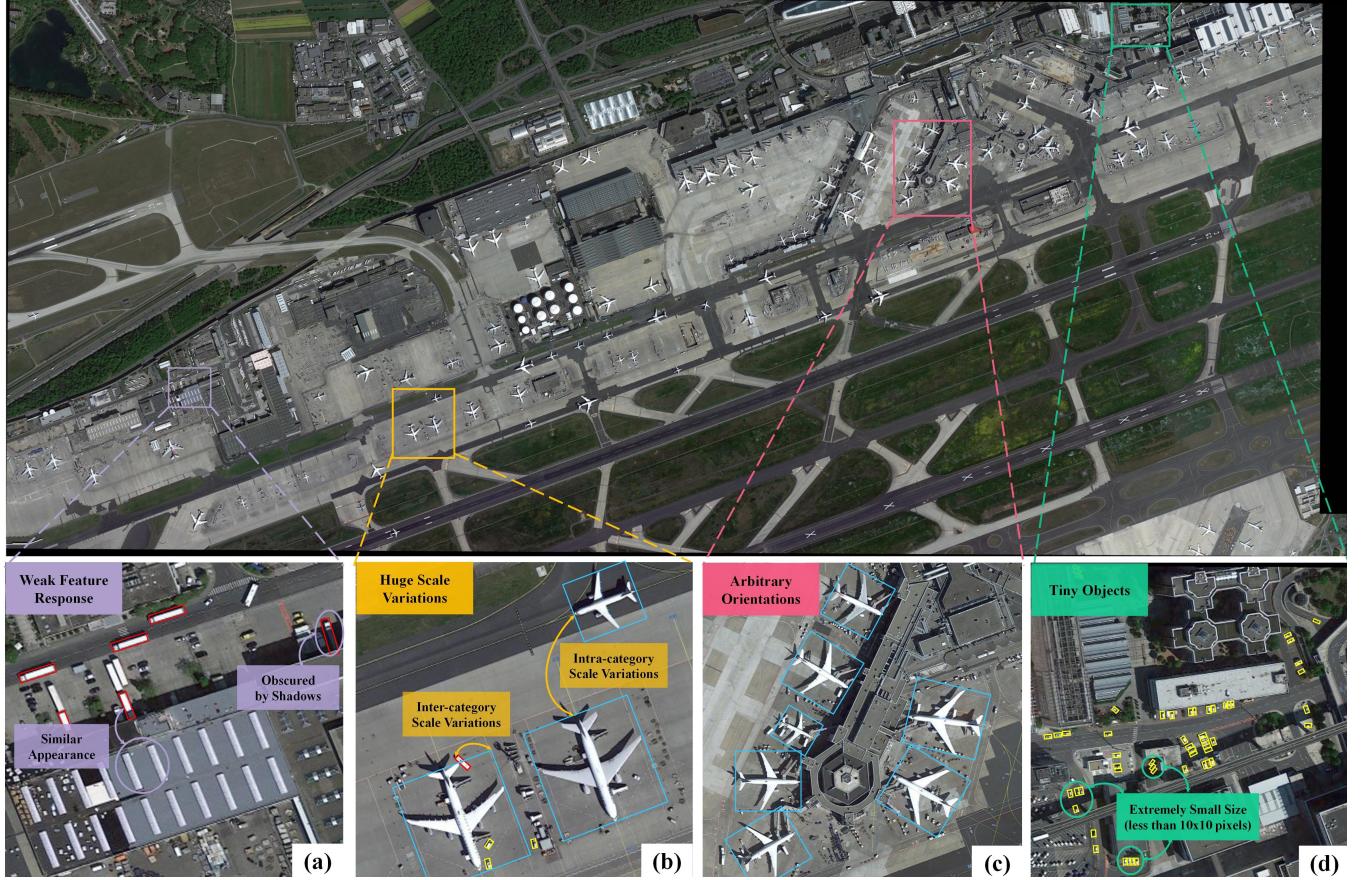
Fig. 1. A typical example of remote sensing images. (a) The complex context and massive background noises lead to weak feature responses of objects. (b) Huge scale variations exist in both inter-category and intra-category objects. (c) Objects are distributed with arbitrary orientations. (d) Tiny objects tend to exhibit extremely small scales.

Later, machine learning approaches [18]–[21] have become mainstream in RSOD, which view object detection as a classification task. Concretely, the machine learning model first searches a set of object proposals from the input image and extracts the texture, context, and other features of these object proposals. Then, it employs an independent classifier to identify the object categories in these object proposals. However, shallow learning based features from the machine learning approaches significantly restrict the representations of objects, especially in more challenging scenarios. Besides, the machine learning based object detection methods cannot be trained in an end-to-end manner, which is no longer applicable in the era of remote sensing big data.

Recently, deep learning techniques [22] have demonstrated powerful feature representation capabilities from massive amounts of data, and the state-of-the-art detectors [23]–[26] in computer vision achieve object detection ability that rivals that of humans [27]. Drawing on the advanced progress of deep learning techniques, various deep learning based methods have dominated the RSOD and led to remarkable break-throughs in detection performance. Compared to the traditional methods, deep neural network architecture can extract high-level semantic features and obtain much more robust feature representations of objects. In addition, the efficient end-to-end training manner and automated feature extraction fashion make

the deep learning based object detection methods more suitable for RSOD in the remote sensing big data era.

Along with the prevalence of RSOD, a number of geospatial object detection surveys [9], [28]–[34] have been published in recent years. For example, Cheng *et al.* [29] reviewed the early development of RSOD. Han *et al.* [9] focused on small and weak object detection in RSIs. In [30], the authors reviewed airplane detection methods. Li *et al.* [31] conducted a thorough survey on deep learning based detectors in the remote sensing community according to various improvement strategies. Besides, some work [28], [33], [34] mainly focused on publishing novel benchmark datasets for RSOD and briefly reviewed object detection methods in the field of remote sensing. Compared with previous works, this survey provides a comprehensive analysis of the major challenges in RSOD based on the characteristics of geospatial objects and system-atically categorizes and summarizes the deep learning based remote sensing object detectors according to these challenges. Moreover, more than 300 papers on RSOD are reviewed in this work, leading to a more comprehensive and systematic survey.

Fig. 2 shows the taxonomy of object detection methods in this review. According to the major challenges in RSOD, we divide the current deep learning based RSOD methods into five main categories: multi-scale object detection, rotated

Fig. 2. Structured taxonomy of the deep learning based RSOD methods in this review. A hierarchical division is adopted to detailed describe each sub-category.

object detection, weak object detection, tiny object detection, and object detection with limited supervision. In each category, we further summarize the sub-categories based on the improvement strategies or learning paradigms designed for the category-specific challenges. For multi-scale object detection, we mainly review the three widely used methods: data augmentation strategy, multi-scale feature representation, and high-quality multi-scale anchor generation. With regard to rotated object detection, we mainly focus on the rotated bounding box representation and rotation-insensitive feature learning. For weak object detection, we divide it into two classes: background noise suppressing and related context mining. As for tiny object detection, we detail it into three streams: discriminative feature extraction, super-resolution reconstruction, and improved detection metrics. According to the learning paradigms, we divide object detection with limited supervision into weakly-supervised object detection, semi-

supervised object detection, and few-shot object detection. Notably, there are still detailed divisions in each sub-category, as shown in the rounded rectangles in Fig. 2. This hierarchical division provides a systematic review and summarization of existing methods. It helps researchers understand RSOD more comprehensively and facilitate further progress, which is the main purpose of this review.

In summary, the main contributions of this review are as follows:

- We comprehensively analyze the major challenges in RSOD based on the characteristics of geospatial objects, including huge scale variations, arbitrary orientations, weak feature responses, tiny objects, and expensive annotations.
- We systematically summarize the deep learning based object detectors in the remote sensing community and categorize them in a hierarchical manner according to
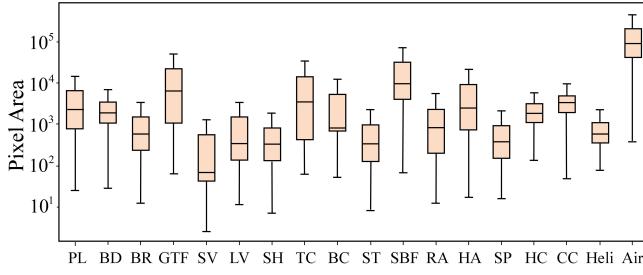
Fig. 3. Scale variations for each category (the short names for categories can be referred to [33]) in the DOTAv2.0 dataset. Huge scale variations exist in both inter-categories and intra-categories.

their motivation.

- We present a forward-looking discussion of future research directions for RSOD to motivate the further progress of RSOD.

## II. MULTI-SCALE OBJECT DETECTION

Due to the different spatial resolutions between RSIs, the huge scale variation is a notoriously challenging problem in RSOD and seriously degrades the detection performance. As depicted in Fig. 3, we present the distribution of object pixel areas for each category in the DOTAv2.0 dataset [33]. Obviously, the scales vary greatly between categories, in which a small vehicle may only contain less than 10 pixel area while an airport exceeds $10^5$ pixel area. Worse still, the huge intra-category scale variations further exacerbate the difficulties of multi-scale object detection. To tackle the huge scale variation problem, current studies are mainly divided into data augmentation, multi-scale feature representation, and multi-scale anchor generation. Fig. 4 gives a brief summary of multi-scale object detection methods.

### A. Data Augmentation

Data augmentation is a simple yet widely applied approach for increasing dataset diversity. As for the scale variation problem in multi-scale object detection, image scaling is a straightforward and effective augmentation method. Zhao *et al.* [35] fed multi-scale image pyramids into multiple networks and fused the output features from these networks to generate multi-scale feature representations. In [36], Azimi *et al.* proposed a combined image cascade and feature pyramid network to extract object features on various scales. Although image pyramids can effectively increase the detection performance for multi-scale objects, the inference time and computational complexity are severely increased. To tackle this problem, Shamsolmoali *et al.* [37] designed a lightweight image pyramid module (LIPM). The proposed LIPM receives multiple down-sampling images to generate multi-scale feature maps and fuses the output multi-scale feature maps with the corresponding scale feature maps from the backbone. Moreover, some modern data augmentation methods (e.g., Moscia and Stitcher [38]) also show remarkable effectiveness in multi-scale object detection, especially for small objects [39]–[41].
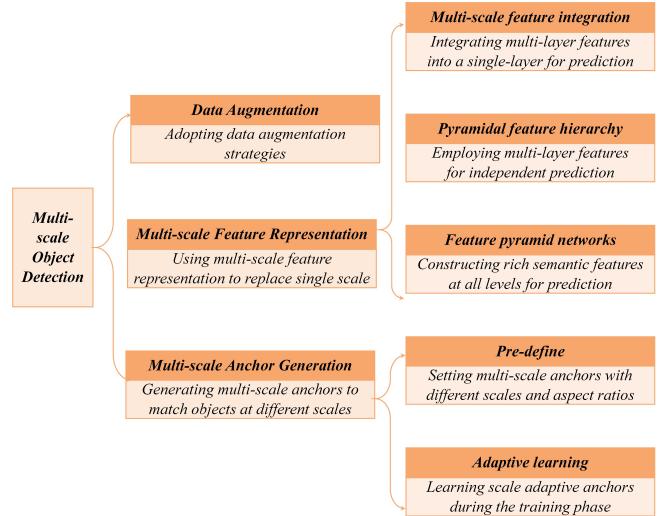


Fig. 4. A brief summary of multi-scale object detection methods.

### B. Multi-scale Feature Representation

Early studies in RSOD usually utilize the last single feature map of the backbone to detect objects, as illustrated in Fig.5(a). However, such a single-scale feature map prediction limits the detector to handle the object with a wide range of scale [42]–[44]. Consequently, multi-scale feature representation methods have been proposed and have become an effective solution to the huge object scale variation problem in RSOD. The current multi-scale feature representation methods are mainly divided into three streams: multi-scale feature integration, pyramidal feature hierarchy, and feature pyramid networks.

*1) Multi-scale Feature Integration:* Convolutional neural networks (CNN) usually adopt a deep hierarchical structure where different level features have different characteristics. The shallow-level features usually contain fine-grained features (e.g., points, edges, and textures of objects) and provide detailed spatial location information, which is more suitable for object localization. In contrast, features from higher-level layers show stronger semantic information and present discriminative information for object classification. To combine the information from different layers and generate the multi-scale representation, some researchers introduced multi-layer feature integration methods that integrate features from multiple layers into a single feature map and perform the detection on this rebuilt feature map [45]–[52]. Fig. 5(b) depicts the structure of multi-layer feature integration methods.

Zhang *et al.* [48] designed a hierarchical robust CNN to extract hierarchical spatial semantic information by fusing multi-scale convolutional features from three different layers and introduced multiple fully connected layers to enhance the rotation and scaling robustness of the network. Considering the different norms between multi-layer features, Lin *et al.* [49] applied an L2 normalization for each feature before integration to maintain stability in the network training stage. Unlike previous multi-scale feature integration at the level of the convolutional layer, Zheng *et al.* [51] designed the HyBlock to build multi-scale feature representation at the intra-layer level.
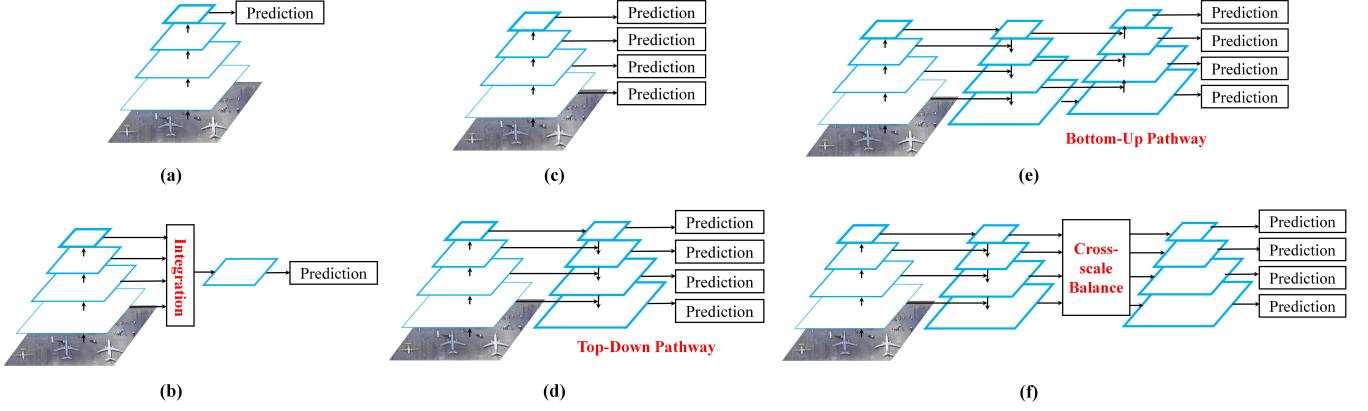
Fig. 5. Single-scale feature representation and six paradigms for multi-scale feature representation. (a) Single-scale feature representation. (b) Multi-scale feature integration. (c) Pyramidal feature hierarchy. (d) Feature pyramid networks. (e) Top-down and Bottom-up. (f) Cross-scale feature balance.

The HyBlock employs the atrous separable convolution with pyramidal receptive fields to learn the hyper-scale features, alleviating the scale-variation issue in RSOD.

*2) Pyramidal Feature Hierarchy:* The key insight behind the pyramidal feature hierarchy is that the features in different layers can encode object information from different scales. For instance, small objects are more likely to appear in shallow layers, while large objects tend to exist in deep layers. Therefore, the pyramidal feature hierarchy employs multiple-layer features for independent prediction to detect objects with a wide scale range, as depicted in Fig. 5(c). SSD [53] is a typical representative of the pyramidal feature hierarchy, which has a wide range of extended applications in both natural scenes [54]–[56] and remote sensing scenes [57]–[63].

To improve the detection performance for small vehicles, Liang [60] *et al.* added an extra scaling branch to the SSD, which consists of a deconvolution module and an average pooling layer. Referring to hierarchical regression layers in SSD, Wang *et al.* [58] introduced the scale-invariant regression layers (SIRLs), where three isolated regression layers are employed to capture the information of full-scale objects. Based on the SIRLs, a novel specific scale joint loss is introduced to accelerate network convergence. In [64], Li *et al.* proposed the HSF-Net that introduces the hierarchical selective filtering layer in both RPN and detection sub-network. Specifically, the hierarchical selective filtering layer employs three convolutional layers with different kernel sizes (e.g., $1 \times 1$, $3 \times 3$, and $5 \times 5$) to obtain multiple receptive field features, which benefits multi-scale ship detection.

*3) Feature Pyramid Networks:* Pyramidal feature hierarchy methods use independent multi-level features for detection and ignore the complementary information between features at different levels, resulting in weak semantic information for low-level features. To tackle this problem, Lin *et al.* [65] proposed the feature pyramid network (FPN). As shown in Fig. 5(d), the FPN introduces a top-down pathway to transfer the rich semantic information from high-level features to shallow-level features, leading to rich semantic features at all levels (please refer to the detailed in [65]). Thanks to the significant improvement of FPN for multi-scale object detection, FPN and its extensions [66]–[68] play a dominant role in multi-scale feature representation.

Considering the extreme aspect ratios of geospatial objects (e.g., bridges, harbors, and airports), Hou *et al.* [69] proposed an asymmetric feature pyramid network (AFPN). The AFPN adopts the asymmetric convolution block to enhance the feature representation regarding the cross-shaped skeleton and improve the performance of large aspect ratio objects. Zhang *et al.* [70] designed a Laplacian feature pyramid network (LFPN) to inject high-frequency information into the multi-scale pyramidal feature representation, which is useful for accurate object detection but has been ignored by previous work. In [71], Zhang *et al.* introduced the high-resolution feature pyramid network (HRFPN) to fully leverage the high-resolution feature representations, leading to precise and robust SAR ship detection. In addition, some researchers integrated the novel feature fusion module [72], [73], attention machine [74]–[77], or dilation convolution layer [78], [79] into the FPN to further obtain a more discriminative multi-scale feature representation.

The FPN introduces a top-down pathway to transfer the high-level semantic information into the shallow layers, while the low-level spatial information is still lost in the top layers after the long-distance propagation in the backbone. Drawing on this problem, Fu *et al.* [80] proposed a feature-fusion architecture (FFA) that integrates an auxiliary bottom-up pathway into the FPN structure to transfer the low-level spatial information to the top layers features via a short path, as depicted in Fig. 5(e). The FFA ensures the detector extracts multi-scale feature pyramids with rich semantic and detailed spatial information. Similarly, in [81], [82], the authors introduced a bidirectional FPN that learns the importance of different level features through learnable parameters and fuses the multi-level features through iteratively top-down and bottom-up pathways.

Different from the above sequential enhancement pathway [80], some studies [83]–[94] adopt a cross-level feature fusion manner. As shown in Fig. 5(f), the cross-level feature fusion methods fully collect the features at all levels to adaptively obtain balanced feature maps. Cheng *et al.* [83] utilized feature concatenation operation to achieve cross-scale feature fusion. Considering that features from different levels should have different contributions to the feature fusion, Fu *et al.* [84]
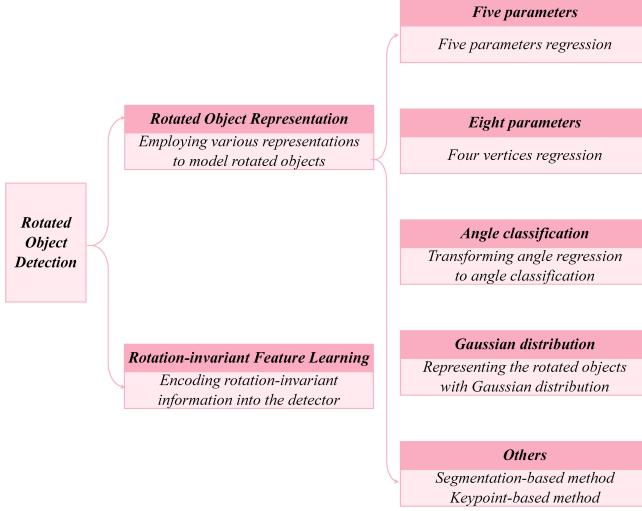
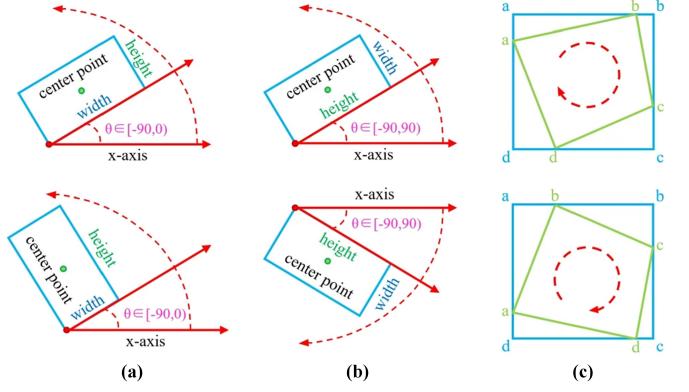Fig. 6. A brief summary of rotated object detection methods.



Fig. 7. Visualization of the five parameters representation and eight parameters representation methods for rotated objects [106].

## III. ROTATED OBJECT DETECTION

Arbitrary orientation of objects is another major challenge in RSOD. Since the objects in RSIs are acquired from a bird's eye view, they exhibit the property of arbitrary orientations, so the widely used horizontal bounding box (HBB) representation in generic object detection is insufficient to locate rotated objects accurately. Therefore, numerous researchers have focused on the arbitrary orientation property of geospatial objects, which can be summarized into rotated object representation and rotation-invariant feature learning. A brief summary of rotated object detection methods is depicted in Fig. 6.

### A. Rotated Object Representation

Rotated object representation is essential for RSOD to avoid redundant backgrounds and obtain precise detection results. Recent rotated object representation methods can be mainly summarized into several categories: five parameters representation [107]–[116], eight parameters representation [117]–[126], angle classification representation [106], [127]–[129], gaussian distribution representation [130]–[133], and others [134]–[144].

*1) Five Parameters:* The most popular solution is representing objects with a five-parameter method $(x, y, w, h, \theta)$, which simply adds an extra rotation angle parameter $\theta$ on HBB [107]–[115]. The definition of the angular range plays a crucial role in such methods, where two kinds of definitions are derived. Some studies [107]–[112] define $\theta$ as the acute angle to the x-axis and restrict the angular range to $90°$, as shown in Fig. 7(a). As the most representative work, Yang *et al.* [107] followed the five parameters method to detect rotated objects and designed an IoU-aware loss function to tackle the boundary discontinuity problem of rotation angles. Another group [113]–[116] refers to $\theta$ as the angle between the x-axis and the long side, whose range is $180°$, as depicted in Fig. 7(b). Ding *et al.* [114] regressed rotation angles by five-parameter methods and transformed the features of horizontal regions into rotated ones to facilitate rotated object detection.

*2) Eight Parameters:* Different from five-parameter methods, eight-parameter methods [117]–[126] solve the issue of

## C. Multi-scale Anchor Generation

Apart from the data augmentation and multi-scale feature representation methods, multi-scale anchor generation can also tackle the huge object scale variation problem in RSOD. Due to the difference in the scale range of objects in natural and remote sensing scenes, some studies [95]–[104] modify the anchor settings in common object detection to better cover the scales of geospatial objects.

Guo *et al.* [95] injected extra anchors with more scales and aspect ratios into the detector for multi-scale object detection. Dong *et al.* [98] designed more suitable anchor scales based on the statistics of object scales in the training set. Qiu *et al.* [99] extended the original square RoI features into vertical, square, and horizontal RoI features and fused these RoI features to represent objects with different aspect ratios in a more flexible way. The above methods follow a fixed anchor setting, while current studies [100]–[104] have attempted to dynamically learn the anchor during the training phase. Considering the aspect ratio variations between different categories, Hou *et al.* [100] devised a novel self-adaptive aspect ratio anchor (SARA) to adaptively learn an appropriate aspect ratio for each category. The SARA embeds the learnable category-wise aspect ratio values into the regression branch to adaptively update the aspect ratio of each category with the gradient of the location regression loss. Inspired by GA-RPN [105], some researchers [102]–[104] introduced a lightweight sub-network into the detector to adaptively learn the location and shape information of anchors.

proposed level-based attention to learn the unique contribution of features from each level. Thanks to the powerful global information extraction ability of the transformer structure, some work [88], [89] introduced the transformer structures to integrate and refine multi-level features. In [90], Chen *et al.* presented a cascading attention network where position supervision is introduced to enhance the semantic information of multi-level features.
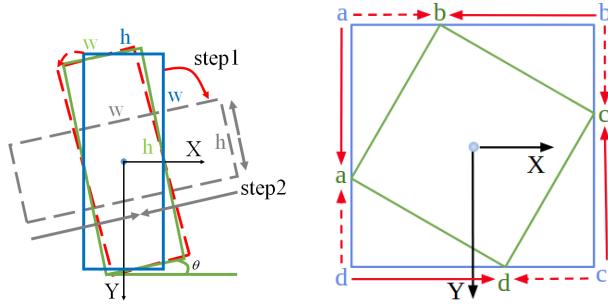
Fig. 8. Boundary discontinuity challenge of five parameters method and eight parameters method [119], [121].

rotated object representation by directly regressing four vertices $\{(a_x, a_y), (b_x, b_y), (c_x, c_y), (d_x, d_y)\}$, as shown in Fig. 7(c). Xia *et al.* [117] first adopted the eight-parameter method for rotated object representation, which directly supervises the detection model by minimizing the difference between each vertex and ground truth coordinates during training. However, the sequence order of these vertices is essential for the eight-parameter method to avoid unstable training. As shown in Fig. 8, it is intuitive that regressing targets from the red dotted arrow is an easier route, but the actual process follows the red solid arrows, which causes the difficulty of model training. To this end, Qian *et al.* [119], [121] proposed a modulated loss function that calculates the losses under different sorted orders and selects the minimum case to learn, efficiently improving the detection performance.

*3) Angle Classification:* To address the issue described in Fig. 8 from the source, many researchers [106], [127]–[129] take a detour from the boundary challenge of regression by transforming the angle prediction problem into an angle classification task. Yang *et al.* [106] proposed the first angle classification method for rotated object detection, which converts the continuous angle into a discrete kind and trains the model with novel circular smooth labels. However, the angle classification head [106] introduces additional parameters and degrades the detector's efficiency. To overcome this, Yang *et al.* [129] improved the [106] with a densely coded label that ensures both the accuracy and efficiency of the model.

*4) Gaussian Distribution:* Although the above methods achieve promising progress, they do not consider the misalignment between the actual detection performance and optimization metric. Most recently, a series of works [130]–[133] aims to handle this challenge by representing rotated objects with Gaussian distribution, as shown in Fig. 9. Specifically, these methods convert rotated objects into a 2D Gaussian distribution $\mathcal{N}(\mu, \Sigma)$ as follows:

$$\Sigma^{1/2} = \mathbf{R}\Lambda\mathbf{R}^\top$$

$$= \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix} \begin{pmatrix} \frac{w}{2} & 0 \\ 0 & \frac{h}{2} \end{pmatrix} \begin{pmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{pmatrix}$$

$$= \begin{pmatrix} \frac{w}{2}\cos^2\theta + \frac{h}{2}\sin^2\theta & \frac{w-h}{2}\cos\theta\sin\theta \\ \frac{w-h}{2}\cos\theta\sin\theta & \frac{w}{2}\sin^2\theta + \frac{h}{2}\cos^2\theta \end{pmatrix}$$
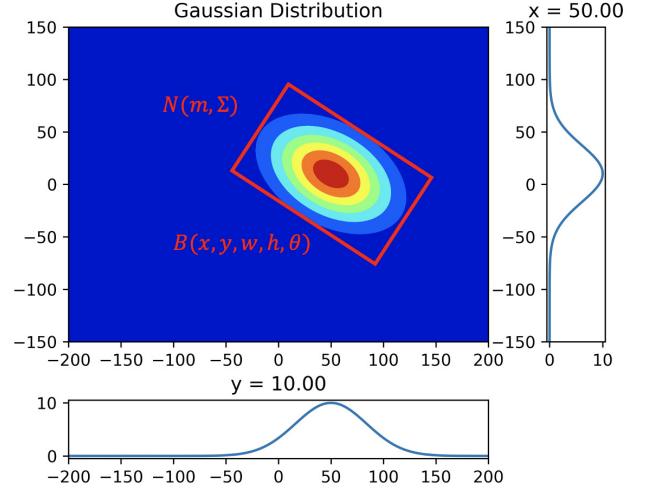
$$\mu = (x, y)^\top$$

(1)



Fig. 9. Visualization of the Gaussian distribution representation methods for rotated objects [130].

where $\mathbf{R}$ represents the rotation matrix, and $\Lambda$ represents the diagonal matrix of eigenvalues. With the Gaussian distribution representation in Eq. 1, the IoU between two rotated objects can be simplified as a distance estimation between two distributions. Besides, the Gaussian distribution representation discards the definition of angular boundary and effectively solves the angular boundary problem. Yang *et al.* [130] proposed a novel metric with Gaussian Wasserstein distance (GWD) for measuring the distance between distributions, which achieves remarkable performance by efficiently approximating the rotation IoU. Based on this, Yang *et al.* [131] introduced a Kullback-Leibler divergence (KLD) metric to enhance its scale invariance.

*5) Others:* Some researchers solve the rotated object representation by other approaches, such as segmentation-based [134]–[136] and keypoint-based [137]–[144]. The representative one in segmentation-based methods is Mask OBB [134], which deploys the segmentation method on each horizontal proposal to obtain the pixel-level object region and produce the minimum external rectangle as a rotated bounding box. On the other side, Wei *et al.* [142] adopted a keypoint-based representation for rotated objects, which locates the object center and leverages a pair of middle lines to represent the whole object. In addition, Yang *et al.* [145] proposed the first rotated object detector supervised by horizontal box annotations, which adopts the self-supervised learning of two different views to predict the angles of rotated objects.

### B. Rotation-invariant Feature Learning

Rotation-invariant features indicate the features remain consistent under any rotation transformations. Thus, rotation-invariant feature learning of objects is a crucial research field to tackle the arbitrary orientation problem in rotated object detection. To this purpose, many researchers proposed a series of methods for learning the rotational invariance of objects [146]–[157], which significantly improves rotated object detection in RSIs.

TABLE I
DETECTION PERFORMANCE OF ROTATED OBJECT DETECTION METHODS
ON THE DOTAV1.0 DATASET WITH ROTATED ANNOTATIONS.

| Models | Backbone | Methods | mAP(%) |
|---|---|---|---|
| SCRDet [107] | R-101-FPN | Five parameters | 72.61 |
| O$^2$Det [142] | H-104 | Keypoint-based | 72.8 |
| R$^3$Det [108] | R-101-FPN | Five parameters | 73.79 |
| S$^2$ANet [156] | R-50-FPN | Rotation-invariant feature | 74.12 |
| RoI Transformer [114] | R-50-FPN | Five parameters | 74.61 |
| Mask OBB [134] | R-50-FPN | Segmentation-based | 74.86 |
| Gliding Vertex [120] | R-101-FPN | Four vertices | 75.02 |
| DCL [128] | R-152-FPN | Angle classification | 75.54 |
| ReDet [155] | ReR50-ReFPN | Rotation-invariant feature | 76.25 |
| Oriented R-CNN [124] | R-101-FPN | Four vertices | 76.28 |
| R$^3$Det-KLD [131] | R-50-FPN | Gaussian distribution | 77.36 |

Cheng *et al.* [146] proposed the first rotation-invariant object detector to precisely recognize objects by using rotation-insensitive features, which enforces the features of objects to be consistent at different rotation angles. Later, Cheng *et al.* [148], [149] employed the rotation-invariant and fisher discrimination regularizers to encourage the detector to learn both rotation-invariant and discriminative features. In [150], [151], Wu *et al.* analyzed object rotation invariance under polar coordinates in the Fourier domain and designed a spatial-frequency channel feature extraction module to obtain the rotation-invariant features. Considering misalignment between axis-aligned convolutional features and rotated objects, Han *et al.* [156] proposed an oriented detection module that adopts a novel alignment convolution operation to learn the orientation information. In [155], Han *et al.* further devised a rotation-equivariant detector to explicitly encode rotation equivariance and rotation invariance. Besides, some researchers [80], [157] extended the RPN with a series of predefined rotated anchors to cope with the arbitrary orientation characteristics of geospatial objects.

We summarize the detection performance of milestone rotated object detection methods in Table I.

## IV. WEAK OBJECT DETECTION

Objects of interest in RSIs are typically embedded in complex scenes with intricate object spatial patterns and massive background noise. The complex context and background noise severely harm the feature representation of objects of interest, resulting in weak feature responses to objects of interest. Thus, many existing works have concentrated on improving the feature representation of the objects of interest, which can be grouped into two streams: suppressing background noise and mining related context information. A brief summary of weak object detection methods is given in Fig. 10.

### A. Suppressing background noise

This type of method aims to strengthen the weak response of the object region in the feature map by weakening the response of background regions. It can be mainly divided into two categories: implicit learning and explicit supervision.
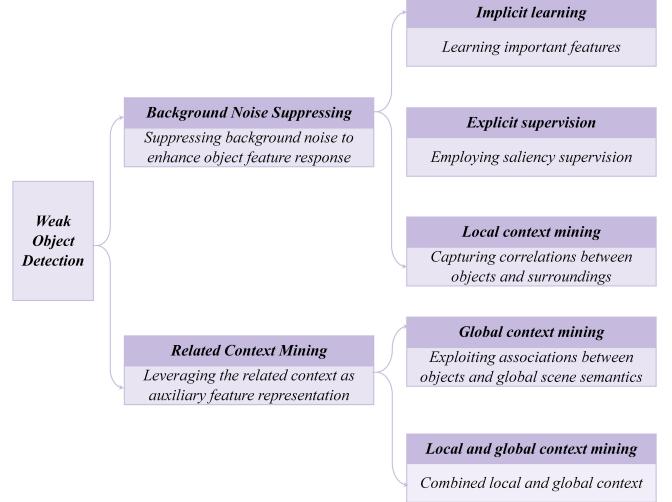


Fig. 10. A brief summary of weak object detection methods.

*1) Implicit Learning:* Implicit learning methods employ carefully designed modules into the detector to adaptively learn important features and suppress redundant features during the training phase, thereby reducing the background noise interference.

In machine learning, dimensionality reduction can effectively learn compact feature representation and suppress irrelevant features. Drawing on the above property, Ye *et al.* [158] proposed a feature filtration module that captures the low-dimensional feature maps by consecutive bottleneck layers to filter background noise interference. Inspired by the selective focus of human visual perception, the attention mechanism has been proposed and heavily researched [159]–[161]. The attention mechanism redistributes the feature importance during the network learning phase to enhance important features and suppress redundant information. Consequently, the attention mechanism has also been widely introduced in RSOD to tackle the background noise interference problem [57], [162]–[170]. In [162], Huang *et al.* emphasized the importance of patch-patch dependencies for RSOD and designed a novel non-local perceptual pyramidal attention (NP-Attention). The NP-Attention learns spatial multi-scale non-local dependencies and channel dependencies to enable the detector to concentrate on the object region rather than the background. Considering the strong scattering interference of the land area in SAR images, Sun *et al.* [163] presented a ship attention module to highlight the feature representation of ships and reduce the false alarm from the land area. Moreover, a series of attention mechanisms devised for RSOD (e.g., spatial shuffle-group enhance attention [165], multi-scale spatial and channel-wise attention [166], discrete wavelet multi-scale attention [167], etc.) have demonstrated their effectiveness in suppressing background noise.

*2) Explicit Supervision:* Unlike implicit learning methods, the explicit supervision approach employs auxiliary saliency supervision information to explicitly guide the detector to highlight the foreground regions and weaken the background.

Li *et al.* [171] employed the region contrast method to obtain the saliency map and construct the saliency feature
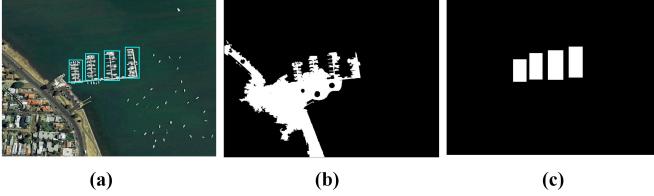
Fig. 11. (a) Input image. (b) Saliency map generated by the saliency detection method [173]. (c) Object-level saliency map.

pyramid by fusing the multi-scale feature maps with the saliency map. In [172], Lei *et al.* extracted the saliency map with the saliency detection method [173] and proposed a saliency reconstruction network. The saliency reconstruction network utilizes the saliency map as pixel-level supervision to guide the training of the detector to strengthen saliency regions in feature maps. The above saliency detection methods are typically unsupervised, and the generated saliency map may contain non-object regions, as shown in Fig. 11(b), providing inaccurate guidance to the detector. Therefore, later works [107], [134], [174]–[180] transformed the box-level annotation into object-level saliency guidance information (as shown in Fig. 11(c)) to generate more accurate saliency supervision. Yang *et al.* [107] designed a pixel attention network that employs object-level saliency supervision to enhance the object cues and weaken the background information. In [175], Zhang *et al.* proposed the FoRDet to exploit object-level saliency supervision in a more concise way. Concretely, the proposed FoRDet leverages the prediction of foreground regions in the coarse stage (supervised under box-level annotation) to enhance the feature representation of the foreground regions in the refined stage.

### B. Mining related context information

Context information typically refers to the spatial and semantic relations between the object and the surrounding environment or scene. This context information can provide auxiliary feature representations to the object that could not be clearly distinguished. Thus, mining context information can effectively solve the weak feature responses problem in RSOD. According to the category of context information, existing methods are mainly classified into local and global context information mining.

*1) Local Context Information Mining:* Local context information refers to the correlation between the object and its surrounding environment in visual information and spatial distribution [147], [181]–[187]. Zhang *et al.* [181] generated multiple local context regions by scaling the original region proposal into three different sizes and proposed a contextual bidirectional enhancement module to fuse the local context features with object features. The context-aware convolutional neural network (CA-CNN) [182] employed a Context-RoI mining layer to extract context information surrounding objects. The Context-RoI for an object is first generated by merging a series of filtered proposals around the object and then fused with the object RoI as the final object feature representation for classification and regression. In [183], Ma

*et al.* exploit the gated recurrent units (GRUs) to fuse object features with the local context information, leading to a more discriminative feature representation for the object. Graph Convolutional Networks (GCN) have recently shown better performance in object-object relationship reasoning. Hence, Tian *et al.* [184], [185] constructed spatial and semantic graphs to model and learn the contextual relationships between objects.

*2) Global Context Information Mining:* Global context information exploits the association between the object and the scene [188]–[195], e.g., vehicles generally locate on roads and ships typically appear at sea. Chen *et al.* [188] extracted the scene context information from the global image feature with the RoI-Align operation and fused it with the object-level RoI features to strengthen the relationship between the object and the scene. Liu *et al.* [192] designed a scene auxiliary detection head that exploits the scene context information under scene-level supervision. The scene auxiliary detection head embeds the predicted scene vector into the classification branch to fuse the object-level features with scene-level context information. In [193], Tao *et al.* presented a scene context-driven vehicle detection approach. Specifically, a pre-trained scene classifier is introduced to classify each image patch into three scene categories, then the scene-specific vehicle detectors are employed to achieve preliminary detection results, and finally the detection results are further optimized with the scene contextual information.

Considering the complementarity of local and global context information, Zhang *et al.* [196] proposed a CAD-Net to mine both local and global context information. CAD-Net employed a pyramid local context network to learn the object-level local context information and designed a global context network to extract scene-level global context information. In [103], Teng *et al.* proposed a GLNet to collect context information from global to local so as to achieve a robust and accurate detector for RSIs. Besides, some studies [197]–[199] also introduced the ASPP [200] or RFB module [54] to leverage both local and global context information.

## V. TINY OBJECT DETECTION

The typical ground sampling distance (GSD) for RSIs is 1-3 meters, which means that even large objects (e.g., airplanes, ships, and storage tanks) can only occupy less than $16 \times 16$ pixels. Besides, even in high-resolution RSIs with a GSD of 0.25m, a vehicle with a dimension of $3 \times 1.5m^2$ only covers 72 pixels ($12 \times 6$). This prevalence of tiny objects in RSIs further increases the difficulty of RSOD. Current studies on tiny object detection are mainly grouped into discriminative feature learning, super-resolution based methods, and improved detection metrics. The tiny object detection methods are briefly summarized in Fig. 12.

### A. Discriminative Feature Learning

The extremely small scales (less than $16 \times 16$ pixels) of tiny objects make them exhibit limited appearance information, which poses serious challenges for detectors to learn the features of tiny objects. To tackle the above problem, many
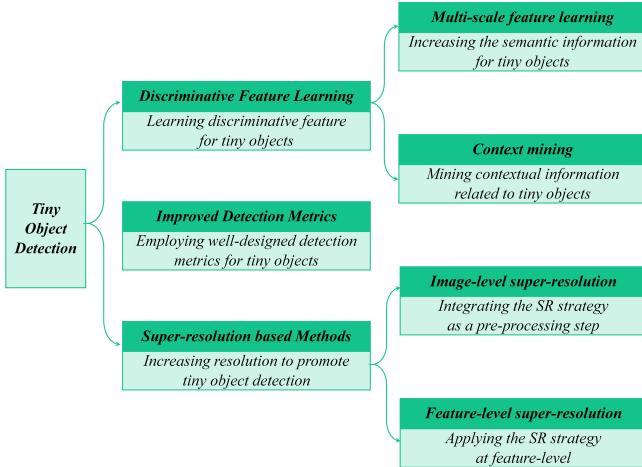
Fig. 12. A brief summary of tiny object detection methods.



Fig. 13. A comparison between (a) the IoU-Deviation Curve and (b) the NWD Deviation Curve [10]. Please refer to [10] for detail.

researchers focus on improving the discriminative feature learning ability for tiny objects [201]–[208].

Since tiny object mainly exists in shallow features and lacks high-level semantic information [65], some literature [201]–[203] introduces top-down structures to fuse high-level semantic information into shallow features to strengthen the semantic information for tiny objects. Considering the limited appearance information of the tiny objects, some studies [204]–[208] establish the connection between the tiny object and the surrounding contextual information through the self-attention mechanism or dilated convolution to enhance the feature discriminative of tiny objects. Notably, some previously mentioned studies on multi-scale feature learning and context information mining also demonstrate remarkable effectiveness in tiny object detection.

### B. Super-resolution based Method

The extremely small scale is a crucial issue for tiny object detection, so increasing the resolution of images is an intuitive solution to promote the detection performance of tiny objects. Some methods [209]–[212] employ the super-resolution strategies as a pre-processing step into the detection pipeline to enlarge the resolution of input images. For example, Rabbi *et al.* [211] emphasized the importance of edge information for tiny object detection and proposed an edge-enhanced super-resolution generative adversarial network (GAN) to generate visually pleasing high-resolution RSIs with detailed edge information. Wu *et al.* [212] developed a Point-to-Region detection framework for tiny objects. The Point-to-Region framework first obtains the proposal regions with key-point prediction and then employs a multi-task GAN to perform the super-resolution on the proposal regions and detect the tiny objects in these proposal regions. However, the high-resolution image generated by super-resolution brings extra computational complexity to the detection pipeline. Drawing on this problem, [213] and [214] employ the super-resolution strategy at the feature level to acquire discriminative feature representation for tiny objects and effectively save computational resources.
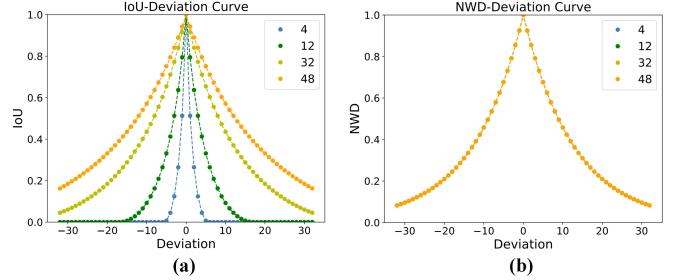
### C. Improved Detection Metrics for Tiny Object

Unlike the first two types of methods, recent advanced works [10], [215]–[222] assert that the current prevailing detection paradigms are unsuitable for tiny object detection and inevitably hinder tiny object detection performance. Pang *et al.* [215] argued that the excessive down-sampling operations in modern detectors leads to the loss of tiny objects on the feature map and proposed a zoom-out and zoom-in structure to enlarge the feature map. In [218], Yan *et al.* adjusted the IoU threshold in the label assignment to increase the positive assigned anchors for tiny objects, facilitating the learning of tiny objects. Dong *et al.* [219] devised the Sig-NMS to reduce the suppression of tiny objects by large and medium objects in traditional non-maximum suppression (NMS).

In [10], Xu *et al.* pointed out that the IoU metric is unsuitable for tiny object detection. As shown in Fig. 13, the IoU metric is sensitive to slight location offsets. Besides, the IoU-based label assignment suffers from a severe scale imbalance problem, where tiny objects tend to be assigned with insufficient positive samples. To solve these problems, Xu *et al.* [10] designed a normalized Wasserstein distance (NWD) to replace the IoU metric. The NWD models the tiny objects as 2D Gaussian distributions and utilizes the normalized Wasserstein distance between Gaussian distributions to represent the location relationship between tiny objects, detailed in [10]. Compared with the IoU metric, the proposed NWD metric is smooth to location deviations and has the characteristic of scale balance, as depicted in Fig. 13(b). In [222], Xu *et al.* further proposed the receptive field distance (RFLA) for tiny object detection and achieved state-of-the-art performance.

### VI. OBJECT DETECTION WITH LIMITED SUPERVISION

In recent years, the widely used deep learning based detectors in RSIs heavily rely on large-scale datasets with high-quality annotations to achieve state-of-the-art performance. However, collecting volumes of well-labeled data is considerably expensive and time-consuming (e.g., a bounding box annotation would cost about 10 seconds), which leads to a data-limited or annotation-limited scenario in RSOD [11]. This lack of sufficient supervised information seriously degrades the detection performance. To tackle this problem, researchers have explored various tasks in RSOD with limited supervision. We summarize the previous research into three main types: weakly-supervised object detection, semi-supervised object detection, and few-shot object detection. Fig. 14 provides
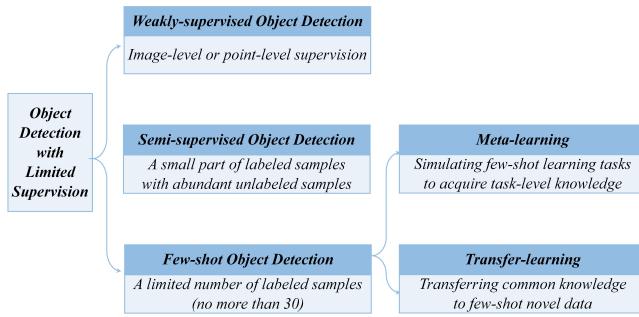
Fig. 14. A brief summary of object detection methods with limited supervision.



Fig. 15. The two-step paradigm of recent WSOD methods [229]–[241].

a brief summary of object detection methods with limited supervision.

## A. Weakly-supervised Object detection

Compared to fully supervised object detection, weakly-supervised object detection (WSOD) only contains weakly supervised information. Formally, WSOD consists of a training data set $\mathcal{D}_{train} = \{(X_i, y_i)\}_{i=1}^{I}$ , where $X_i = \{x_1, ..., x_{m_i}\}$ is a collection of training samples, termed as bag, $m_i$ is the total number of training samples in the bag, and $y_i$ is the weakly supervised information (e.g., image-level labels [223] or point-level labels [224]) of $X_i$. Effectively transferring image-level supervision to object-level labels is the key challenge in WSOD [225].

Han *et al.* [226] introduced a deep Boltzmann machine to learn the high-level features of objects and proposed a weakly-supervised learning framework based on Bayesian principles for remote sensing WSOD. Li *et al.* [227] exploited the mutual information between scene pairs to learn the discriminative convolutional weights and employed a multi-scale category activation map to locate geospatial objects.

Motivated by the remarkable performance of WSDDN [228], a series of remote sensing WSOD methods [229]–[241] are proposed. As shown in Fig. 15, the paradigm of the current WSOD methods usually consists of two steps, which first constructs a multiple instance learning model (MIL) to find contributing proposals to the image classification task as pseudo-labels and then utilizes them to train the detector. Yao *et al.* [229] introduced a dynamic curriculum learning strategy where the detector progressively improves detection performance through an easy-to-hard training process. Feng *et al.* [231] designed a progressive contextual instance refinement method that suppresses low-quality object parts and highlights the whole object by leveraging surrounding context information. Wang *et al.* [233] introduced the spatial and appearance relation graph into WSOD, which propagates high-quality label information to mine more possible objects. In [240], Feng *et al.* argued that existing remote sensing WSOD methods ignored the arbitrary orientations of geospatial objects, resulting in rotation-sensitive object detectors. To address this problem, Feng *et al.* [240] proposed a RINet, which brings rotation-invariant yet diverse feature learning for WSOD by employing rotation-invariant learning and multi-instance mining.
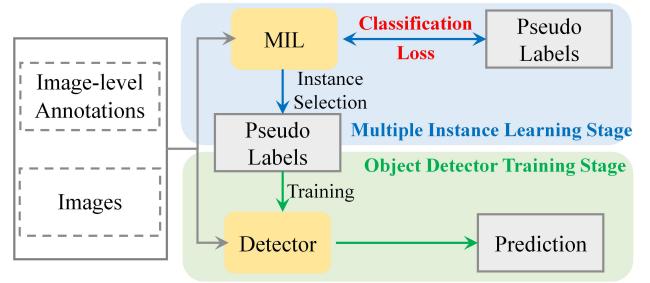
TABLE II
PERFORMANCE OF WEAKLY SUPERVISED OBJECT DETECTION METHODS ON THE NWPU VHR-10.V2 AND DIOR DATASETS.

| Models | NWPU VHR-10 | | DIOR | |
|---|---|---|---|---|
| | CorLoc(%) | mAP(%) | CorLoc(%) | mAP(%) |
| WSDDN [228] | 35.24 | 35.12 | 32.44 | 13.26 |
| DCL [229] | 69.65 | 52.11 | 42.23 | 20.19 |
| DPLG [230] | 61.50 | 53.60 | - | - |
| PCIR [231] | 71.87 | 54.97 | 46.12 | 24.92 |
| MIGL [233] | 70.16 | 55.95 | 46.80 | 25.11 |
| TCANet [234] | 72.76 | 58.82 | 48.41 | 25.82 |
| SAENet [235] | 73.46 | 60.72 | 49.42 | 27.10 |
| OS-DES [236] | 73.68 | 61.49 | 49.92 | 27.52 |
| SPG+MELM [239] | 73.41 | 62.80 | 48.30 | 25.77 |
| RINet [240] | - | 70.4 | 52.8 | 28.3 |
| MOL [241] | 75.96 | 75.46 | 50.66 | 29.21 |

We summarize the performance of milestone WSOD methods in Table II, where the correct localization metric (CorLoc) [242] is adopted to evaluate the localization performance.

## B. Semi-supervised Object detection

Semi-supervised Object detection (SSOD) typically contains only a small part (no more than 50%) of well-labeled samples $\mathcal{D}_{labeled} = \{(x_i, y_i)\}_{i=1}^{I_{labeled}}$, difficult to construct a reliable supervised detector, and has a large number of unlabeled samples $\mathcal{D}_{unlabeled} = \{(x_j)\}_{j=1}^{I_{unlabeled}}$. SSOD aims to improve detection performance under scarce supervised information by learning the latent information from volume unlabeled samples.

Hou *et al.* [243] proposed a SCLANet for semi-supervised SAR ship detection. The SCLANet employs adversarial learning between labeled and unlabeled samples to exploit the unlabeled sample information and adopts consistency learning for unlabeled samples to enhance the robustness of the network. The pseudo-label generation mechanism is also a widely used approach for semi-supervised object detection [244]–[248], and the typical paradigm is shown in Fig. 16. First, a pretrained detector learned from scare labeled samples are used to predict unlabeled samples, then the pseudo labels with higher confidence scores are selected as the trusted part, and finally, the model is retrained with the labeled and pseudo-labeled samples. Wu *et al.* [246] proposed a self-paced curriculum learning that follows an "easy to hard" scheme to select more reliable pseudo labels. Zhong *et al.* [245] adopt an active learning strategy in which high-scored predictions are
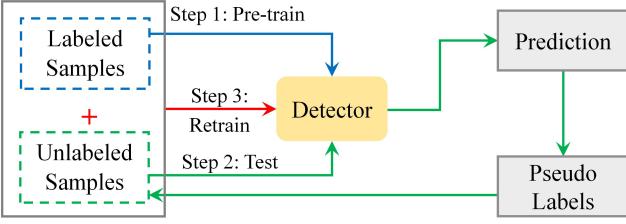
Fig. 16. The pipeline of pseudo-label generation mechanism in SSOD.



Fig. 17. The two-stage training pipeline of FSOD.

manually adjusted by experts to obtain refined pseudo labels. Chen *et al.* [247] employed teacher-student mutual learning to fully leverage unlabeled samples and iteratively generate higher-quality pseudo-labels.

In addition, some studies [249]–[253] have worked on weakly semi-supervised object detection, in which the unlabeled samples are replaced with weakly annotated samples. Du *et al.* [251], [252] employed a large number of image-level labeled samples to improve SAR vehicle detection performance under scarce box-level labeled samples. Chen *et al.* [253] adopted a small portion of pixel-level labeled samples and a dominant amount of box-level labeled samples to boost the performance in label-scarce instance segmentation.

### C. Few-shot Object Detection

Few-shot object detection (FSOD) refers to detecting novel classes with only a limited number (no more than 30) of samples. Generally, FSOD contains a base class dataset with abundant samples $\mathcal{D}_{\text{base}} = \{(x_i, y_i), y_i \in C_{base}\}_{i=1}^{I_{base}}$ and a novel class dataset with only $K$-shot samples $\mathcal{D}_{\text{novel}} = \{(x_j, y_j), y_j \in C_{novel}\}_{j=1}^{C_{novel}*K}$. Note that $C_{base}$ and $C_{novel}$ are disjointed. As depicted in Fig. 17, a typical FSOD paradigm consists of a two-stage training pipeline where the base training stage establishes prior knowledge with abundant base class samples, and the few-shot fine-tuning stage leverages the prior knowledge to facilitate the learning of few-shot novel concepts. The research on remote sensing FSOD mainly focuses on meta-learning methods [254]–[259] and transfer-learning methods [260]–[269].

The meta-learning based methods acquire task-level knowledge by simulating a series of few-shot learning tasks and generalize this knowledge to tackle few-shot learning of novel classes. Li *et al.* [255] first employed meta-learning for remote sensing FSOD and achieved satisfactory detection performance with only 1 to 10 labeled samples. Later, a series of meta-learning based few-shot detectors have been developed in the remote sensing community [254]–[259]. For example, Cheng *et al.* [254] proposed a Prototype-CNN to generate better foreground proposals and class-aware RoI features for remote sensing FSOD by learning class-specific prototypes. Wang *et al.* [258] presented a meta-metric training paradigm to enable the few-shot learner with flexible scalability for fast adaptation to the few-shot novel tasks.

Transfer-learning based methods aim at fine-tuning the common knowledge learned from the abundant annotated data to the few-shot novel data and typically consist of a base training stage and a few-shot fine-tuning stage. Huang *et al.*
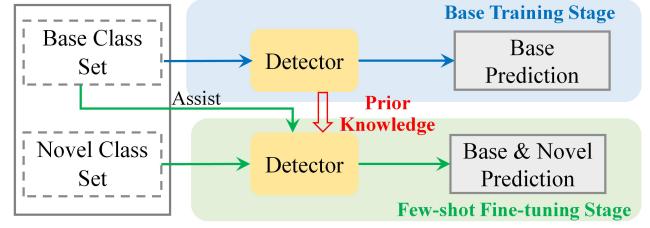
[266] proposed a balanced fine-tuning strategy to alleviate the number imbalance problem between novel class samples and base class samples. Zhou *et al.* [265] introduced proposal-level contrast learning in the fine-tuning phase to learn more robust feature representations in few-shot scenarios. Compared with the meta-learning based methods, the transfer-learning based method has a simpler and memory-efficient training paradigm.

## VII. DATASETS AND EVALUATION METRICS

### A. Datasets Introduction and Selection

Datasets have played an indispensable role throughout the development of object detection in RSIs. On the one hand, datasets serve as a common ground for the performance evaluation and comparison of detectors. On the other hand, datasets push researchers to address increasingly challenging problems in the RSOD field. In the past decade, several datasets with different attributes have been released to facilitate the development of RSOD, as shown in Table III. In this section, we mainly introduce 10 widely used datasets with specific characteristics.

**NWPU VHR-10** [18]. This dataset is a multi-class geospatial object detection dataset. It contains 3,775 HBB annotated instances in ten categories: airplane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, and vehicle. There are 800 very high-resolution RSIs, consisting of 715 color images from Google Earth and 85 pan-sharpened color infrared images from Vaihingen data. The image resolutions range from 0.5 to 2 m.

**VEDAI** [272]. VEDAI is a fine-grained vehicle detection dataset that contains five fine-grained vehicle categories: camping car, car, pick-up, tractor, truck, and van. There are 1,210 images and 3,700 instances in the VEDAI dataset, and the size of each image is $1,024 \times 1,024$. The small area and the arbitrary orientation of vehicles are the main challenges in the VEDAI dataset.

**UCAS-AOD** [274]. The UCAS-AOD dataset includes 910 images and 6,029 objects, where 3,210 aircraft are contained in 600 images and 2,819 vehicles are contained in 310 images. All images are acquired from Google Earth with an image size of approximately $1,000 \times 1,000$.

**HRSC** [276]. The HRSC dataset is widely used for arbitrary orientation ship detection and consists of 1,070 images and 2,976 instances with OBB annotation. The images are captured from Google Earth, containing offshore and inshore scenes. The image sizes vary from $300 \times 300$ and $1,500 \times 900$, and the image resolutions range from 2 to 0.4 m.

TABLE III
COMPARISONS OF WIDELY USED DATASETS IN THE FIELD OF RSOD. HBB AND OBB REFER TO *horizontal bounding box* AND *oriented bounding box*, RESPECTIVELY. * STANDS FOR *the average image width*.

| Dataset | Source | Annotation | Categories | Instances | Images | Image width | Resolution | Year |
|---|---|---|---|---|---|---|---|---|
| TAS [270] | Google Earth | HBB | 1 | 1,319 | 30 | 792 | - | 2008 |
| SZTAKI-INRIA [271] | Quick Bird, IKONOS and Google Earth | OBB | 1 | 665 | 9 | ∼800 | 0.5-1m | 2012 |
| NWPU VHR-10 [18] | Google Earth | HBB | 10 | 3,651 | 800 | ∼1,000 | 0.3-2m | 2014 |
| VEDAI [272] | Utah AGRC | OBB | 9 | 2,950 | 1,268 | 1,024 | 0.125m | 2015 |
| DLR 3k [273] | DLR 3K camera system | OBB | 8 | 14,235 | 20 | 5,616 | 0.13m | 2015 |
| UCAS-AOD [274] | Google Earth | OBB | 2 | 6,029 | 910 | ∼1,000 | 0.3-2m | 2015 |
| COWC [275] | Multiple Sources | Point | 1 | 32,716 | 53 | 2,000−19,000 | 0.15m | 2016 |
| HRSC [276] | Google Earth | OBB | 26 | 2,976 | 1,061 | ∼1,100 | 0.4-2m | 2016 |
| RSOD [43] | Google Earth and Tianditu | HBB | 4 | 6,950 | 976 | ∼1,000 | 0.3-3m | 2017 |
| SSDD [277] | RadarSat-2, TerraSARX and Sentinel-1 | HBB | 1 | 2,456 | 1,160 | 500 | 1-15m | 2017 |
| LEVIR [278] | Google Earth | HBB | 3 | 11,000 | 22,000 | 800−600 | 0.2-1m | 2018 |
| xView [2] | Worldview-3 | HBB | 60 | 1,000,000 | 1,413 | ∼3,000 | 0.3m | 2018 |
| DOTA-v1.0 [117] | Google Earth, JL-1, and GF-2 | HBB and OBB | 15 | 188,282 | 2,806 | 800−13,000 | 0.1-1m | 2018 |
| HRRSD [48] | Google Earth and Baidu Map | HBB | 13 | 55,740 | 21,761 | 152−10,569 | 0.15-1.2m | 2019 |
| DIOR [28] | Google Earth | HBB | 20 | 190,288 | 23,463 | 800 | 0.5-30m | 2019 |
| AIR-SARShip-1.0 [279] | Gaofen-3 | HBB | 1 | 3,000 | 31 | 3,000 | 1m and 3m | 2019 |
| MAR20 [280] | Google Earth | HBB and OBB | 20 | 22,341 | 3,824 | ∼800 | - | 2020 |
| FGSD [281] | Google Earth | OBB | 43 | 5,634 | 2,612 | 930 | 0.12-1.93m | 2020 |
| DOSR [282] | Google Earth | OBB | 20 | 6,172 | 1,066 | 600-1,300 | 0.5-2.5m | 2021 |
| AI-TOD [283] | Multiple Sources | HBB | 8 | 700,621 | 28,036 | 800 | - | 2021 |
| FAIR1M [34] | Gaofen satellites and Google Earth | OBB | 37 | 1,020,579 | 42,796 | 600-10,000 | 0.3-0.8m | 2021 |
| DOTA-v2.0 [33] | Google Earth, JL-1, GF-2 and airborne images | HBB and OBB | 18 | 1,793,658 | 11,268 | 800-20,000 | 0.1-4.5m | 2021 |
| SODA-A [284] | Google Earth | OBB | 9 | 800,203 | 2,510 | 4,761×2,777* | - | 2022 |

TABLE IV
DATASET SELECTION GUIDELINES IN RSOD FOR DIFFERENT CHALLENGES AND SCENARIOS.

| Scenarios | Datasets | Methods |
|---|---|---|
| Multi-scale Objects | DOTA, DIOR, FAIR1M | HyNet [51], FFA [80] |
| Rotated Objects | DOTA, HRSC | KLD [131], ReDet [155] |
| Weak Objects | DOTA, DIOR, FAIR1M | RECNN [172], CADNet [196] |
| Tiny Objects | SODA-A, AI-TOD | NWD [10], FSANet [216] |
| Weakly Supervision | NWPU VHR-10, DIOR | RINet [240], MOL [241] |
| Few-shot Supervision | NWPU VHR-10, DIOR | P-CNN [254], G-FSDet [269] |
| Fine-grained Objects | DOSR, FAIR1M | RBFPN [94], EIRNet [282] |
| SAR image Objects | SSDD, AIR-SARShip | SSPNet [163], HyperLiNet [285] |
| Specific Objects | HRSC, MAR20 | GRS-Det [135], COLOR [245] |

**SSDD** [277]. SSDD is the first open dataset for SAR image ship detection and contains 1,160 SAR images and 2,456 ships. The SAR images in the SSDD dataset are collected from different sensors with resolutions from 1m to 15 m and have different polarizations (HH, VV, VH, and HV). Subsequently, the author further refines and enriches the SSDD dataset into three different types to satisfy the current research of SAR ship detection [286].

**xView** [2]. The xView dataset is one of the largest publicly available datasets in ROSD, with approximately 1 million labeled objects across 60 fine-grained classes. Compared to other RSOD datasets, the images in xView dataset are collected from WorldView-3 at 0.3m ground sample distance, providing higher resolution images. Moreover, the xView dataset cover over 1,400 km$^2$ of the earth's surface, which leads to higher diversity.

**DOTA** [117]. DOTA is a large-scale dataset consisting of 188,282 objects annotated with both HBB and OBB. All objects are divided into 15 categories: plane, ship, storage tank, baseball diamond, tennis court, swimming pool, ground track field, harbor, bridge, large vehicle, small vehicle, helicopter, roundabout, soccer ball field, and basketball court. The images in this dataset are collected from Google Earth, JL-1 satellite, and GF-2 satellite with a spatial resolution of 0.1 to 1 m. Recently, the latest DOTAv2.0 [33] has been publicly available, which contains over 1.7 million objects of 18 categories.

**DIOR** [28]. DIOR is an object detection dataset for optical RSIs. There are 23,463 optical images in this dataset with a spatial resolution of 0.5 to 30m. The total number of objects in the dataset is 192,472, and all objects are labeled with HBB. The categories of objects are as follows: airplane, airport, baseball field, basketball court, bridge, chimney, dam, expressway service area, expressway toll station, harbor, golf course, ground track field, overpass, ship, stadium, storage tank, tennis court, train station, vehicle, and windmill.

**FAIR1M** [34]. FAIR1M is a more challenging dataset for fine-grained object detection in RSIs, including 5 categories and 37 subcategories. There are more than 40,000 images and more than 1 million objects annotated by oriented bounding boxes. The images are acquired from multiple platforms with a resolution of 0.3 m to 0.8 m and are spread across different countries and regions. The fine-grained categories, massive numbers of objects, large ranges of sizes and orientations, and diverse scenes make the FAIR1M more challenging.

**SODA-A** [284]. SODA-A is a recently released dataset designed for tiny object detection in RSIs. This dataset consists of 2,510 images with an average image size of $4,761 \times 2,777$, and 800,203 objects with OBB annotation. All objects are divided into four subsets (i.e., extremely small, relatively small, generally small, and normal) based on their area ranges. There are nine categories in this dataset, including airplane, he-
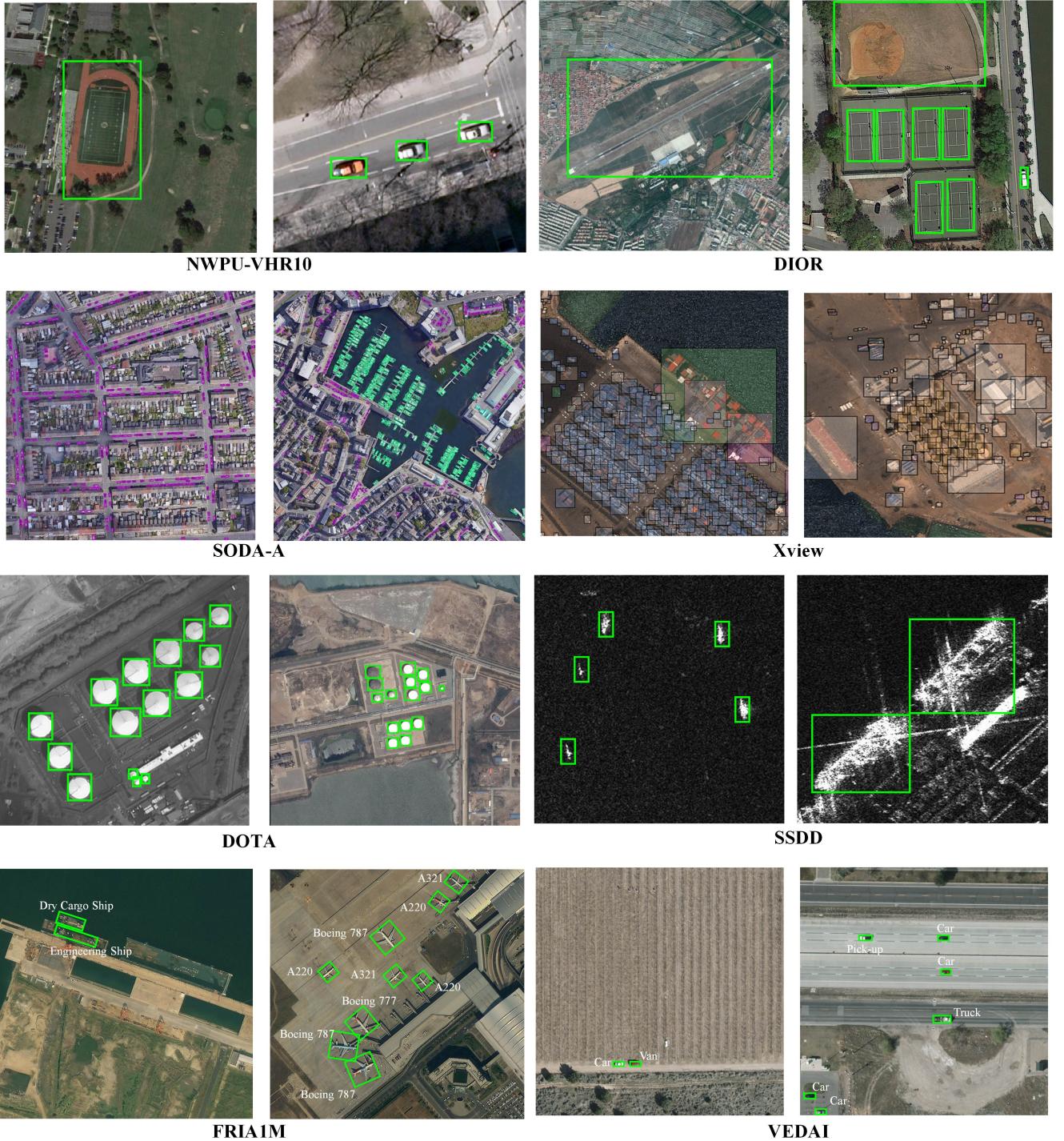
Fig. 18. Visualization of different RSOD datasets. Diverse resolutions, massive instances, multi-sensor images, and fine-grained categories are typical characteristics of RSOD datasets.

licopter, small-vehicle, large-vehicle, ship, container, storage-tank, swimming-pool, and windmill.

The above review shows that the early published datasets generally have limited samples. For example, NWPU VHR-10 [18] only contains 10 categories and 3,651 instances, and UCAC-AOD [274] consists of 2 categories with 6,029 instances. In recent years, researchers have not only introduced massive amounts of data and fine-grained level objects but also collected data from multi-sensor, various resolutions, and

diverse scenes (e.g., DOTA [117], DIOR [28], FAIR1M [34]) to satisfy the practical applications in RSOD. Fig. 18 depicts the typical samples of different RSOD datasets.

We also provide the dataset selection guidelines in Table IV to help researchers select proper datasets and methods for different challenges and scenarios. Notably, only the image-level annotations of the datasets are available for the weakly supervision scenario. As for the few-shot supervision scenario, there are only $K$-shot box-level annotated samples for each
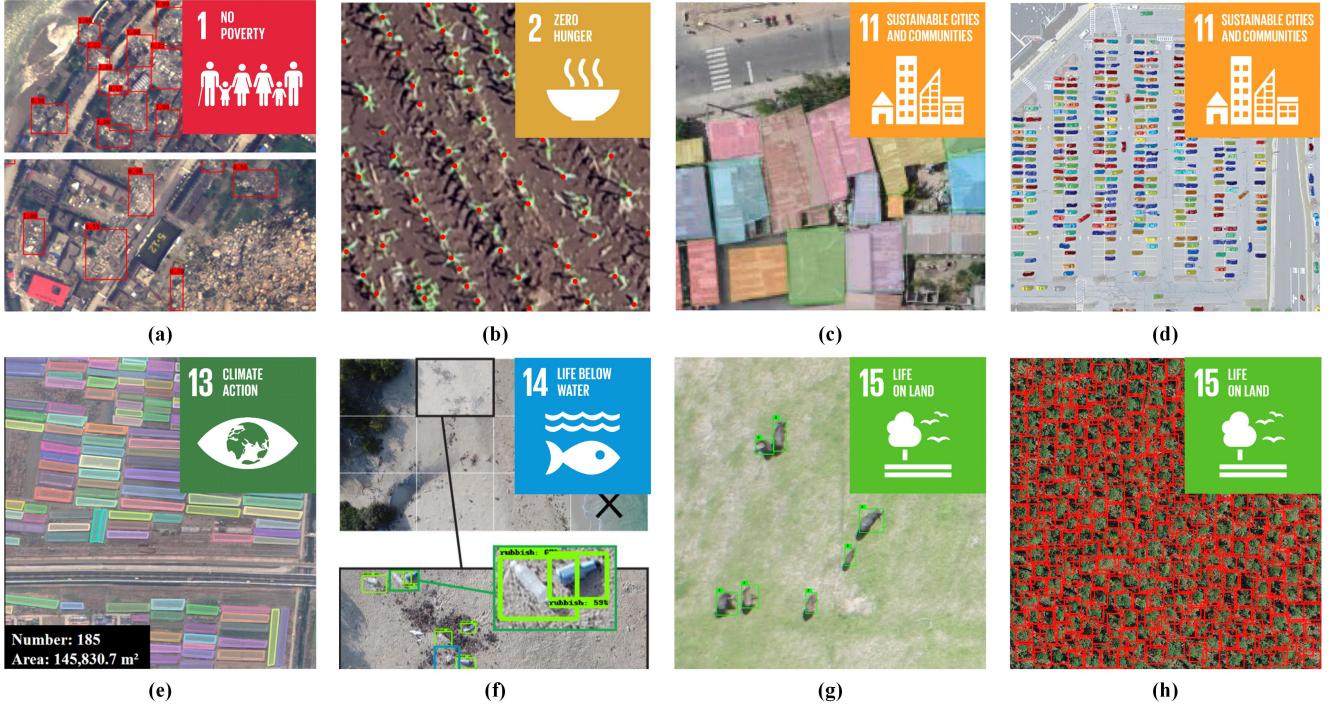
Fig. 19. The widespread applications of RSOD make substantial contributions to implementing of SDGs and improving society. (a) Collapsed buildings detection in Post-Earthquake for disaster assessment. (b) Corn plant detection for precision agriculture. (c-d) Building and vehicle detection for sustainable cities and communities. (e) Solar photovoltaic detection for climate change mitigation. (f) Litter detection along the shore for ocean conservation. (g) African mammals detection for wildlife surveillance. (h) Single tree detection for forest ecosystem protection.

novel class, where $K$ is set to $\{3, 5, 10, 20, 30\}$.

### B. Evaluation Metrics

In addition to the dataset, the evaluation metrics are equally important. Generally, the inference speed and the detection accuracy are the two commonly adopted metrics for evaluating the performance of detectors.

**Frames Per Second** (FPS) is a standard metric for inference speed evaluation that indicates the number of images that the detector can detect per second. Notably, both the image size and hardware devices can influence the inference speed.

**Average Precision** (AP) is the most commonly used metric for detection accuracy. Given a test image $I$, let $\{(b_i, c_i, p_i)\}_{i=1}^{N}$ denotes the prediction detections, where $b_i$ is the predicted box, $c_i$ is the predicted label, and $p_i$ is the confidence score. Let $\{(b_j^{gt}, c_j^{gt})\}_{j=1}^{M}$ refers to the ground truth annotations on the test image $I$, where $b_j^{gt}$ is the ground-truth box, and $c_j^{gt}$ is the ground truth category. A prediction detection $(b_i, c_i, p_i)$ is assigned as a True Positive (TP) for ground truth annotation $(b_j^{gt}, c_j^{gt})$, if it meets both of the following criteria:

- The confidence score $p_i$ is greater than the confidence threshold $t$, and the predicted label is same as the ground truth label $c_j^{gt}$.
- The IoU between the predicted box $b_i$ and the ground truth box $b_j^{gt}$ is larger than the IoU threshold $\varepsilon$. The IoU is calculated as follows:

$$\text{IoU}\,(b, b^g) = \frac{\text{area}\,(b \cap b^g)}{\text{area}\,(b \cup b^g)} \qquad (2)$$

where $area(b_i \cap b_j^{gt})$ and $area(b_i \cup b_j^{gt})$ stand for the intersection and union area of the predicted box and ground truth box.

Otherwise, it is considered to be a False Positive (FP). It is worth noting that multiple prediction detections may match the same ground truth annotation according to the above criteria, but only the prediction detection with the highest confidence score is assigned as a TP, and the rest are FPs [287].

Based on TP and FP detections, the Precision (P) and Recall (R) can be computed as Eq. 3 and Eq. 4.

$$P = \frac{TP}{TP + FP} \qquad (3)$$

$$R = \frac{TP}{TP + FN} \qquad (4)$$

where *FN* denotes the number of false negatives. The precision measures the fraction of true positives of the prediction detections and the recall measures the fraction of positives that are correctly detected. However, the above two evaluation metrics only reflect the single aspect of detection performance.

Taking into account both precision and recall, AP provides a comprehensive evaluation of detection performance and is calculated individually for each class. For a given class, the Precision/Recall Curve (PRC) is drawn according to the detection of maximum Precision at each Recall, and the AP summarises the shape of the PRC [287]. For multi-class object detection, the mean of the AP values for all classes, termed $mAP$, is adopted to evaluate the overall detection accuracy.

The early studies mainly employ a fixed IoU based AP metric (i.e., $AP_{50}$) [18], [28], [117], where the IoU threshold

$\varepsilon$ is given as 0.5. This low IoU threshold exhibits a high tolerance for bounding box deviations and fails to satisfy the high localization accuracy requirements. Later, some works [130], [131], [284] introduce a novel evaluation metric, named $AP_{50:95}$, which averages the AP over 10 IoU thresholds from 0.5 to 0.95 with an interval of 0.05. The $AP_{50:95}$ considers higher IoU thresholds and encourages more accurate localization.

As the cornerstone of evaluation metrics in RSOD, AP has various extensions for different specific tasks. In the few-shot learning scenario, $AP_{novel}$ and $AP_{base}$ are two critical metrics to evaluate the performance of few-shot detectors, where $AP_{novel}$ and $AP_{base}$ represent detection performance on the novel class and base class, respectively. An excellent few-shot detector should achieve satisfactory performance in the novel class and avoid performance degradation in the base class [269]. In the incremental detection of remote sensing objects, $AP_{old}$ and $AP_{inc}$ are employed to evaluate the performance of the old and incremental classes on different incremental tasks. In addition, the harmonic mean is also a vital evaluation metric for incremental object detection [288], which provides a comprehensive performance evaluation of both old and incremental classes, as described by Eq. 5:

$$HM = \frac{2AP_{old}AP_{inc}}{AP_{old} + AP_{inc}} \tag{5}$$

## VIII. APPLICATIONS

Deep learning techniques have injected significant innovations into RSOD, leading to an effective way to automatically identify objects of interest from voluminous RSIs. Therefore, RSOD methods have been applied in a rich diversity of practice scenarios that significantly support the implementation of Sustainable Development Goals (SDGs) and the improvement of society [289]–[291], as depicted in Fig. 19.

### A. Disaster Management

Natural disasters pose a serious threat to the safety of human life and property. A quick and precise understanding of disaster impact and extent of damage is critical to disaster management. RSOD methods can accurately identify ground objects from a bird's-eye view of the disaster-affected area, providing a novel potential for disaster management [292]–[296]. Guan *et al.* [293] proposed a novel instance segmentation model to accurately detect fire in a complex environment, which can be applied to the forest fire disaster response. Ma *et al.* [295] designed a real-time detection method for collapsed building assessment in Post-Earthquake.

### B. Precision Agriculture

With the unprecedented and still-expanding population, ensuring agricultural production is a fundamental obstacle to feeding the growing population. RSOD has the ability to monitor crop growth and estimate food production, promoting further progress for precision agriculture [297]–[302]. Pang *et al.* [298] used RSIs for early-season maize detection and achieved an accurate estimation of emergence rates. Chen *et al.* [302] designed an automatic strawberry flower detection system to monitor the growth cycle of strawberry fields.

### C. Sustainable Cities and Communities

Half of the global population now lives in cities, and this population will keep growing in the coming decades. Sustainable cities and communities are the goals of modern city development, in which RSOD can make a significant impact. For instance, building and vehicle detection [303]–[306] can help estimate population density distribution and transport traffic statistics, providing suggestions for city development planning. Infrastructure distribution detection [307] can assist in disaster assessment and early warning in the city environment.

### D. Climate Action

The ongoing climate change forces humans to face the daunting challenge of the climate crisis. Some researchers [308]–[310] employed object detection methods for automatically mapping tundra ice-wedge polygon to document and analyze the effects of climate warming on the Arctic region. Besides, RSOD can produce statistics on the number and spatial distribution of solar panels and wind turbines [311]–[314], facilitating the mitigation of greenhouse gas emissions.

### E. Ocean Conservation

The oceans cover nearly three-quarters of the Earth's surface, and more than 3 billion people depend on the diverse life of the oceans and coasts. The ocean is gradually deteriorating due to pollution, and the RSOD can provide powerful support for ocean conservation [315]. Several works applied detection methods for litter detection along shores [316], floating plastic detection at sea [317], deep-sea debris detection [318], etc. Another important application is ship detection [135], [136], which can help monitor illegal fishing activities.

### F. Wildlife Surveillance

A global loss of biodiversity is observed at all levels, and object detection in combination with RSIs provides a novel perspective for wildlife conservation [319]–[323]. Delplanque *et al.* [322] adopted the deep learning based detector for multi-species detection and identification of African mammals. Kellenberger *et al.* [323] designed a weakly supervised wildlife detection framework that only requires image-level labels to identify wildlife.

### G. Forest Ecosystem Protection

The forest ecosystem plays an important role in ecological protection, climate regulation, and carbon cycling. Understanding the condition of trees is essential for forest ecosystem protection [324]–[328]. Safonova *et al.* [326] analyzed the shape, texture, and color of the detected trees' crowns to determine their damage stage, providing a more efficient way to assess forest health. Sani-Mohammed *et al.* [328] utilized an instance segmentation approach to map the standing dead trees, which is imperative for forest ecosystem management and protection.

## IX. Future Directions

Apart from the five RSOD research topics mentioned in this survey, there is still much work to be done in this field. Therefore, we present a forward-looking discussion of future directions to further improve and enhance the detectors in remote sensing scenes.

### A. Unified detection framework for large-scale remote sensing images

Benefiting from the development of remote sensing technology, high-resolution large-scale RSIs (e.g., over $10,000 \times 10,000$ pixels) can be easily obtained. However, limited by the GPU memory, the current mainstream RSOD methods fail to directly perform object detection in large-scale RSIs but adopt a sliding window strategy, mainly including sliding window cropping, patch prediction, and results merging. On the one hand, this sliding window framework requires complex data pre-processing and post-processing, compared with the unified detection framework. On the other hand, the objects usually occupy a small area of the RSIs, and the invalid calculation of the massive backgrounds leads to increasing computation time and memory consumption. Some studies [215], [329], [330] proposed a coarse-to-fine detection framework for object detection in large-scale RSIs. This framework first locates the regions of interest by filtering out meaningless regions and then achieves accurate detection from these filtered regions.

### B. Detection with Multi-modal remote sensing images

Restricted by the sensor imaging mechanism, the detectors based on the single-modal RSIs often have detection performance deviations, which are difficult to meet in practical applications [331]. In contrast, the multi-modal RSIs from different sensors have their characteristics. For instance, hyperspectral images contain high spectral resolution and fine-grained spectral features, SAR images provide abundant texture information, and optical images exhibit high spatial resolution with rich detailed information. The integrated processing of multi-modal RSIs can improve the interpretation ability of the scene and obtain a more objective and comprehensive understanding of the geospatial objects [332]–[334], providing the possibility to further improve the detection performance of RSOD.

### C. Domain adaptation object detection in remote sensing images

Due to the diversity of remote sensing satellite sensors, resolutions, and bands, as well as the influence of weather conditions, seasons, and geospatial regions [6], RSIs collected from different satellites are generally drawn from similar but not identical distributions. Such distribution differences (also called the domain gap) severely restrict the generalization performance of the detector. Recent studies on domain adaptation object detection [335]–[338] have proposed to tackle the domain gap problem. However, these studies only focus on the domain adaptation detectors in the single-modal, while the cross-modal domain adaptation object detection (e.g., from optical images to SAR images [339], [340]) is a more challenging and worthwhile topic to investigate.

### D. Incremental detection of remote sensing objects

The real-world environment is dynamic and open, where the number of categories evolved over time. However, mainstream detectors require both old and new data to retrain the model when meeting new categories, resulting in high computational costs. Recently, incremental learning has been considered the most promising way to solve this problem, which can learn new knowledge without forgetting old knowledge with only new data [341]. Incremental learning has been preliminarily explored in the remote sensing community [342]–[345]. For example, Chen *et al.* [342] integrated knowledge distillation into FPN and detection heads to learn new concepts while maintaining the old ones. More thorough research is still needed in incremental RSOD to meet the dynamic learning task in practical application.

### E. Self-supervised pre-trained models for remote sensing scenes

Current RSOD methods are always initialized with the ImageNet [346] pre-trained weights. However, there is an inevitable domain gap between the natural and remote sensing scenes, probably limiting the performance of RSOD. Recently, the self-supervised pre-training approaches have received extensive attention and shown excellent performance in the classification and downstream tasks in the nature scenes. Benefiting from the rapid advances in remote sensing technology, the abundant remote sensing data [347], [348] also provide sufficient data support for self-supervised pre-training. Some researchers [349]–[353] have initially demonstrated the effectiveness of remote sensing pre-training on representative downstream tasks. Therefore, exploring the self-supervised pre-training models based on multi-source remote sensing data deserves further research.

### F. Compact and efficient object detection architectures

Most existing airborne and satellite-borne satellites require sending back the remote sensing data for interpretation, leading to additional resource overheads. Thus, it is essential to investigate compact and efficient detectors for airborne and satellite-borne platforms to reduce resource consumption in data transmission. Drawing on this demand, some researchers have proposed lightweight detectors through model design [285], [354], [355], network pruning [356], [357], and knowledge distillation [358]–[360]. However, these detectors still rely heavily on high-performance GPUs and cannot be deployed on airborne and satellite-borne satellites. Therefore, designing compact and efficient object detection architectures for limited resources scenarios remains challenging.

## X. Conclusion

Object detection has been a fundamental but challenging research topic in the remote sensing community. Thanks to the rapid development of deep learning techniques, RSOD has received considerable attention and gained remarkable achievements in the past decade. In this review, we present a systematic review and summarization of existing deep learning

based methods in RSOD. Firstly, we summarized the five main challenges in RSOD according to the characteristics of geospatial objects and categorized the methods into five streams: multi-scale object detection, rotated object detection, weak object detection, tiny object detection, and object detection with limited supervision. Then, we adopted a systematic hierarchical division to review and summarize the methods in each category. Next, we introduced the typical benchmark datasets, evaluation metrics, and practical applications in the RSOD field. Finally, considering the limitations of existing RSOD methods, we discussed some promising directions for further research.

Given this time of high-speed technical evolution in RSOD, we believe this survey can help researchers to achieve a more comprehensive understanding of the main topics in this field and to find potential directions for future research.

## REFERENCES

[1] N. Gorelick, M. Hancher, M. Dixon, S. Ilyushchenko, D. Thau, and R. Moore, "Google earth engine: Planetary-scale geospatial analysis for everyone," *Remote Sens. Environ.*, vol. 202, pp. 18–27, 2017.

[2] D. Lam, R. Kuzma, K. McGee, S. Dooley, M. Laielli, M. Klaric, Y. Bulatov, and B. McCord, "xview: Objects in context in overhead imagery," 2018. [Online]. Available: http://arxiv.org/abs/1802.07856

[3] Z. Li, H. Shen, H. Li, G. Xia, P. Gamba, and L. Zhang, "Multi-feature combined cloud and cloud shadow detection in gaofen-1 wide field of view imagery," *Remote Sens. Environ.*, vol. 191, pp. 342–358, 2017.

[4] S. Zhang, R. Wu, K. Xu, J. Wang, and W. Sun, "R-cnn-based ship detection from high resolution remote sensing imagery," *Remote Sens.*, vol. 11, no. 6, p. 631, 2019.

[5] Y. Wang, C. Wang, H. Zhang, Y. Dong, and S. Wei, "Automatic ship detection based on retinanet using multi-resolution gaofen-3 imagery," *Remote Sens.*, vol. 11, no. 5, p. 531, 2019.

[6] X. X. Zhu, D. Tuia, L. Mou, G.-S. Xia, L. Zhang, F. Xu, and F. Fraundorfer, "Deep learning in remote sensing: A comprehensive review and list of resources," *IEEE Geosci. Remote Sens. Mag.*, vol. 5, no. 4, pp. 8–36, 2017.

[7] L. Zhang, L. Zhang, and B. Du, "Deep learning for remote sensing data: A technical tutorial on the state of the art," *IEEE Geosci. Remote Sens. Mag.*, vol. 4, no. 2, pp. 22–40, 2016.

[8] L. Zhang and L. Zhang, "Artificial intelligence for remote sensing data analysis: A review of challenges and opportunities," *IEEE Geosci. Remote Sens. Mag.*, vol. 10, no. 2, pp. 270–294, 2022.

[9] W. Han, J. Chen, L. Wang, R. Feng, F. Li, L. Wu, T. Tian, and J. Yan, "Methods for small, weak object detection in optical high-resolution remote sensing images: A survey of advances and challenges," *IEEE Geosci. Remote Sens. Mag.*, vol. 9, no. 4, pp. 8–34, 2021.

[10] C. Xu, J. Wang, W. Yang, H. Yu, L. Yu, and G.-S. Xia, "Detecting tiny objects in aerial images: A normalized wasserstein distance and a new benchmark," *ISPRS J. Photogrammetry Remote Sens.*, vol. 190, pp. 79–93, 2022.

[11] J. Yue, L. Fang, P. Ghamisi, W. Xie, J. Li, J. Chanussot, and A. Plaza, "Optical remote sensing image understanding with weak supervision: Concepts, methods, and perspectives," *IEEE Geosci. Remote Sens. Mag.*, vol. 10, no. 2, pp. 250–269, 2022.

[12] C. Xu and H. Duan, "Artificial bee colony (abc) optimized edge potential function (epf) approach to target recognition for low-altitude aircraft," *Pattern Recognit. Lett.*, vol. 31, no. 13, pp. 1759–1772, 2010.

[13] X. Sun, H. Wang, and K. Fu, "Automatic detection of geospatial objects using taxonomic semantics," *IEEE Geosci. Remote Sens. Lett.*, vol. 7, no. 1, pp. 23–27, 2010.

[14] Y. Lin, H. He, Z. Yin, and F. Chen, "Rotation-invariant object detection in remote sensing images based on radial-gradient angle," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 4, pp. 746–750, 2015.

[15] H. Moon, R. Chellappa, and A. Rosenfeld, "Performance analysis of a simple vehicle detection algorithm," *Image Vis. Comput.*, vol. 20, no. 1, pp. 1–13, 2002.

[16] S. Leninisha and K. Vani, "Water flow based geometric active deformable model for road network," *ISPRS J. Photogrammetry Remote Sens.*, vol. 102, pp. 140–147, 2015.

[17] D. Chaudhuri and A. Samal, "An automatic bridge detection technique for multispectral images," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 9, pp. 2720–2727, 2008.

[18] G. Cheng, J. Han, P. Zhou, and L. Guo, "Multi-class geospatial object detection and geographic image classification based on collection of part detectors," *ISPRS J. Photogrammetry Remote Sens.*, vol. 98, pp. 119–132, 2014.

[19] L. Zhang, L. Zhang, D. Tao, and X. Huang, "Sparse transfer manifold embedding for hyperspectral target detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 2, pp. 1030–1043, 2013.

[20] J. Han, P. Zhou, D. Zhang, G. Cheng, L. Guo, Z. Liu, S. Bu, and J. Wu, "Efficient, simultaneous detection of multi-class geospatial targets based on visual saliency modeling and discriminative learning of sparse coding," *ISPRS J. Photogrammetry Remote Sens.*, vol. 89, pp. 37–48, 2014.

[21] H. Sun, X. Sun, H. Wang, Y. Li, and X. Li, "Automatic target detection in high-resolution remote sensing images using spatial sparse coding bag-of-words model," *IEEE Geosci. Remote Sens. Lett.*, vol. 9, no. 1, pp. 109–113, 2011.

[22] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[23] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," in *Proc. Annu. Conf. Neural Inf. Process. Syst*, 2015, pp. 91–99.

[24] J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 6517–6525.

[25] T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017, pp. 2999–3007.

[26] Z. Tian, C. Shen, H. Chen, and T. He, "FCOS: fully convolutional one-stage object detection," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 9626–9635.

[27] L. Liu, W. Ouyang, X. Wang, P. W. Fieguth, J. Chen, X. Liu, and M. Pietikäinen, "Deep learning for generic object detection: A survey," *Int. J. Comput. Vis.*, vol. 128, no. 2, pp. 261–318, 2020.

[28] K. Li, G. Wan, G. Cheng, L. Meng, and J. Han, "Object detection in optical remote sensing images: A survey and a new benchmark," *ISPRS J. Photogrammetry Remote Sens.*, vol. 159, pp. 296–307, 2020.

[29] G. Cheng and J. Han, "A survey on object detection in optical remote sensing images," *ISPRS J. Photogrammetry Remote Sens.*, vol. 117, pp. 11–28, 2016.

[30] U. Alganci, M. Soydas, and E. Sertel, "Comparative research on deep learning approaches for airplane detection from very high-resolution satellite images," *Remote Sens.*, vol. 12, no. 3, p. 458, 2020.

[31] Z. Li, Y. Wang, N. Zhang, Y. Zhang, Z. Zhao, D. Xu, G. Ben, and Y. Gao, "Deep learning-based object detection techniques for remote sensing images: A survey," *Remote Sens.*, vol. 14, no. 10, p. 2385, 2022.

[32] J. Kang, S. Tariq, H. Oh, and S. S. Woo, "A survey of deep learning-based object detection methods and datasets for overhead imagery," *IEEE Access*, vol. 10, pp. 20 118–20 134, 2022.

[33] J. Ding, N. Xue, G. Xia, X. Bai, W. Yang, M. Y. Yang, S. J. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang, "Object detection in aerial images: A large-scale benchmark and challenges," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 11, pp. 7778–7796, 2022.

[34] X. Sun, P. Wang, Z. Yan, F. Xu, R. Wang, W. Diao, J. Chen, J. Li, Y. Feng, T. Xu, M. Weinmann, S. Hinz, C. Wang, and K. Fu, "Fair1m: A benchmark dataset for fine-grained object recognition in high-resolution remote sensing imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 184, pp. 116–130, 2022.

[35] W. Zhao, W. Ma, L. Jiao, P. Chen, S. Yang, and B. Hou, "Multi-scale image block-level F-CNN for remote sensing images object detection," *IEEE Access*, vol. 7, pp. 43 607–43 621, 2019.

[36] S. M. Azimi, E. Vig, R. Bahmanyar, M. Körner, and P. Reinartz, "Towards multi-class object detection in unconstrained remote sensing imagery," in *Asian Conference on Computer Vision*, vol. 11363, 2018, pp. 150–165.

[37] P. Shamsolmoali, M. Zareapoor, J. Chanussot, H. Zhou, and J. Yang, "Rotation equivariant feature image pyramid network for object detection in optical remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3112481

[38] Y. Chen, P. Zhang, Z. Li, Y. Li, X. Zhang, G. Meng, S. Xiang, J. Sun, and J. Jia, "Stitcher: Feedback-driven data provider for object detection," 2020. [Online]. Available: https://arxiv.org/abs/2004.12432

[39] X. Xu, X. Zhang, and T. Zhang, "Lite-yolov5: A lightweight deep learning detector for on-board ship detection in large-scene sentinel-1 SAR images," *Remote Sens.*, vol. 14, no. 4, p. 1018, 2022. [Online]. Available: https://doi.org/10.3390/rs14041018

[40] N. Su, Z. Huang, Y. Yan, C. Zhao, and S. Zhou, "Detect larger at once: Large-area remote-sensing image arbitrary-oriented ship detection," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2022.3144485

[41] B. Zhao, Y. Wu, X. Guan, L. Gao, and B. Zhang, "An improved aggregated-mosaic method for the sparse object detection of remote sensing imagery," *Remote Sens.*, vol. 13, no. 13, p. 2602, 2021.

[42] X. Han, Y. Zhong, and L. Zhang, "An efficient and robust integrated geospatial object detection framework for high spatial resolution remote sensing imagery," *Remote Sens.*, vol. 9, no. 7, p. 666, 2017.

[43] Y. Long, Y. Gong, Z. Xiao, and Q. Liu, "Accurate object localization in remote sensing images based on convolutional neural networks," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 5, pp. 2486–2498, 2017.

[44] Y. Zhong, X. Han, and L. Zhang, "Multi-class geospatial object detection based on a position-sensitive balancing framework for high spatial resolution remote sensing imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 138, pp. 281–294, 2018.

[45] P. Ding, Y. Zhang, W.-J. Deng, P. Jia, and A. Kuijper, "A light and faster regional convolutional neural network for object detection in optical remote sensing images," *ISPRS J. Photogrammetry Remote Sens.*, vol. 141, pp. 208–218, 2018.

[46] W. Liu, L. Ma, and H. Chen, "Arbitrary-oriented ship detection framework in optical remote-sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 15, no. 6, pp. 937–941, 2018.

[47] W. Liu, L. Ma, J. Wang, and H. Chen, "Detection of multiclass objects in optical remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 5, pp. 791–795, 2019.

[48] Y. Zhang, Y. Yuan, Y. Feng, and X. Lu, "Hierarchical and robust convolutional neural network for very high-resolution remote sensing object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 8, pp. 5535–5548, 2019.

[49] Z. Lin, K. Ji, X. Leng, and G. Kuang, "Squeeze and excitation rank faster R-CNN for ship detection in SAR images," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 5, pp. 751–755, 2019.

[50] 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 J. Photogrammetry Remote Sens.*, vol. 145, pp. 3–22, 2018.

[51] Z. Zheng, Y. Zhong, A. Ma, X. Han, J. Zhao, Y. Liu, and L. Zhang, "Hynet: Hyper-scale object detection network framework for multiple spatial resolution remote sensing imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 166, pp. 1–14, 2020.

[52] Y. Ren, C. Zhu, and S. Xiao, "Deformable faster r-cnn with aggregating multi-layer features for partially occluded object detection in optical remote sensing images," *Remote Sens.*, vol. 10, no. 9, p. 1470, 2018.

[53] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *in Proc. Euro. Conf. Comput. Vis.* Springer, 2016, pp. 21–37.

[54] S. Liu, D. Huang, and Y. Wang, "Receptive field block net for accurate and fast object detection," in *in Proc. Euro. Conf. Comput. Vis.*, 2018, pp. 385–400.

[55] Z. Shen, Z. Liu, J. Li, Y.-G. Jiang, Y. Chen, and X. Xue, "Dsod: Learning deeply supervised object detectors from scratch," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2017, pp. 1919–1927.

[56] Z. Zhang, S. Qiao, C. Xie, W. Shen, B. Wang, and A. L. Yuille, "Single-shot object detection with enriched semantics," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 5813–5821.

[57] X. Lu, J. Ji, Z. Xing, and Q. Miao, "Attention and feature fusion SSD for remote sensing object detection," *IEEE Trans. Instrum. Meas.*, vol. 70, pp. 1–9, 2021.

[58] G. Wang, Y. Zhuang, H. Chen, X. Liu, T. Zhang, L. Li, S. Dong, and Q. Sang, "Fsod-net: Full-scale object detection from optical remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–18, 2022.

[59] B. Hou, Z. Ren, W. Zhao, Q. Wu, and L. Jiao, "Object detection in high-resolution panchromatic images using deep models and spatial template matching," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 2, pp. 956–970, 2020.

[60] X. Liang, J. Zhang, L. Zhuo, Y. Li, and Q. Tian, "Small object detection in unmanned aerial vehicle images using feature fusion and scaling-based single shot detector with spatial context analysis," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 30, no. 6, pp. 1758–1770, 2020.

[61] Z. Wang, L. Du, J. Mao, B. Liu, and D. Yang, "Sar target detection based on ssd with data augmentation and transfer learning," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 1, pp. 150–154, 2018.

[62] S. Bao, X. Zhong, R. Zhu, X. Zhang, Z. Li, and M. Li, "Single shot anchor refinement network for oriented object detection in optical remote sensing imagery," *IEEE Access*, vol. 7, pp. 87 150–87 161, 2019.

[63] T. Xu, X. Sun, W. Diao, L. Zhao, K. Fu, and H. Wang, "ASSD: feature aligned single-shot detection for multiscale objects in aerial imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–17, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3089170

[64] Q. Li, L. Mou, Q. Liu, Y. Wang, and X. X. Zhu, "Hsf-net: Multiscale deep feature embedding for ship detection in optical remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 12, pp. 7147–7161, 2018.

[65] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 2117–2125.

[66] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2018, pp. 8759–8768.

[67] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, "Libra R-CNN: towards balanced learning for object detection," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 821–830.

[68] M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020, pp. 10 781–10 790.

[69] L. Hou, K. Lu, and J. Xue, "Refined one-stage oriented object detection method for remote sensing images," *IEEE Trans. Image Process.*, vol. 31, pp. 1545–1558, 2022.

[70] W. Zhang, L. Jiao, Y. Li, Z. Huang, and H. Wang, "Laplacian feature pyramid network for object detection in VHR optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3072488

[71] S. Wei, H. Su, J. Ming, C. Wang, M. Yan, D. Kumar, J. Shi, and X. Zhang, "Precise and robust ship detection for high-resolution SAR imagery based on hr-sdnet," *Remote Sens.*, vol. 12, no. 1, p. 167, 2020.

[72] G. Cheng, M. He, H. Hong, X. Yao, X. Qian, and L. Guo, "Guiding clean features for object detection in remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022.

[73] J. Jiao, Y. Zhang, H. Sun, X. Yang, X. Gao, W. Hong, K. Fu, and X. Sun, "A densely connected end-to-end neural network for multiscale and multiscene SAR ship detection," *IEEE Access*, vol. 6, pp. 20 881–20 892, 2018.

[74] Q. Guo, H. Wang, and F. Xu, "Scattering enhanced attention pyramid network for aircraft detection in SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 9, pp. 7570–7587, 2021.

[75] Y. Li, Q. Huang, X. Pei, L. Jiao, and R. Shang, "Radet: Refine feature pyramid network and multi-layer attention network for arbitrary-oriented object detection of remote sensing images," *Remote Sens.*, vol. 12, no. 3, p. 389, 2020.

[76] L. Shi, L. Kuang, X. Xu, B. Pan, and Z. Shi, "Canet: Centerness-aware network for object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, 2022.

[77] R. Yang, Z. Pan, X. Jia, L. Zhang, and Y. Deng, "A novel cnn-based detector for ship detection based on rotatable bounding box in SAR images," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 1938–1958, 2021.

[78] Y. Zhao, L. Zhao, B. Xiong, and G. Kuang, "Attention receptive pyramid network for ship detection in SAR images," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 13, pp. 2738–2756, 2020.

[79] X. Yang, X. Zhang, N. Wang, and X. Gao, "A robust one-stage detector for multiscale ship detection with complex background in massive SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–12, 2022.

[80] K. Fu, Z. Chang, Y. Zhang, G. Xu, K. Zhang, and X. Sun, "Rotation-aware and multi-scale convolutional neural network for object detection in remote sensing images," *ISPRS J. Photogrammetry Remote Sens.*, vol. 161, pp. 294–308, 2020.

[81] W. Huang, G. Li, B. Jin, Q. Chen, J. Yin, and L. Huang, "Scenario context-aware-based bidirectional feature pyramid network for remote sensing target detection," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2021.3135935

[82] V. Chalavadi, J. Prudviraj, R. Datla, C. S. Babu, and K. M. C, "msodanet: A network for multi-scale object detection in aerial images using hierarchical dilated convolutions," *Pattern Recognit.*, vol. 126, p. 108548, 2022.

[83] G. Cheng, Y. Si, H. Hong, X. Yao, and L. Guo, "Cross-scale feature fusion for object detection in optical remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 18, no. 3, pp. 431–435, 2021.

[84] J. Fu, X. Sun, Z. Wang, and K. Fu, "An anchor-free method based on feature balancing and refinement network for multiscale ship detection in SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 2, pp. 1331–1344, 2021.

[85] Y. Liu, Q. Li, Y. Yuan, Q. Du, and Q. Wang, "Abnet: Adaptive balanced network for multiscale object detection in remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3133956

[86] H. Guo, X. Yang, N. Wang, B. Song, and X. Gao, "A rotational libra R-CNN method for ship detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 8, pp. 5772–5781, 2020.

[87] T. Zhang, Y. Zhuang, G. Wang, S. Dong, H. Chen, and L. Li, "Multiscale semantic fusion-guided fractal convolutional object detection network for optical remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–20, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3108476

[88] Y. Zheng, P. Sun, Z. Zhou, W. Xu, and Q. Ren, "Adt-det: Adaptive dynamic refined single-stage transformer detector for arbitrary-oriented object detection in satellite optical imagery," *Remote Sens.*, vol. 13, no. 13, p. 2623, 2021.

[89] Z. Wei, D. Liang, D. Zhang, L. Zhang, Q. Geng, M. Wei, and H. Zhou, "Learning calibrated-guidance for object detection in aerial images," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 15, pp. 2721–2733, 2022.

[90] L. Chen, C. Liu, F. Chang, S. Li, and Z. Nie, "Adaptive multi-level feature fusion and attention-based network for arbitrary-oriented object detection in remote sensing imagery," *Neurocomputing*, vol. 451, pp. 67–80, 2021.

[91] X. Sun, P. Wang, C. Wang, Y. Liu, and K. Fu, "Pbnet: Part-based convolutional neural network for complex composite object detection in remote sensing imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 173, pp. 50–65, 2021.

[92] T. Zhang, X. Zhang, C. Liu, J. Shi, S. Wei, I. Ahmad, X. Zhan, Y. Zhou, D. Pan, J. Li, and H. Su, "Balance learning for ship detection from synthetic aperture radar remote sensing imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 182, pp. 190–207, 2021.

[93] T. Zhang, X. Zhang, and X. Ke, "Quad-fpn: A novel quad feature pyramid network for sar ship detection," *Remote Sens.*, vol. 13, no. 14, p. 2771, 2021.

[94] J. Song, L. Miao, Q. Ming, Z. Zhou, and Y. Dong, "Fine-grained object detection in remote sensing images via adaptive label assignment and refined-balanced feature pyramid network," *IEEE J. Sel. Top. Appl. Earth Obs. Remote. Sens.*, vol. 16, pp. 71–82, 2023.

[95] W. Guo, W. Yang, H. Zhang, and G. Hua, "Geospatial object detection in high resolution satellite images based on multi-scale convolutional neural network," *Remote Sens.*, vol. 10, no. 1, p. 131, 2018.

[96] S. Zhang, G. He, H. Chen, N. Jing, and Q. Wang, "Scale adaptive proposal network for object detection in remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 6, pp. 864–868, 2019.

[97] C. Li, C. Xu, Z. Cui, D. Wang, T. Zhang, and J. Yang, "Feature-attentioned object detection in remote sensing imagery," in *Proc. IEEE Int. Conf. Image Process. Conf. (ICIP)*. IEEE, 2019, pp. 3886–3890.

[98] Z. Dong, M. Wang, Y. Wang, Y. Zhu, and Z. Zhang, "Object detection in high resolution remote sensing imagery based on convolutional neural networks with suitable object scale features," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 3, pp. 2104–2114, 2020.

[99] H. Qiu, H. Li, Q. Wu, F. Meng, K. N. Ngan, and H. Shi, "A$^2$rmnet: Adaptively aspect ratio multi-scale network for object detection in remote sensing images," *Remote Sens.*, vol. 11, no. 13, p. 1594, 2019.

[100] J. Hou, X. Zhu, and X. Yin, "Self-adaptive aspect ratio anchor for oriented object detection in remote sensing images," *Remote Sens.*, vol. 13, no. 7, p. 1318, 2021.

[101] N. Mo, L. Yan, R. Zhu, and H. Xie, "Class-specific anchor based and context-guided multi-class object detection in high resolution remote sensing imagery with a convolutional neural network," *Remote Sens.*, vol. 11, no. 3, p. 272, 2019.

[102] Z. Tian, R. Zhan, J. Hu, W. Wang, Z. He, and Z. Zhuang, "Generating anchor boxes based on attention mechanism for object detection in remote sensing images," *Remote Sens.*, vol. 12, no. 15, p. 2416, 2020.

[103] Z. Teng, Y. Duan, Y. Liu, B. Zhang, and J. Fan, "Global to local: Clip-lstm-based object detection from remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3064840

[104] Y. Yu, H. Guan, D. Li, T. Gu, E. Tang, and A. Li, "Orientation guided anchoring for geospatial object detection from remote sensing imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 160, pp. 67–82, 2020.

[105] J. Wang, K. Chen, S. Yang, C. C. Loy, and D. Lin, "Region proposal by guided anchoring," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 2965–2974.

[106] X. Yang and J. Yan, "On the arbitrary-oriented object detection: Classification based approaches revisited," *Int. J. Comput. Vis.*, vol. 130, no. 5, pp. 1340–1365, 2022.

[107] X. Yang, J. Yang, J. Yan, Y. Zhang, T. Zhang, Z. Guo, X. Sun, and K. Fu, "Scrdet: Towards more robust detection for small, cluttered and rotated objects," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 8232–8241.

[108] X. Yang, J. Yan, Z. Feng, and T. He, "R3det: Refined single-stage detector with feature refinement for rotating object," in *Pro. AAAI Conf. Artific. Intell.*, vol. 35, no. 4, 2021, pp. 3163–3171.

[109] X. Yang, H. Sun, K. Fu, J. Yang, X. Sun, M. Yan, and Z. Guo, "Automatic ship detection in remote sensing images from google earth of complex scenes based on multiscale rotation dense feature pyramid networks," *Remote Sens.*, vol. 10, no. 1, p. 132, 2018.

[110] X. Yang, H. Sun, X. Sun, M. Yan, Z. Guo, and K. Fu, "Position detection and direction prediction for arbitrary-oriented ships via multitask rotation region convolutional neural network," *IEEE Access*, vol. 6, pp. 50 839–50 849, 2018.

[111] Q. Ming, L. Miao, Z. Zhou, and Y. Dong, "Cfc-net: A critical feature capturing network for arbitrary-oriented object detection in remote-sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3095186

[112] Q. Ming, Z. Zhou, L. Miao, H. Zhang, and L. Li, "Dynamic anchor learning for arbitrary-oriented object detection," in *Pro. AAAI Conf. Artific. Intell.*, 2021, pp. 2355–2363.

[113] Y. Zhu, J. Du, and X. Wu, "Adaptive period embedding for representing oriented objects in aerial images," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 10, pp. 7247–7257, 2020.

[114] J. Ding, N. Xue, Y. Long, G.-S. Xia, and Q. Lu, "Learning roi transformer for oriented object detection in aerial images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 2849–2858.

[115] Q. An, Z. Pan, L. Liu, and H. You, "Drbox-v2: An improved detector with rotatable boxes for target detection in SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 11, pp. 8333–8349, 2019.

[116] Q. Li, L. Mou, Q. Xu, Y. Zhang, and X. X. Zhu, "R$^3$-net: A deep network for multi-oriented vehicle detection in aerial images and videos," 2018. [Online]. Available: http://arxiv.org/abs/1808.05560

[117] G. Xia, X. Bai, J. Ding, Z. Zhu, S. J. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang, "DOTA: A large-scale dataset for object detection in aerial images," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 3974–3983.

[118] Y. Liu, S. Zhang, L. Jin, L. Xie, Y. Wu, and Z. Wang, "Omnidirectional scene text detection with sequential-free box discretization," *arXiv preprint arXiv:1906.02371*, 2019.

[119] W. Qian, X. Yang, S. Peng, J. Yan, and Y. Guo, "Learning modulated loss for rotated object detection," in *in Proc. AAAI Conf. Artific. Intell.*, vol. 35, no. 3, 2021, pp. 2458–2466.

[120] Y. Xu, M. Fu, Q. Wang, Y. Wang, K. Chen, G.-S. Xia, and X. Bai, "Gliding vertex on the horizontal bounding box for multi-oriented object detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 4, pp. 1452–1459, 2021.

[121] W. Qian, X. Yang, S. Peng, X. Zhang, and J. Yan, "Rsdet++: Point-based modulated loss for more accurate rotated object detection," *IEEE*

*Trans. Circuits Syst. Video Technol.*, vol. 32, no. 11, pp. 7869–7879, 2022.

[122] J. Luo, Y. Hu, and J. Li, "Surround-net: A multi-branch arbitrary-oriented detector for remote sensing," *Remote Sens.*, vol. 14, no. 7, p. 1751, 2022.

[123] Q. Song, F. Yang, L. Yang, C. Liu, M. Hu, and L. Xia, "Learning point-guided localization for detection in remote sensing images," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 1084–1094, 2021.

[124] X. Xie, G. Cheng, J. Wang, X. Yao, and J. Han, "Oriented r-cnn for object detection," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2021, pp. 3520–3529.

[125] Y. Yao, G. Cheng, G. Wang, S. Li, P. Zhou, X. Xie, and J. Han, "On improving bounding box representations for oriented object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–11, 2023. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3231340

[126] Q. Ming, L. Miao, Z. Zhou, X. Yang, and Y. Dong, "Optimization for arbitrary-oriented object detection via representation invariance loss," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2021.3115110

[127] X. Yang and J. Yan, "Arbitrary-oriented object detection with circular smooth label," in *Proc. Euro. Conf. Comput. Vis.* Springer, 2020, pp. 677–694.

[128] X. Yang, L. Hou, Y. Zhou, W. Wang, and J. Yan, "Dense label encoding for boundary discontinuity free rotation detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 15 819–15 829.

[129] J. Wang, F. Li, and H. Bi, "Gaussian focal loss: Learning distribution polarized angle prediction for rotated object detection in aerial images," *IEEE Trans. Geosci. Remote Sens.*, 2022.

[130] X. Yang, J. Yan, Q. Ming, W. Wang, X. Zhang, and Q. Tian, "Rethinking rotated object detection with gaussian wasserstein distance loss," in *Proc. Int. Conf. Machine Learn*, 2021, pp. 11 830–11 841.

[131] X. Yang, X. Yang, J. Yang, Q. Ming, W. Wang, Q. Tian, and J. Yan, "Learning high-precision bounding box for rotated object detection via kullback-leibler divergence," vol. 34, pp. 18 381–18 394, 2021.

[132] X. Yang, Y. Zhou, G. Zhang, J. Yang, W. Wang, J. Yan, X. Zhang, and Q. Tian, "The kfiou loss for rotated object detection," *arXiv preprint arXiv:2201.12558*, 2022.

[133] X. Yang, G. Zhang, X. Yang, Y. Zhou, W. Wang, J. Tang, T. He, and J. Yan, "Detecting rotated objects as gaussian distributions and its 3-d generalization," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2022. [Online]. Available: https://doi.org/10.1109/TPAMI.2022.3197152.

[134] J. Wang, J. Ding, H. Guo, W. Cheng, T. Pan, and W. Yang, "Mask obb: A semantic attention-based mask oriented bounding box representation for multi-category object detection in aerial images," *Remote Sens.*, vol. 11, no. 24, p. 2930, 2019.

[135] X. Zhang, G. Wang, P. Zhu, T. Zhang, C. Li, and L. Jiao, "Grs-det: An anchor-free rotation ship detector based on gaussian-mask in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 4, pp. 3518–3531, 2020.

[136] Y. Yang, X. Tang, Y. Cheung, X. Zhang, F. Liu, J. Ma, and L. Jiao, "Ar²det: An accurate and real-time rotational one-stage ship detector in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3092433

[137] F. Zhang, X. Wang, S. Zhou, Y. Wang, and Y. Hou, "Arbitrary-oriented ship detection through center-head point extraction," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2021.

[138] J. Yi, P. Wu, B. Liu, Q. Huang, H. Qu, and D. Metaxas, "Oriented object detection in aerial images with box boundary-aware vectors," in *Proc. IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, 2021, pp. 2150–2159.

[139] Z. Xiao, L. Qian, W. Shao, X. Tan, and K. Wang, "Axis learning for orientated objects detection in aerial images," *Remote Sens.*, vol. 12, no. 6, p. 908, 2020.

[140] X. He, S. Ma, L. He, L. Ru, and C. Wang, "Learning rotated inscribed ellipse for oriented object detection in remote sensing images," *Remote Sens.*, vol. 13, no. 18, p. 3622, 2021.

[141] K. Fu, Z. Chang, Y. Zhang, and X. Sun, "Point-based estimator for arbitrary-oriented object detection in aerial images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 5, pp. 4370–4387, 2020.

[142] H. Wei, Y. Zhang, Z. Chang, H. Li, H. Wang, and X. Sun, "Oriented objects as pairs of middle lines," *ISPRS J. Photogrammetry Remote Sens.*, vol. 169, pp. 268–279, 2020.

[143] L. Zhou, H. Wei, H. Li, W. Zhao, Y. Zhang, and Y. Zhang, "Arbitrary-oriented object detection in remote sensing images based on polar coordinates," *IEEE Access*, vol. 8, pp. 223 373–223 384, 2020.

[144] X. Zheng, W. Zhang, L. Huan, J. Gong, and H. Zhang, "Apronet: Detecting objects with precise orientation from aerial images," *ISPRS J. Photogrammetry Remote Sens.*, vol. 181, pp. 99–112, 2021.

[145] X. Yang, G. Zhang, W. Li, X. Wang, Y. Zhou, and J. Yan, "H2rbox: Horizontal box annotation is all you need for oriented object detection," 2022. [Online]. Available: https://doi.org/10.48550/arXiv.2210.06742

[146] G. Cheng, P. Zhou, and J. Han, "Learning rotation-invariant convolutional neural networks for object detection in vhr optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 12, pp. 7405–7415, 2016.

[147] K. Li, G. Cheng, S. Bu, and X. You, "Rotation-insensitive and context-augmented object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 4, pp. 2337–2348, 2017.

[148] G. Cheng, P. Zhou, and J. Han, "Rifd-cnn: Rotation-invariant and fisher discriminative convolutional neural networks for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 2884–2893.

[149] G. Cheng, J. Han, P. Zhou, and D. Xu, "Learning rotation-invariant and fisher discriminative convolutional neural networks for object detection," *IEEE Trans. Image Process.*, vol. 28, no. 1, pp. 265–278, 2019.

[150] X. Wu, D. Hong, J. Tian, J. Chanussot, W. Li, and R. Tao, "Orsim detector: A novel object detection framework in optical remote sensing imagery using spatial-frequency channel features," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 7, pp. 5146–5158, 2019.

[151] X. Wu, D. Hong, J. Chanussot, Y. Xu, R. Tao, and Y. Wang, "Fourier-based rotation-invariant feature boosting: An efficient framework for geospatial object detection," *IEEE Geosci. Remote Sens. Lett.*, vol. 17, no. 2, pp. 302–306, 2019.

[152] G. Wang, X. Wang, B. Fan, and C. Pan, "Feature extraction by rotation-invariant matrix representation for object detection in aerial image," *IEEE Geosci. Remote Sens. Lett.*, vol. 14, no. 6, pp. 851–855, 2017.

[153] X. Wu, D. Hong, P. Ghamisi, W. Li, and R. Tao, "Msri-ccf: Multi-scale and rotation-insensitive convolutional channel features for geospatial object detection," *Remote Sens.*, vol. 10, no. 12, p. 1990, 2018.

[154] M. Zand, A. Etemad, and M. Greenspan, "Oriented bounding boxes for small and freely rotated objects," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–15, 2021.

[155] J. Han, J. Ding, N. Xue, and G.-S. Xia, "Redet: A rotation-equivariant detector for aerial object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 2786–2795.

[156] J. Han, J. Ding, J. Li, and G. Xia, "Align deep features for oriented object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–11, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3062048

[157] X. Yao, H. Shen, X. Feng, G. Cheng, and J. Han, "R2ipoints: Pursuing rotation-insensitive point representation for aerial object detection," *IEEE Trans. Geosci. Remote Sens.*, 2022.

[158] X. Ye, F. Xiong, J. Lu, J. Zhou, and Y. Qian, "$\mathcal{R}^3$-net: Feature fusion and filtration network for object detection in optical remote sensing images," *Remote Sens.*, vol. 12, no. 24, p. 4027, 2020.

[159] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 7132–7141.

[160] X. Li, W. Wang, X. Hu, and J. Yang, "Selective kernel networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 510–519.

[161] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 3146–3154.

[162] Z. Huang, W. Li, X. Xia, X. Wu, Z. Cai, and R. Tao, "A novel nonlocal-aware pyramid and multiscale multitask refinement detector for object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–20, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3059450

[163] Y. Sun, X. Sun, Z. Wang, and K. Fu, "Oriented ship detection based on strong scattering points network in large-scale SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–18, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3130117

[164] W. Ma, N. Li, H. Zhu, L. Jiao, X. Tang, Y. Guo, and B. Hou, "Feature split-merge-enhancement network for remote sensing object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–17, 2022.

[165] Z. Cui, X. Wang, N. Liu, Z. Cao, and J. Yang, "Ship detection in large-scale SAR images via spatial shuffle-group enhance attention," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 1, pp. 379–391, 2021.

[166] J. Chen, L. Wan, J. Zhu, G. Xu, and M. Deng, "Multi-scale spatial and channel-wise attention for improving object detection in remote

sensing imagery," *IEEE Geosci. Remote Sens. Lett.*, vol. 17, no. 4, pp. 681–685, 2020.

[167] J. Bai, J. Ren, Y. Yang, Z. Xiao, W. Yu, V. Havyarimana, and L. Jiao, "Object detection in large-scale remote-sensing images based on time-frequency analysis and feature optimization," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3119344

[168] J. Hu, X. Zhi, S. Jiang, H. Tang, W. Zhang, and L. Bruzzone, "Supervised multi-scale attention-guided ship detection in optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3206306

[169] Y. Guo, X. Tong, X. Xu, S. Liu, Y. Feng, and H. Xie, "An anchor-free network with density map and attention mechanism for multiscale object detection in aerial images," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2022.3207178

[170] D. Yu and S. Ji, "A new spatial-oriented object detection framework for remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3127232

[171] C. Li, B. Luo, H. Hong, X. Su, Y. Wang, J. Liu, C. Wang, J. Zhang, and L. Wei, "Object detection based on global-local saliency constraint in aerial images," *Remote Sens.*, vol. 12, no. 9, p. 1435, 2020.

[172] J. Lei, X. Luo, L. Fang, M. Wang, and Y. Gu, "Region-enhanced convolutional neural network for object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 8, pp. 5693–5702, 2020.

[173] Y. Yuan, C. Li, J. Kim, W. Cai, and D. D. Feng, "Reversion correction and regularized random walk ranking for saliency detection," *IEEE Trans. Image Process.*, vol. 27, no. 3, pp. 1311–1322, 2018.

[174] C. Xu, C. Li, Z. Cui, T. Zhang, and J. Yang, "Hierarchical semantic propagation for object detection in remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 6, pp. 4353–4364, 2020.

[175] T. Zhang, X. Zhang, P. Zhu, P. Chen, X. Tang, C. Li, and L. Jiao, "Foreground refinement network for rotated object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, 2021.

[176] J. Wang, W. Yang, H. Li, H. Zhang, and G. Xia, "Learning center probability map for detecting objects in aerial images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 5, pp. 4307–4323, 2021.

[177] Z. Fang, J. Ren, H. Sun, S. Marshall, J. Han, and H. Zhao, "Safdet: A semi-anchor-free detector for effective detection of oriented objects in aerial images," *Remote Sens.*, vol. 12, no. 19, p. 3225, 2020.

[178] Z. Ren, Y. Tang, Z. He, L. Tian, Y. Yang, and W. Zhang, "Ship detection in high-resolution optical remote sensing images aided by saliency information," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3173610

[179] H. Qu, L. Shen, W. Guo, and J. Wang, "Ships detection in SAR images based on anchor-free model with mask guidance features," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 15, pp. 666–675, 2022.

[180] S. Liu, L. Zhang, H. Lu, and Y. He, "Center-boundary dual attention for oriented object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3069056

[181] J. Zhang, C. Xie, X. Xu, Z. Shi, and B. Pan, "A contextual bidirectional enhancement method for remote sensing image object detection," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 13, pp. 4518–4531, 2020.

[182] Y. Gong, Z. Xiao, X. Tan, H. Sui, C. Xu, H. Duan, and D. Li, "Context-aware convolutional neural network for object detection in VHR remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 1, pp. 34–44, 2020.

[183] W. Ma, Q. Guo, Y. Wu, W. Zhao, X. Zhang, and L. Jiao, "A novel multi-model decision fusion network for object detection in remote sensing images," *Remote Sens.*, vol. 11, no. 7, p. 737, 2019.

[184] S. Tian, L. Kang, X. Xing, Z. Li, L. Zhao, C. Fan, and Y. Zhang, "Siamese graph embedding network for object detection in remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 18, no. 4, pp. 602–606, 2021.

[185] S. Tian, L. Kang, X. Xing, J. Tian, C. Fan, and Y. Zhang, "A relation-augmented embedded graph attention network for remote sensing object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–18, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3073269

[186] Y. Wu, K. Zhang, J. Wang, Y. Wang, Q. Wang, and Q. Li, "Cdd-net: A context-driven detection network for multiclass object detection,"

[187] Y. Han, J. Liao, T. Lu, T. Pu, and Z. Peng, "Kcpnet: Knowledge-driven context perception networks for ship detection in infrared imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–19, 2023. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3233401

[188] C. Chen, W. Gong, Y. Chen, and W. Li, "Object detection in remote sensing images based on a scene-contextual feature pyramid network," *Remote Sens.*, vol. 11, no. 3, p. 339, 2019.

[189] Z. Wu, B. Hou, B. Ren, Z. Ren, S. Wang, and L. Jiao, "A deep detection network based on interaction of instance segmentation and object detection for SAR images," *Remote Sens.*, vol. 13, no. 13, p. 2582, 2021.

[190] Y. Wu, K. Zhang, J. Wang, Y. Wang, Q. Wang, and X. Li, "Gcwnet: A global context-weaving network for object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–12, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3155899

[191] G. Shi, J. Zhang, J. Liu, C. Zhang, C. Zhou, and S. Yang, "Global context-augmented objection detection in VHR optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 12, pp. 10 604–10 617, 2021.

[192] J. Liu, S. Li, C. Zhou, X. Cao, Y. Gao, and B. Wang, "Sraf-net: A scene-relevant anchor-free object detection network in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3124959

[193] C. Tao, L. Mi, Y. Li, J. Qi, Y. Xiao, and J. Zhang, "Scene context-driven vehicle detection in high-resolution aerial images," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 10, pp. 7339–7351, 2019.

[194] K. Zhang, Y. Wu, J. Wang, Y. Wang, and Q. Wang, "Semantic context-aware network for multiscale object detection in remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2021.3067313

[195] M. Wang, Q. Li, Y. Gu, L. Fang, and X. X. Zhu, "Scaf-net: Scene context attention-based fusion network for vehicle detection in aerial imagery," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2021.3107281

[196] G. Zhang, S. Lu, and W. Zhang, "Cad-net: A context-aware detection network for objects in remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 12, pp. 10 015–10 024, 2019.

[197] E. Liu, Y. Zheng, B. Pan, X. Xu, and Z. Shi, "Dcl-net: Augmenting the capability of classification and localization for remote sensing object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 9, pp. 7933–7944, 2021.

[198] Y. Feng, W. Diao, X. Sun, M. Yan, and X. Gao, "Towards automated ship detection and category recognition from high-resolution aerial images," *Remote Sens.*, vol. 11, no. 16, p. 1901, 2019.

[199] P. Wang, X. Sun, W. Diao, and K. Fu, "FMSSD: feature-merged single-shot detection for multiscale objects in large-scale remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 5, pp. 3377–3390, 2020.

[200] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, 2018.

[201] Y. Bai, R. Li, S. Gou, C. Zhang, Y. Chen, and Z. Zheng, "Cross-connected bidirectional pyramid network for infrared small-dim target detection," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2022.3145577

[202] Y. Li, Q. Huang, X. Pei, Y. Chen, L. Jiao, and R. Shang, "Cross-layer attention network for small object detection in remote sensing imagery," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 2148–2161, 2021.

[203] H. Gong, T. Mu, Q. Li, H. Dai, C. Li, Z. He, W. Wang, F. Han, A. Tuniyazi, H. Li, X. Lang, Z. Li, and B. Wang, "Swin-transformer-enabled yolov5 with attention mechanism for small object detection on satellite images," *Remote Sens.*, vol. 14, no. 12, p. 2861, 2022.

[204] J. Qu, C. Su, Z. Zhang, and A. Razi, "Dilated convolution and feature fusion SSD network for small object detection in remote sensing images," *IEEE Access*, vol. 8, pp. 82 832–82 843, 2020.

[205] T. Ma, Z. Yang, J. Wang, S. Sun, X. Ren, and U. Ahmad, "Infrared small target detection network with generate label and feature mapping," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2022.3140432

[206] W. Han, A. Kuerban, Y. Yang, Z. Huang, B. Liu, and J. Gao, "Multi-vision network for accurate and real-time small object detection in optical remote sensing images," *IEEE Geosci. Remote*

*Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2020.3044422

[207] Q. Hou, Z. Wang, F. Tan, Y. Zhao, H. Zheng, and W. Zhang, "Ristdnet: Robust infrared small target detection network," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2021.3050828

[208] X. Lu, Y. Zhang, Y. Yuan, and Y. Feng, "Gated and axis-concentrated localization network for remote sensing object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 1, pp. 179–192, 2020.

[209] L. Courtrai, M. Pham, and S. Lefèvre, "Small object detection in remote sensing images based on super-resolution with auxiliary generative adversarial networks," *Remote Sens.*, vol. 12, no. 19, p. 3152, 2020.

[210] S. M. A. Bashir and Y. Wang, "Small object detection in remote sensing images with residual feature aggregation-based super-resolution and object detector network," *Remote Sens.*, vol. 13, no. 9, p. 1854, 2021.

[211] J. Rabbi, N. Ray, M. Schubert, S. Chowdhury, and D. Chao, "Small-object detection in remote sensing images with end-to-end edge-enhanced GAN and object detector network," *Remote Sens.*, vol. 12, no. 9, p. 1432, 2020.

[212] J. Wu and S. Xu, "From point to region: Accurate and efficient hierarchical small object detection in low-resolution remote sensing images," *Remote Sens.*, vol. 13, no. 13, p. 2620, 2021.

[213] J. Li, Z. Zhang, Y. Tian, Y. Xu, Y. Wen, and S. Wang, "Target-guided feature super-resolution for vehicle detection in remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2021.3112172

[214] J. Chen, K. Chen, H. Chen, Z. Zou, and Z. Shi, "A degraded reconstruction enhancement-based method for tiny ship detection in remote sensing images with a new large-scale dataset," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3180894

[215] 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 Trans. Geosci. Remote Sens.*, vol. 57, no. 8, pp. 5512–5524, 2019.

[216] J. Wu, Z. Pan, B. Lei, and Y. Hu, "Fsanet: Feature-and-spatial-aligned network for tiny object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–17, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3205052

[217] M. Pham, L. Courtrai, C. Friguet, S. Lefèvre, and A. Baussard, "Yolo-fine: One-stage detector of small objects under various backgrounds in remote sensing images," *Remote Sens.*, vol. 12, no. 15, p. 2501, 2020.

[218] J. Yan, H. Wang, M. Yan, W. Diao, X. Sun, and H. Li, "Iou-adaptive deformable R-CNN: make full use of iou for multi-class object detection in remote sensing imagery," *Remote Sens.*, vol. 11, no. 3, p. 286, 2019.

[219] R. Dong, D. Xu, J. Zhao, L. Jiao, and J. An, "Sig-nms-based faster R-CNN combining transfer learning for small target detection in VHR optical remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 11, pp. 8534–8545, 2019.

[220] Z. Shu, X. Hu, and J. Sun, "Center-point-guided proposal generation for detection of small and dense buildings in aerial imagery," *IEEE Geosci. Remote Sens. Lett.*, vol. 15, no. 7, pp. 1100–1104, 2018.

[221] C. Xu, J. Wang, W. Yang, and L. Yu, "Dot distance for tiny object detection in aerial images," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. Workshops.* IEEE, 2021, pp. 1192–1201.

[222] C. Xu, J. Wang, W. Yang, H. Yu, L. Yu, and G. Xia, "RFLA: gaussian receptive field based label assignment for tiny object detection," pp. 526–543, 2022.

[223] F. Zhang, B. Du, L. Zhang, and M. Xu, "Weakly supervised learning based on coupled convolutional neural networks for aircraft detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 9, pp. 5553–5563, 2016.

[224] Y. Li, B. He, F. Melgani, and T. Long, "Point-based weakly supervised learning for object detection in high spatial resolution remote sensing images," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 5361–5371, 2021.

[225] D. Zhang, J. Han, G. Cheng, Z. Liu, S. Bu, and L. Guo, "Weakly supervised learning for target detection in remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 4, pp. 701–705, 2015.

[226] J. Han, D. Zhang, G. Cheng, L. Guo, and J. Ren, "Object detection in optical remote sensing images based on weakly supervised learning and high-level feature learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 6, pp. 3325–3337, 2015.

[227] Y. Li, Y. Zhang, X. Huang, and A. L. Yuille, "Deep networks under scene-level supervision for multi-class geospatial object detection from remote sensing images," *ISPRS J. Photogrammetry Remote Sens.*, vol. 146, pp. 182–196, 2018.

[228] H. Bilen and A. Vedaldi, "Weakly supervised deep detection networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 2846–2854.

[229] X. Yao, X. Feng, J. Han, G. Cheng, and L. Guo, "Automatic weakly supervised object detection from high spatial resolution remote sensing images via dynamic curriculum learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 1, pp. 675–685, 2021.

[230] H. Wang, H. Li, W. Qian, W. Diao, L. Zhao, J. Zhang, and D. Zhang, "Dynamic pseudo-label generation for weakly supervised object detection in remote sensing images," *Remote Sens.*, vol. 13, no. 8, p. 1461, 2021.

[231] X. Feng, J. Han, X. Yao, and G. Cheng, "Progressive contextual instance refinement for weakly supervised object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 11, pp. 8002–8012, 2020.

[232] P. Shamsolmoali, J. Chanussot, M. Zareapoor, H. Zhou, and J. Yang, "Multipatch feature pyramid network for weakly supervised object detection in optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3106442

[233] B. Wang, Y. Zhao, and X. Li, "Multiple instance graph learning for weakly supervised remote sensing object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–12, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3123231

[234] X. Feng, J. Han, X. Yao, and G. Cheng, "Tcanet: Triple context-aware network for weakly supervised object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 8, pp. 6946–6955, 2021.

[235] X. Feng, X. Yao, G. Cheng, J. Han, and J. Han, "Saenet: Self-supervised adversarial and equivariant network for weakly supervised object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–11, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3105575

[236] X. Qian, Y. Huo, G. Cheng, X. Yao, K. Li, H. Ren, and W. Wang, "Incorporating the completeness and difficulty of proposals into weakly supervised object detection in remote sensing images," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 15, pp. 1902–1911, 2022.

[237] W. Qian, Z. Yan, Z. Zhu, and W. Yin, "Weakly supervised part-based method for combined object detection in remote sensing imagery," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 15, pp. 5024–5036, 2022.

[238] S. Chen, D. Shao, X. Shu, C. Zhang, and J. Wang, "Fcc-net: A full-coverage collaborative network for weakly supervised remote sensing object detection," *Electronics*, vol. 9, no. 9, p. 1356, 2020.

[239] G. Cheng, X. Xie, W. Chen, X. Feng, X. Yao, and J. Han, "Self-guided proposal generation for weakly supervised object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–11, 2022.

[240] X. Feng, X. Yao, G. Cheng, and J. Han, "Weakly supervised rotation-invariant aerial object detection network," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR).* IEEE, 2022, pp. 14 126–14 135.

[241] G. Wang, X. Zhang, Z. Peng, X. Jia, X. Tang, and L. Jiao, "Mol: Towards accurate weakly supervised remote sensing object detection via multi-view noisy learning," *ISPRS J. Photogrammetry Remote Sens.*, vol. 196, pp. 457–470, 2023.

[242] T. Deselaers, B. Alexe, and V. Ferrari, "Weakly supervised localization and learning with generic knowledge," *Int. J. Comput. Vis.*, vol. 100, no. 3, pp. 275–293, 2012.

[243] B. Hou, Z. Wu, B. Ren, Z. Li, X. Guo, S. Wang, and L. Jiao, "A neural network based on consistency learning and adversarial learning for semisupervised synthetic aperture radar ship detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022.

[244] Z. Song, J. Yang, D. Zhang, S. Wang, and Z. Li, "Semi-supervised dim and small infrared ship detection network based on haar wavelet," *IEEE Access*, vol. 9, pp. 29 686–29 695, 2021.

[245] Y. Zhong, Z. Zheng, A. Ma, X. Lu, and L. Zhang, "COLOR: cycling, offline learning, and online representation framework for airport and airplane detection using GF-2 satellite images," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 12, pp. 8438–8449, 2020.

[246] Y. Wu, W. Zhao, R. Zhang, and F. Jiang, "Amr-net: Arbitrary-oriented ship detection using attention module, multi-scale feature fusion and rotation pseudo-label," *IEEE Access*, vol. 9, pp. 68 208–68 222, 2021.

[247] S. Chen, R. Zhan, W. Wang, and J. Zhang, "Domain adaptation for semi-supervised ship detection in SAR images," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2022.3171789

[248] Z. Zhang, Z. Feng, and S. Yang, "Semi-supervised object detection framework with object first mixup for remote sensing images," in *International Geoscience and Remote Sensing Symposium*. IEEE, 2021, pp. 2596–2599.

[249] B. Xue and N. Tong, "DIOD: fast and efficient weakly semi-supervised deep complex ISAR object detection," *IEEE Trans. Cybern.*, vol. 49, no. 11, pp. 3991–4003, 2019.

[250] L. Liao, L. Du, and Y. Guo, "Semi-supervised SAR target detection based on an improved faster R-CNN," *Remote Sens.*, vol. 14, no. 1, p. 143, 2022. [Online]. Available: https://doi.org/10.3390/rs14010143

[251] Y. Du, L. Du, Y. Guo, and Y. Shi, "Semisupervised sar ship detection network via scene characteristic learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–17, 2023. [Online]. Available: https://doi.org/10.1109/TGRS.2023.3235859

[252] D. Wei, Y. Du, L. Du, and L. Li, "Target detection network for SAR images based on semi-supervised learning and attention mechanism," *Remote Sens.*, vol. 13, no. 14, p. 2686, 2021.

[253] L. Chen, Y. Fu, S. You, and H. Liu, "Efficient hybrid supervision for instance segmentation in aerial images," *Remote Sens.*, vol. 13, no. 2, p. 252, 2021.

[254] G. Cheng, B. Yan, P. Shi, K. Li, X. Yao, L. Guo, and J. Han, "Prototype-cnn for few-shot object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–10, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3078507

[255] X. Li, J. Deng, and Y. Fang, "Few-shot object detection on remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3051383

[256] L. Li, X. Yao, G. Cheng, M. Xu, J. Han, and J. Han, "Solo-to-collaborative dual-attention network for one-shot object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–11, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3091003

[257] H. Zhang, X. Zhang, G. Meng, C. Guo, and Z. Jiang, "Few-shot multi-class ship detection in remote sensing images using attention feature map and multi-relation detector," *Remote Sens.*, vol. 14, no. 12, p. 2790, 2022.

[258] B. Wang, Z. Wang, X. Sun, H. Wang, and K. Fu, "Dmml-net: Deep metametric learning for few-shot geographic object segmentation in remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–18, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3116672

[259] J. Li, Y. Tian, Y. Xu, X. Hu, Z. Zhang, H. Wang, and Y. Xiao, "Mm-rcnn: Toward few-shot object detection in remote sensing images with meta memory," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–14, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3228612

[260] Z. Zhao, P. Tang, L. Zhao, and Z. Zhang, "Few-shot object detection of remote sensing images via two-stage fine-tuning," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2021.3116858

[261] Y. Zhou, H. Hu, J. Zhao, H. Zhu, R. Yao, and W. Du, "Few-shot object detection via context-aware aggregation for remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2022.3171257

[262] Y. Wang, C. Xu, C. Liu, and Z. Li, "Context information refinement for few-shot object detection in remote sensing images," *Remote Sens.*, vol. 14, no. 14, p. 3255, 2022.

[263] Z. Zhou, S. Li, W. Guo, and Y. Gu, "Few-shot aircraft detection in satellite videos based on feature scale selection pyramid and proposal contrastive learning," *Remote Sens.*, vol. 14, no. 18, p. 4581, 2022.

[264] S. Chen, J. Zhang, R. Zhan, R. Zhu, and W. Wang, "Few shot object detection for SAR images via feature enhancement and dynamic relationship modeling," *Remote Sens.*, vol. 14, no. 15, p. 3669, 2022.

[265] S. Liu, Y. You, H. Su, G. Meng, W. Yang, and F. Liu, "Few-shot object detection in remote sensing image interpretation: Opportunities and challenges," *Remote Sens.*, vol. 14, no. 18, p. 4435, 2022.

[266] X. Huang, B. He, M. Tong, D. Wang, and C. He, "Few-shot object detection on remote sensing images via shared attention module and balanced fine-tuning strategy," *Remote Sens.*, vol. 13, no. 19, p. 3816, 2021.

[267] Z. Xiao, J. Qi, W. Xue, and P. Zhong, "Few-shot object detection with self-adaptive attention network for remote sensing images," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 4854–4865, 2021.

[268] S. Wolf, J. Meier, L. Sommer, and J. Beyerer, "Double head predictor based few-shot object detection for aerial imagery," in *Proc. IEEE Int. Conf. Comput. Vis. Workshops*. IEEE, 2021, pp. 721–731.

[269] T. Zhang, X. Zhang, P. Zhu, X. Jia, X. Tang, and L. Jiao, "Generalized few-shot object detection in remote sensing images," *ISPRS J. Photogrammetry Remote Sens.*, vol. 195, pp. 353–364, 2023.

[270] G. Heitz and D. Koller, "Learning spatial context: Using stuff to find things," in *Proc. Euro. Conf. Comput. Vis.,*, vol. 5302. Springer, 2008, pp. 30–43.

[271] C. Benedek, X. Descombes, and J. Zerubia, "Building development monitoring in multitemporal remotely sensed image pairs with stochastic birth-death dynamics," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 1, pp. 33–50, 2012.

[272] S. Razakarivony and F. Jurie, "Vehicle detection in aerial imagery : A small target detection benchmark," *J. Vis. Commun. Image Represent.*, vol. 34, pp. 187–203, 2016.

[273] K. Liu and G. Máttyus, "Fast multiclass vehicle detection on aerial images," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 9, pp. 1938–1942, 2015.

[274] H. Zhu, X. Chen, W. Dai, K. Fu, Q. Ye, and J. Jiao, "Orientation robust object detection in aerial images using deep convolutional neural network," in *Proc. IEEE Int. Conf. Image Process.*, 2015, pp. 3735–3739.

[275] T. N. Mundhenk, G. Konjevod, W. A. Sakla, and K. Boakye, "A large contextual dataset for classification, detection and counting of cars with deep learning," in *Proc. Euro. Conf. Comput. Vis.,*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., vol. 9907, 2016, pp. 785–800.

[276] Z. Liu, H. Wang, L. Weng, and Y. Yang, "Ship rotated bounding box space for ship extraction from high-resolution optical satellite images with complex backgrounds," *IEEE Geosci. Remote Sens. Lett.*, vol. 13, no. 8, pp. 1074–1078, 2016.

[277] J. Li, C. Qu, and J. Shao, "Ship detection in sar images based on an improved faster r-cnn," in *SAR in Big Data Era: Models, Methods and Applications (BIGSARDATA)*. IEEE, 2017, pp. 1–6.

[278] Z. Zou and Z. Shi, "Random access memories: A new paradigm for target detection in high resolution aerial remote sensing images," *IEEE Trans. Image Process.*, vol. 27, no. 3, pp. 1100–1111, 2018.

[279] X. Sun, Z. Wang, Y. Sun, W. Diao, Y. Zhang, and K. Fu, "Air-sarship-1.0: High-resolution sar ship detection dataset," *Journal of Radars*, vol. 8, no. 6, pp. 852–863, 2019.

[280] W. Yu, G. Cheng, M. Wang, Y. Yao, X. Xie, X. Yao, and J. Han, "Mar20: A benchmark for military aircraft recognition in remote sensing images," *National Remote Sensing Bulletin*, pp. 1–11, 2022.

[281] K. Chen, M. Wu, J. Liu, and C. Zhang, "FGSD: A dataset for fine-grained ship detection in high resolution satellite images," 2020. [Online]. Available: https://arxiv.org/abs/2003.06832

[282] Y. Han, X. Yang, T. Pu, and Z. Peng, "Fine-grained recognition for oriented ship against complex scenes in optical remote sensing images," *IEEE Trans. Geosci. Remote. Sens.*, vol. 60, pp. 1–18, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3123666

[283] J. Wang, W. Yang, H. Guo, R. Zhang, and G. Xia, "Tiny object detection in aerial images," in *Proc. Int. Conf. Pattern Recognit*. IEEE, 2020, pp. 3791–3798.

[284] G. Cheng, X. Yuan, X. Yao, K. Yan, Q. Zeng, and J. Han, "Towards large-scale small object detection: Survey and benchmarks," 2022. [Online]. Available: https://doi.org/10.48550/arXiv.2207.14096

[285] T. Zhang, X. Zhang, J. Shi, and S. Wei, "Hyperli-net: A hyper-light deep learning network for high-accurate and high-speed ship detection from synthetic aperture radar imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 167, pp. 123–153, 2020.

[286] T. Zhang, X. Zhang, J. Li, X. Xu, B. Wang, X. Zhan, Y. Xu, X. Ke, T. Zeng, H. Su, I. Ahmad, D. Pan, C. Liu, Y. Zhou, J. Shi, and S. Wei, "SAR ship detection dataset (SSDD): official release and comprehensive data analysis," *Remote Sens.*, vol. 13, no. 18, p. 3690, 2021.

[287] M. Everingham, L. V. Gool, C. K. I. Williams, J. M. Winn, and A. Zisserman, "The pascal visual object classes (VOC) challenge," *Int. J. Comput. Vis.*, vol. 88, no. 2, pp. 303–338, 2010.

[288] A. G. Menezes, G. de Moura, C. Alves, and A. C. P. L. F. de Carvalho, "Continual object detection: A review of definitions, strategies, and challenges," *Neural Networks*, vol. 161, pp. 476–493, 2023.

[289] C. Persello, J. D. Wegner, R. Hänsch, D. Tuia, P. Ghamisi, M. Koeva, and G. Camps-Valls, "Deep learning and earth observation to support the sustainable development goals: Current approaches, open challenges, and future opportunities," *IEEE Geosci. Remote Sens. Mag.*, vol. 10, no. 2, pp. 172–200, 2022.

[290] T. Hoeser, F. Bachofer, and C. Kuenzer, "Object detection and image segmentation with deep learning on earth observation data: A review—part ii: Applications," *Remote Sen.*, vol. 12, no. 18, p. 3053, 2020.

[291] L. Ma, Y. Liu, X. Zhang, Y. Ye, G. Yin, and B. A. Johnson, "Deep learning in remote sensing applications: A meta-analysis and review," *ISPRS J. Photogrammetry Remote Sens.*, vol. 152, pp. 166–177, 2019.

[292] P. Barmpoutis, P. Papaioannou, K. Dimitropoulos, and N. Grammalidis, "A review on early forest fire detection systems using optical remote sensing," *Sensors*, vol. 20, no. 22, p. 6442, 2020.

[293] Z. Guan, X. Miao, Y. Mu, Q. Sun, Q. Ye, and D. Gao, "Forest fire segmentation from aerial imagery data using an improved instance segmentation model," *Remote. Sens.*, vol. 14, no. 13, p. 3159, 2022.

[294] Z. Zheng, Y. Zhong, J. Wang, A. Ma, and L. Zhang, "Building damage assessment for rapid disaster response with a deep object-based semantic change detection framework: From natural disasters to man-made disasters," *Remote Sens. Environ.*, vol. 265, p. 112636, 2021.

[295] H. Ma, Y. Liu, Y. Ren, and J. Yu, "Detection of collapsed buildings in post-earthquake remote sensing images based on the improved yolov3," *Remote. Sens.*, vol. 12, no. 1, p. 44, 2020.

[296] Y. Pi, N. D. Nath, and A. H. Behzadan, "Convolutional neural networks for object detection in aerial imagery for disaster response and recovery," *Adv. Eng. Informatics*, vol. 43, p. 101009, 2020.

[297] M. Weiss, F. Jacob, and G. Duveiller, "Remote sensing for agricultural applications: A meta-review," *Remote Sens. Environ.*, vol. 236, p. 111402, 2020.

[298] Y. Pang, Y. Shi, S. Gao, F. Jiang, A. N. V. Sivakumar, L. Thompson, J. D. Luck, and C. Liu, "Improved crop row detection with deep neural network for early-season maize stand count in UAV imagery," *Comput. Electron. Agric.*, vol. 178, p. 105766, 2020.

[299] C. Mota-Delfin, G. de Jesús López-Canteñs, I. L. L. Cruz, E. Romantchik-Kriuchkova, and J. C. Olguín-Rojas, "Detection and counting of corn plants in the presence of weeds with convolutional neural networks," *Remote. Sens.*, vol. 14, no. 19, p. 4892, 2022.

[300] L. P. Osco, M. d. S. de Arruda, D. N. Gonçalves, A. Dias, J. Batistoti, M. de Souza, F. D. G. Gomes, A. P. M. Ramos, L. A. de Castro Jorge, V. Liesenberg *et al.*, "A cnn approach to simultaneously count plants and detect plantation-rows from uav imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 174, pp. 1–17, 2021.

[301] M. M. Anuar, A. A. Halin, T. Perumal, and B. Kalantar, "Aerial imagery paddy seedlings inspection using deep learning," *Remote. Sens.*, vol. 14, no. 2, p. 274, 2022.

[302] Y. Chen, W. S. Lee, H. Gan, N. Peres, C. W. Fraisse, Y. Zhang, and Y. He, "Strawberry yield prediction based on a deep neural network using high-resolution aerial orthoimages," *Remote. Sens.*, vol. 11, no. 13, p. 1584, 2019.

[303] W. Zhao, C. Persello, and A. Stein, "Building outline delineation: From aerial images to polygons with an improved end-to-end learning framework," *ISPRS J. Photogrammetry Remote Sens.*, vol. 175, pp. 119–131, 2021.

[304] Z. Li, J. D. Wegner, and A. Lucchi, "Topological map extraction from overhead images," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 1715–1724.

[305] L. Mou and X. X. Zhu, "Vehicle instance segmentation from aerial image and video using a multitask learning residual fully convolutional network," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 11, pp. 6699–6711, 2018.

[306] J. Zhang, X. Zhang, Z. Huang, X. Cheng, J. Feng, and L. Jiao, "Bidirectional multiple object tracking based on trajectory criteria in satellite videos," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–14, 2023.

[307] H. Kim and Y. Ham, "Participatory sensing-based geospatial localization of distant objects for disaster preparedness in urban built environments," *Automation in Construction*, vol. 107, p. 102960, 2019.

[308] M. A. E. Bhuiyan, C. Witharana, and A. K. Liljedahl, "Use of very high spatial resolution commercial satellite imagery and deep learning to automatically map ice-wedge polygons across tundra vegetation types," *J. Imaging*, vol. 6, no. 12, p. 137, 2020.

[309] W. Zhang, A. K. Liljedahl, M. Kanevskiy, H. E. Epstein, B. M. Jones, M. T. Jorgenson, and K. Kent, "Transferability of the deep learning mask R-CNN model for automated mapping of ice-wedge polygons in high-resolution satellite and UAV images," *Remote. Sens.*, vol. 12, no. 7, p. 1085, 2020.

[310] C. Witharana, M. A. E. Bhuiyan, A. K. Liljedahl, M. Kanevskiy, M. T. Jorgenson, B. M. Jones, R. Daanen, H. E. Epstein, C. G. Griffin, K. Kent, and M. K. W. Jones, "An object-based approach for mapping tundra ice-wedge polygon troughs from very high spatial resolution optical satellite imagery," *Remote. Sens.*, vol. 13, no. 4, p. 558, 2021.

[311] J. Yu, Z. Wang, A. Majumdar, and R. Rajagopal, "Deepsolar: A machine learning framework to efficiently construct a solar deployment database in the united states," *Joule*, vol. 2, no. 12, pp. 2605–2617, 2018.

[312] J. M. Malof, K. Bradbury, L. M. Collins, and R. G. Newell, "Automatic detection of solar photovoltaic arrays in high resolution aerial imagery," *Applied energy*, vol. 183, pp. 229–240, 2016.

[313] W. Zhang, G. Wang, J. Qi, G. Wang, and T. Zhang, "Research on the extraction of wind turbine all over the china based on domestic satellite remote sensing data," in *International Geoscience and Remote Sensing Symposium*. IEEE, 2021, pp. 4167–4170.

[314] W. Hu, T. Feldman, Y. J. Ou, N. Tarn, B. Ye, Y. Xu, J. M. Malof, and K. Bradbury, "Wind turbine detection with synthetic overhead imagery," in *International Geoscience and Remote Sensing Symposium*. IEEE, 2021, pp. 4908–4911.

[315] T. Jia, Z. Kapelan, R. de Vries, P. Vriend, E. C. Peereboom, I. Okkerman, and R. Taormina, "Deep learning for detecting macroplastic litter in water bodies: A review," *Water Research*, vol. 231, p. 119632, 2023.

[316] C. Martin, Q. Zhang, D. Zhai, X. Zhang, and C. M. Duarte, "Enabling a large-scale assessment of litter along saudi arabian red sea shores by combining drones and machine learning," *Environmental Pollution*, vol. 277, p. 116730, 2021.

[317] K. Themistocleous, C. Papoutsa, S. C. Michaelides, and D. G. Hadjimitsis, "Investigating detection of floating plastic litter from space using sentinel-2 imagery," *Remote. Sens.*, vol. 12, no. 16, p. 2648, 2020.

[318] B. Xue, B. Huang, W. Wei, G. Chen, H. Li, N. Zhao, and H. Zhang, "An efficient deep-sea debris detection method using deep neural networks," *IEEE J. Sel. Top. Appl. Earth Obs. Remote. Sens.*, vol. 14, pp. 12 348–12 360, 2021.

[319] J. Peng, D. Wang, X. Liao, Q. Shao, Z. Sun, H. Yue, and H. Ye, "Wild animal survey using uas imagery and deep learning: modified faster r-cnn for kiang detection in tibetan plateau," *ISPRS J. Photogrammetry Remote Sens.*, vol. 169, pp. 364–376, 2020.

[320] N. Rey, M. Volpi, S. Joost, and D. Tuia, "Detecting animals in african savanna with uavs and the crowds," *Remote Sens. Environ.*, vol. 200, pp. 341–351, 2017.

[321] B. Kellenberger, D. Marcos, and D. Tuia, "Detecting mammals in uav images: Best practices to address a substantially imbalanced dataset with deep learning," *Remote Sens. Environ.*, vol. 216, pp. 139–153, 2018.

[322] A. Delplanque, S. Foucher, P. Lejeune, J. Linchant, and J. Théau, "Multispecies detection and identification of african mammals in aerial imagery using convolutional neural networks," *Remote Sensing in Ecology and Conservation*, vol. 8, no. 2, pp. 166–179, 2022.

[323] D. Wang, Q. Shao, and H. Yue, "Surveying wild animals from satellites, manned aircraft and unmanned aerial systems (uass): A review," *Remote Sen.*, vol. 11, no. 11, p. 1308, 2019.

[324] T. Kattenborn, J. Leitloff, F. Schiefer, and S. Hinz, "Review on convolutional neural networks (cnn) in vegetation remote sensing," *ISPRS J. Photogrammetry Remote Sens.*, vol. 173, pp. 24–49, 2021.

[325] T. Dong, Y. Shen, J. Zhang, Y. Ye, and J. Fan, "Progressive cascaded convolutional neural networks for single tree detection with google earth imagery," *Remote. Sens.*, vol. 11, no. 15, p. 1786, 2019.

[326] A. Safonova, S. Tabik, D. Alcaraz-Segura, A. Rubtsov, Y. Maglinets, and F. Herrera, "Detection of fir trees (*Abies sibirica*) damaged by the bark beetle in unmanned aerial vehicle images with deep learning," *Remote. Sens.*, vol. 11, no. 6, p. 643, 2019.

[327] Z. Hao, L. Lin, C. J. Post, E. A. Mikhailova, M. Li, Y. Chen, K. Yu, and J. Liu, "Automated tree-crown and height detection in a young forest plantation using mask region-based convolutional neural network (mask r-cnn)," *ISPRS J. Photogrammetry Remote Sens.*, vol. 178, pp. 112–123, 2021.

[328] A. Sani-Mohammed, W. Yao, and M. Heurich, "Instance segmentation of standing dead trees in dense forest from aerial imagery using deep learning," *ISPRS O. J. Photogrammetry Remote Sens.*, vol. 6, p. 100024, 2022.

[329] A. V. Etten, "You only look twice: Rapid multi-scale object detection in satellite imagery," 2018. [Online]. Available: http://arxiv.org/abs/1805.09512

[330] Q. Lin, J. Zhao, G. Fu, and Z. Yuan, "Crpn-sfnet: A high-performance object detector on large-scale remote sensing images," *IEEE Trans. Neural Networks Learn. Syst.*, vol. 33, no. 1, pp. 416–429, 2022.

[331] D. Hong, L. Gao, N. Yokoya, J. Yao, J. Chanussot, Q. Du, and B. Zhang, "More diverse means better: Multimodal deep learning meets remote-sensing imagery classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 5, pp. 4340–4354, 2021.

[332] D. Hong, N. Yokoya, G.-S. Xia, J. Chanussot, and X. X. Zhu, "X-modalnet: A semi-supervised deep cross-modal network for classifica-

tion of remote sensing data," *ISPRS J. Photogrammetry Remote Sens.*, vol. 167, pp. 12–23, 2020.

[333] M. Segal-Rozenhaimer, A. Li, K. Das, and V. Chirayath, "Cloud detection algorithm for multi-modal satellite imagery using convolutional neural-networks (cnn)," *Remote Sens. Environ.*, vol. 237, p. 111446, 2020.

[334] Y. Shendryk, Y. Rist, C. Ticehurst, and P. Thorburn, "Deep learning for multi-modal classification of cloud, shadow and land cover scenes in planetscope and sentinel-2 imagery," *ISPRS J. Photogrammetry Remote Sens.*, vol. 157, pp. 124–136, 2019.

[335] Y. Shi, L. Du, and Y. Guo, "Unsupervised domain adaptation for SAR target detection," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 6372–6385, 2021.

[336] Y. Zhu, X. Sun, W. Diao, H. Li, and K. Fu, "Rfa-net: Reconstructed feature alignment network for domain adaptation object detection in remote sensing imagery," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 15, pp. 5689–5703, 2022.

[337] T. Xu, X. Sun, W. Diao, L. Zhao, K. Fu, and H. Wang, "Fada: Feature aligned domain adaptive object detection in remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022.

[338] Y. Koga, H. Miyazaki, and R. Shibasaki, "A method for vehicle detection in high-resolution satellite images that uses a region-based object detector and unsupervised domain adaptation," *Remote Sens.*, vol. 12, no. 3, p. 575, 2020.

[339] Y. Shi, L. Du, Y. Guo, and Y. Du, "Unsupervised domain adaptation based on progressive transfer for ship detection: From optical to SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–17, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3185298

[340] P. Zhang, H. Xu, T. Tian, P. Gao, L. Li, T. Zhao, N. Zhang, and J. Tian, "Sefepnet: Scale expansion and feature enhancement pyramid network for SAR aircraft detection with small sample dataset," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 15, pp. 3365–3375, 2022.

[341] S. Dang, Z. Cao, Z. Cui, Y. Pi, and N. Liu, "Open set incremental learning for automatic target recognition," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 7, pp. 4445–4456, 2019.

[342] J. Chen, S. Wang, L. Chen, H. Cai, and Y. Qian, "Incremental detection of remote sensing objects with feature pyramid and knowledge distillation," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–13, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2020.3042554

[343] X. Chen, J. Jiang, Z. Li, H. Qi, Q. Li, J. Liu, L. Zheng, M. Liu, and Y. Deng, "An online continual object detector on VHR remote sensing images with class imbalance," *Eng. Appl. Artif. Intell.*, vol. 117, no. Part, p. 105549, 2023.

[344] J. Li, X. Sun, W. Diao, P. Wang, Y. Feng, X. Lu, and G. Xu, "Class-incremental learning network for small objects enhancing of semantic segmentation in aerial imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–20, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3124303

[345] W. Liu, X. Nie, B. Zhang, and X. Sun, "Incremental learning with open-set recognition for remote sensing image scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3173995

[346] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2009, pp. 248–255.

[347] Y. Long, G. Xia, S. Li, W. Yang, M. Y. Yang, X. X. Zhu, L. Zhang, and D. Li, "On creating benchmark dataset for aerial image interpretation: Reviews, guidances, and million-aid," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 4205–4230, 2021.

[348] G. A. Christie, N. Fendley, J. Wilson, and R. Mukherjee, "Functional map of the world," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 6172–6180.

[349] D. Wang, J. Zhang, B. Du, G.-S. Xia, and D. Tao, "An empirical study of remote sensing pretraining," *IEEE Trans. Geosci. Remote Sens.*, pp. 1–1, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3176603

[350] W. Li, K. Chen, H. Chen, and Z. Shi, "Geographical knowledge-driven representation learning for remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3115569

[351] X. Sun, P. Wang, W. Lu, Z. Zhu, X. Lu, Q. He, J. Li, X. Rong, Z. Yang, H. Chang, Q. He, G. Yang, R. Wang, J. Lu, and K. Fu, "Ringmo: A remote sensing foundation model with masked image modeling," *IEEE Trans. Geosci. Remote Sens.*, pp. 1–1, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3194732

[352] A. Fuller, K. Millard, and J. R. Green, "Satvit: Pretraining transformers for earth observation," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022. [Online]. Available: https://doi.org/10.1109/LGRS.2022.3201489

[353] D. Wang, Q. Zhang, Y. Xu, J. Zhang, B. Du, D. Tao, and L. Zhang, "Advancing plain vision transformer towards remote sensing foundation model," *IEEE Trans. Geosci. Remote Sens.*, pp. 1–1, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3222818

[354] T. Zhang and X. Zhang, "Shipdenet-20: An only 20 convolution layers and <1-mb lightweight SAR ship detector," *IEEE Geosci. Remote Sens. Lett.*, vol. 18, no. 7, pp. 1234–1238, 2021.

[355] T. Zhang, X. Zhang, J. Shi, and S. Wei, "Depthwise separable convolution neural network for high-speed SAR ship detection," *Remote Sens.*, vol. 11, no. 21, p. 2483, 2019.

[356] Z. Wang, L. Du, and Y. Li, "Boosting lightweight cnns through network pruning and knowledge distillation for SAR target recognition," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 8386–8397, 2021.

[357] S. Chen, R. Zhan, W. Wang, and J. Zhang, "Learning slimming SAR ship object detector through network pruning and knowledge distillation," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 1267–1282, 2021.

[358] Y. Zhang, Z. Yan, X. Sun, W. Diao, K. Fu, and L. Wang, "Learning efficient and accurate detectors with dynamic knowledge distillation in remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–19, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2021.3130443

[359] Y. Yang, X. Sun, W. Diao, H. Li, Y. Wu, X. Li, and K. Fu, "Adaptive knowledge distillation for lightweight remote sensing object detectors optimizing," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–15, 2022. [Online]. Available: https://doi.org/10.1109/TGRS.2022.3175213

[360] C. Li, G. Cheng, G. Wang, P. Zhou, and J. Han, "Instance-aware distillation for efficient object detection in remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–11, 2023. [Online]. Available: https://doi.org/10.1109/TGRS.2023.3238801