

Threatening Patch Attacks on Object Detection in Optical Remote Sensing Images

Xuxiang Sun, Gong Cheng, *Member, IEEE*, Lei Pei, Hongda Li, and Junwei Han, *Fellow, IEEE*

Abstract—Advanced Patch Attacks (PAs) on object detection in natural images have pointed out the great safety vulnerability in methods based on deep neural networks. However, little attention has been paid to this topic in Optical Remote Sensing Images (O-RSIs). To this end, we focus on this research, *i.e.*, PAs on object detection in O-RSIs, and propose a more Threatening PA without the scarification of the visual quality, dubbed TPA. Specifically, to address the problem of inconsistency between local and global landscapes in existing patch selection schemes, we propose leveraging the First-Order Difference (FOD) of the objective function before and after masking to select the sub-patches to be attacked. Further, considering the problem of gradient inundation when applying existing coordinate-based loss to PAs directly, we design an IoU-based objective function specific for PAs, dubbed Bounding box Drifting Loss (BDL), which pushes the detected bounding boxes far from the initial ones until there are no intersections between them. Finally, on two widely used benchmarks, *i.e.*, DIOR and DOTA, comprehensive evaluations of our TPA with four typical detectors (Faster R-CNN, FCOS, RetinaNet, and YOLO-v4) witness its remarkable effectiveness. To the best of our knowledge, this is the first attempt to study the PAs on object detection in O-RSIs, and we hope this work can get our readers interested in studying this topic.

Index Terms—Object detection, Adversarial patch attacks, Remote sensing images.

I. INTRODUCTION

DRAW on the powerful representation ability of Deep Neural Networks (DNNs), a great deal of revolutionary achievements have been made in the aspect of image understanding technology [1]–[10]. Similarly, the technology of understanding Optical Remote Sensing Images (O-RSIs) has made great progress [11]–[17]. Nevertheless, the exposed adversarial vulnerability [18], [19] of DNN leaves great security concerns, hindering their widespread applications.

Facing this security hazard, many researchers devote themselves to studying adversarial robustness [18]–[23]. Recently, the security concerns of DNN-based deep learning methods

This work was supported in part by the National Natural Science Foundation of China under Grant 62136007, in part by the Natural Science Basic Research Program of Shaanxi under Grants 2021JC-16 and 2023-JC-ZD-36, in part by the Fundamental Research Funds for the Central Universities, in part by the Guangdong Basic and Applied Basic Research Foundation under Grant 2021B1515020072, and in part by the Innovation Foundation for Doctor Dissertation of Northwestern Polytechnical University under Grant CX2022054. (Gong Cheng is the corresponding author).

Xuxiang Sun, Gong Cheng, Lei Pei, and Hongda Li are with the School of Automation, Northwestern Polytechnical University, Xi'an 710129, China, and also with the Research and Development Institute of Northwestern Polytechnical University in Shenzhen, Shenzhen 518057, China. (e-mail: gcheng@nwpu.edu.cn).

Junwei Han is with the School of Automation, Northwestern Polytechnical University, Xi'an 710129, China.

Digital Object Identifier xx.xxxx/TGRS.xxxx.xxxxxx



Fig. 1. Visualization of samples drawn from DIOR and DOTA. Here we can see that the distribution of objects in remote sensing images are globally sparse and locally dense.

in O-RSIs have also received progressive attention [24]–[27]. Among them, little attention has been paid to the adversarial vulnerability in O-RSI object detection, an essential and typical research field in O-RSI understanding.

To date, there are three widely studied attack schemes for object detection, *i.e.*, the Full-Scale Attacks (FSAs) [28]–[32], the Patch Attacks (PAs) [33]–[35], and the Adversarial Patches (APs) [36]–[38]. In short, FSAs perturb the whole image, and only the pixels in some specific regions are perturbed in PAs. Compared to FSAs and PAs, the adversarial examples generated by APs are more human-perceptible, and all the targets in an image share the same pattern. The only difference among these APs is their physical parameters (location, angle, scale, *etc.*). Given the intrinsic property of O-RSIs, *i.e.*, the objects in O-RSIs are characterized by globally sparse and locally dense (*cf.*, Fig. 1 for better visualization), PAs exhibit more threats than the other attacks for O-RSIs, on account of their region-efficiency and visual-imperceptibility.

In general, the attack scheme of PAs consists of two critical steps, *i.e.*, the patch selection scheme and the objective function. For the former, recent research [35] proposes to leverage the norm of the gradients passed from the objective function to select the most critical sub-patches. However, adding perturbation is to produce a significant function drop, but the gradient is defined within a tiny neighborhood of data points. Thus, it could not simulate the masking manipulation in PAs. Besides, previous study [23] has shown that the direction of gradients does not always align with the optimal direction (*cf.* Fig. 2 for more details). That is, for high-dimensional nonlinear functions such as DNNs, the local landscape and the global landscape around a data point are usually inconsistent. To this light, we propose a patch selection scheme based on First-Order Difference (FOD), which first calculates the

FOD of the objective function by masking the sub-patches and selects top-k sub-patches with the largest FOD. In this way, we could find the most critical sub-patches within a relatively larger neighborhood.

Besides, the objective functions leveraged in existing PAs [33]–[35] only focus on the classification branch without the attack on the Bounding box (Bbox) regression. Fortunately, a commonly adopted loss [30] in the field of FSAs designs the Coordinate-Based Loss (CBL) to make all the Bboxes cover the entire image. However, the perturbed regions in PAs are not large enough to force the detected Bboxes to cover the whole image, especially for the images with the characteristic of global sparse such as O-RSIs (cf. Fig. 1). Then, the gradient of CBL will exist over the entire attack progress with a larger magnitude than that passed from the classification branch [30]. Consequently, the gradients passed from the classification head may be at risk of being inundated by those passed from the regression branch. To this end, we propose a Bbox Drifting Loss (BDL) to merely reduce the Intersection over Union (IoU) between the detected Bboxes and the initial ones so as to avoid the problem of gradient inundation when applying CBL to PAs directly.

Finally, we validate the effectiveness of our method on DIOR [11] and DOTA [12], respectively. Here, a total of seven typical detectors are utilized to evaluate the general effectiveness of our method. Specifically, four kinds of victim detectors including Faster R-CNN [3], RetinaNet [5], FCOS [4], and Yolo-v4 [6] are leveraged for the evaluations, and we equip the first three detectors with two backbones, *i.e.*, ResNet-50 [8] with Feature Pyramid Networks (FPN) [39] and ResNet-101 [8] with FPN [39]. Throughout the comprehensive evaluations, our TPA achieves the most threatening results.

In summary, our contributions are:

- The threats of patch attacks on object detection in O-RSIs are exhibited for the first time in this paper, which provides the preliminary empirical evidence for the safety concern when applying DNN-based methods to practical deployment.
- We propose FOD patch selection scheme to boost the visual-efficiency of patch attacks. It imitates the attack scheme in PAs by masking the sub-patches and selecting the ones with the highest FOD.
- We propose Bounding box Drifting Loss, an IoU-based objective function specialized for patch attacks. In this way, the gradient passed from the regression branch can be stopped if there are no overlaps between the detected Bboxes and the initial ones.

II. RELATED WORK

In this section, we first provide a brief review of the classical researches on adversarial attacks for image recognition. Later, we will pay much attention to the studies on adversarial attacks for object detection.

A. Adversarial Attacks on Image Recognition

Early studies on adversarial vulnerability mainly focus on the task of image recognition. Among them, the white-box

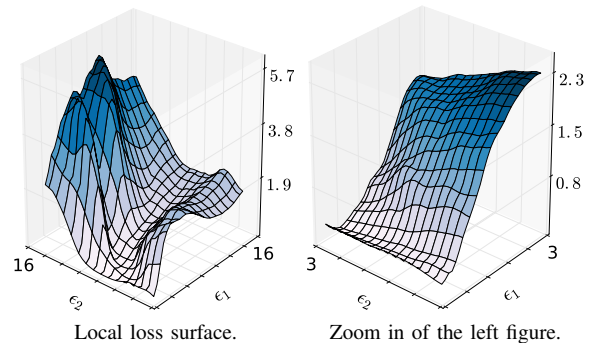


Fig. 2. Illustrations of the local curvature artifacts [23]. It reflects the loss of the adversarially trained Inception-v3 [1] ($v3_{adv}$), with the form of $\mathbf{x}^* = \mathbf{x} + \epsilon_1 \cdot \mathbf{g} + \epsilon_2 \cdot \mathbf{g}^\perp$. Here, \mathbf{g} is the signed gradient of $v3_{adv}$ and \mathbf{g}^\perp is an orthogonal adversarial direction, sampled from Inception-v4 [2].

setting [18]–[22], where the attackers have full access to the victim model, received a wide-spread attention. Later, on the black-box setting, some researchers tried to enhance the transferability via data-augmentation [40]–[42], advanced optimization scheme [42], [43], etc. Besides, some other attack scenarios also received in-depth studies, *e.g.*, query-based black-box attacks [44]–[47] and model stealing [48], [49]. In fact, there is a wide variety of works on this topic. Limited by the space, we could not review them thoroughly in this paper. To this end, we recommend our readers [50] for more comprehensive reviews of the research progress.

B. Adversarial Attacks on Object Detection

Full-Scale Attacks. Xie *et al.* [28] propose a first attempt, dubbed DAG, to attack object detection and segmentation. It sets a wrong category for the target to increase the confidence of negative samples while reducing the confidence of positive samples in an iterative manner. Inspired by Carlini *et al.* [20] and the Expectation over Transformation (EoT) [51], Chen *et al.* [29] propose to add random perturbation over each iteration so as to enhance the robustness of adversarial examples for Faster R-CNN [3]. Later, RAP [30] attacks the Region Proposal Networks (RPN) in two-stage object detectors by destroying both the classification and regression. In addition to the objective function of RAP [30], Zhang *et al.* [31] introduces contextual loss to increase the confidence of the background and inhibits the confidence of the foreground. More recently, Nezami *et al.* [32] proposes to precisely manipulate the pixel of the target object to change its label without affecting the other objects.

Patch Attacks. In this field, Wu *et al.* [33] proposes a diffused patch with the shape of asteroid-like or grid-like and pays more attention to the proposals that escaped from attack. Zhao *et al.* [34] designs heatmap-based and consensus-based algorithms to select patches for the attack. Recently, RPA-tack [35] enhances the threats of PAs in patch selection and optimization schemes. Specifically, it proposes a patch selection based on the gradient feedback and leverages the ensemble learning to improve the attack strength.

Adversarial Patches. Brown *et al.* [36] is the first attempt to attack object detection via a single adversarial patch. It

designs an unrestricted patch with a fixed position to attack the classification branch. Liu *et al.* [37] mislead the victim detector by forcing it to perceive only the stoked rectangular patch at a fixed position. Liu *et al.* [52] proposed a perceptual-sensitive generative adversarial network to synthesize adversarial patches. Moreover, there are a great variety of researches on adversarial patches, including aerial detection [53]–[55] and physical APs for real-world detection [38], [56], [57], etc.

In summary, current works regarding adversarial attacks on object detection pay much attention to FSAs and APs, while the research on PAs has not received widespread attention. However, given the threats posed by the imperceptibility of PAs, it deserves in-depth research. Therefore, we take a closer look at PAs on object detection in this paper.

III. PRELIMINARY

A. Problem Formulation

Considering the general representation, we denote the ground truth of an image \mathbf{x} as:

$$\mathbf{O}(\mathbf{x}) = \{\mathbf{B}_i(\mathbf{x}), \mathbf{y}_i(\mathbf{x})\}, \quad (i = 1, 2, 3 \dots N), \quad (1)$$

where N denotes the number of instance in \mathbf{x} . Here, $\mathbf{B}_i(\mathbf{x}) = \{B_i^x, B_i^y, B_i^w, B_i^h\}$ is the location information with (B_i^x, B_i^y) denote the coordinates of the Bbox center point, and (B_i^w, B_i^h) are the width and height of the Bbox. The class information $\mathbf{y}_i(\mathbf{x}) \in \{1, 2, \dots, C\}$ denotes the label of an instance, where C is the number of categories. In this way, the detected results for a given image \mathbf{x} can be represented as:

$$\tilde{\mathbf{O}}(\mathbf{x}) = \{\tilde{\mathbf{B}}_i(\mathbf{x}), \tilde{\mathbf{P}}_i(\mathbf{x})\}, \quad (2)$$

where $\tilde{\mathbf{P}}_i(\mathbf{x}) = \{\tilde{P}_i^1, \tilde{P}_i^2, \dots, \tilde{P}_i^C\}$ is the class probability vector. Then, the detected class is the index corresponding to the maximum in $\tilde{\mathbf{P}}_i(\mathbf{x})$. That is,

$$\tilde{C}_i(\mathbf{x}) = \operatorname{argmax} \tilde{\mathbf{P}}_i(\mathbf{x}). \quad (3)$$

Based on the above notations, the objective of adversarial attacks on object detection can be formulated as the following:

$$\begin{aligned} \min \|\xi\|_p \\ \text{s.t. } \mathbf{O}(\mathbf{x}_{\text{adv}}) \neq \mathbf{O}(\mathbf{x}) \end{aligned}, \quad (4)$$

where \mathbf{x} are clean images and ξ are the corresponding perturbations. The objective in Eq. (4) is to find the adversarial example with the minimum visual distortion. Since both our method and our competitors are based on the BIM optimizing framework, we choose the ℓ_∞ norm as the visual constraint, *i.e.*, $p = \infty$ in Eq. (4).

Here, the misleading to either the branch of classification or regression can be seen as a successful attack. That is, either $\tilde{C}_i(\mathbf{x}) \neq C_i(\mathbf{x}_{\text{adv}})$ or $\text{IoU}(\tilde{\mathbf{B}}_i(\mathbf{x}), \tilde{\mathbf{B}}_i(\mathbf{x}_{\text{adv}})) < 0.5$ can be seen as a successful attack. Besides, for PAs, the representation of adversarial examples can be formulated as:

$$\mathbf{x}_{\text{adv}} = \mathbf{x} + \mathbf{M}^* \odot \xi^*, \quad (5)$$

where \odot is the element-wise multiplication. \mathbf{M} is the attack map with the same size as \mathbf{x} , which determine the regions to be attacked. Here, ξ^* represents the optimal solution of Eq. (4) and \mathbf{M}^* denotes the optimal results of the attack map, which is fixed over the entire attack progress in general.

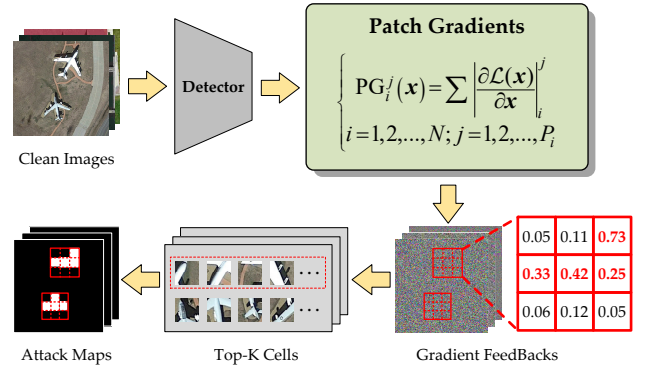


Fig. 3. Illustration of the patch selection scheme in RPAAttack [35].

B. Patch Selection Scheme

In this section, we mainly introduce the patch selection scheme in RPAAttack [35]. For simplicity, we dub it Gradient Feedback (GF). As can be seen in Fig. 3, it first splits the Bbox of each instance evenly to get a series of sub-patches. Then, the sub-patches are ranked w.r.t. the gradient norm, and the top-k sub-patches with the largest gradient norms are selected. Here, \mathcal{L} denotes the objective function leveraged for the optimization of the perturbations and $PG_i^j(\mathbf{x})$ denotes the ℓ_1 norm of the j -th sub-patch in the i -th instance of an image \mathbf{x} .

C. Objective Function

Here, the objective function proposed by [30] to attack the regression branch of a detector is introduced in this section. For simplicity, we dub it the Coordinate-Based Loss (CBL). In general, attacking the regression branch aims at making the detected Bboxes own undesirable shape or position. To this end, RAP assigns large offsets for the detected Bboxes to make them cover the entire image as much as possible. Formally, the CBL is expressed as:

$$\begin{aligned} \mathcal{L}_{\text{CBL}}(\tilde{\mathbf{B}}(\mathbf{x}_{\text{adv}})) = \sum_{j=1}^m z_j ((B_j^x - \tau^x)^2 + (B_j^y - \tau^y)^2 \\ + (B_j^w - \tau^w)^2 + (B_j^h - \tau^h)^2) \end{aligned}, \quad (6)$$

where $\{\tau^x, \tau^y, \tau^w, \tau^h\}$ is the predefined offsets. Besides, z_j is the indicator of j -th proposal and its formulation can be expressed as:

$$\begin{cases} z_j = 1, & \text{if } \begin{cases} \text{IoU}(\tilde{\mathbf{B}}_j(\mathbf{x}), \tilde{\mathbf{B}}_j(\mathbf{x}_{\text{adv}})) > 0.1 \\ \max \tilde{\mathbf{P}}_j(\mathbf{x}_{\text{adv}}) > 0.4 \end{cases} \\ z_j = 0, & \text{otherwise} \end{cases}. \quad (7)$$

In this paper, we follow the setting of RAP [30] to set $\tau^x = \tau^y = \tau^w = \tau^h = 10^5$ in Eq. (7).

IV. METHOD

A. Overview

This paper aims to propose a threatening patch attack that is applicable to multiple target detectors, including Faster R-CNN [3], RetinaNet [5], YOLO-v4 [6], and FCOS [4]. Since the structures of these detector networks are different, the

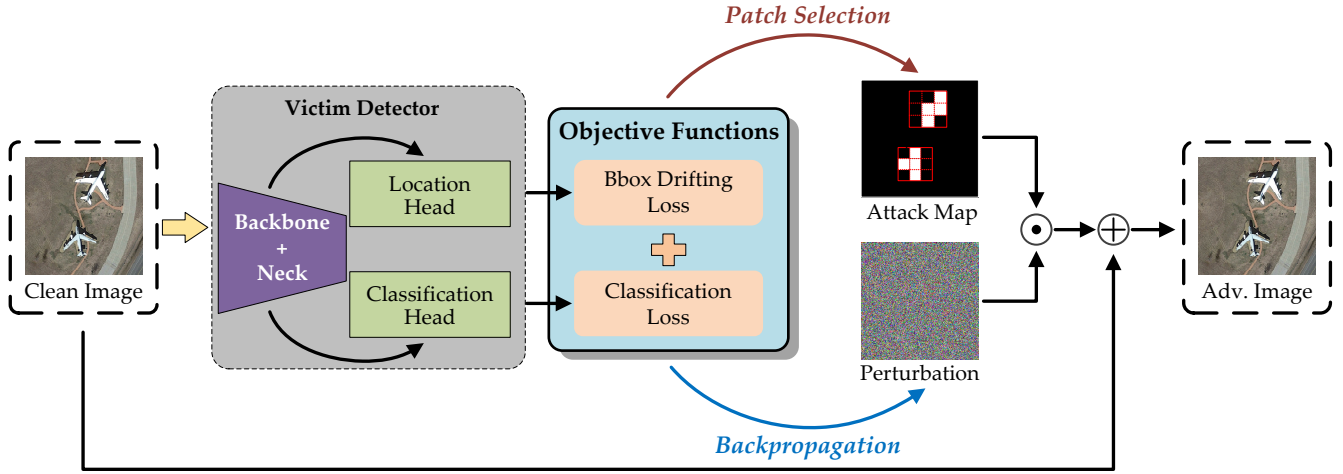


Fig. 4. **The framework of our proposed TPA.** Here, \odot and \oplus denote the element-wise multiplication and addition, respectively. Also, the generated adversarial examples should be clipped to legal interval of image, as shown in Eq. (5). Please refer to Sec. IV-A for more exact illustration.

attack framework is not designed specifically for a certain kind of detector but based on the final prediction results. Here, we depict the framework of our TPA in Fig. 4. Specifically, we first select the regions to acquire the attack map \mathbf{M} via the proposed FOD patch selection scheme, and \mathbf{M} will be fixed over the entire attack progress. Then, similar to BIM [58], we optimize the adversarial example in an iterative manner.

Formally, at each iteration $t + 1$, we have:

$$\begin{aligned} \mathbf{x}_{\text{adv}} &= \hat{\mathbf{x}}_{t+1} \\ &= \mathbf{x} + \mathbf{M} \odot \boldsymbol{\xi}_{t+1}, \end{aligned} \quad (8)$$

where $\boldsymbol{\xi}_{t+1}$ is formulated as:

$$\boldsymbol{\xi}_{t+1} = \text{Clip}_{-\epsilon}^{\epsilon} \{ \text{Clip}_0^1 \{ \hat{\mathbf{x}}_t + \mathbf{M} \odot [\alpha \cdot \text{sign}(\mathbf{g}_{t+1})] \} - \mathbf{x} \}, \quad (9)$$

where ϵ denotes the ℓ_{∞} constraint and α is the attack step size. In our TPA, the attack map \mathbf{M} is a binary mask, which determines where to attack. Besides, the expression of \mathbf{g}_{t+1} in Eq. (9) is:

$$\mathbf{g}_{t+1} = \frac{\nabla_{\hat{\mathbf{x}}_t} \mathcal{L}(\hat{\mathbf{x}}_t)}{\|\nabla_{\hat{\mathbf{x}}_t} \mathcal{L}(\hat{\mathbf{x}}_t)\|_1}. \quad (10)$$

Here, the specific details regarding how to get the attack map \mathbf{M} and the expression of $\mathcal{L}(\mathbf{x})$ will be introduced next.

B. First-Order Difference Patch Selection Scheme

For patch attacks, the selection of sub-patches is a critical factor that affects the attack efficiency. However, as we have mentioned in Sec. I, the advanced patch selection scheme proposed by RPAAttack [35] may suffer from the problem of inconsistency between the local and global landscapes, since the gradient is defined within a small neighborhood around the data point, which is not large enough to explore the global landscapes. To this end, we propose to imitate the “masking” manipulation in patch attacks by covering each sub-patch of the instance and select the sub-patches with the highest feedback.

Specifically, as shown in Fig. 5, we first divide the Bbox of each instance into a grid of $n \times n$. Here, considering the objects in O-RSIs with the same class could have different

sizes, n can vary according to the size of the instance. In this paper, we provide two options for the grid segmentation, *i.e.*, the uniform segmentation scheme and the scale-adaptive segmentation scheme. All the instances in the uniform segmentation share the same setting of n . For simplicity, we denote $U(n)$ as the uniform segmentation scheme with the size of n . For the scale-adaptive segmentation scheme, similar to MS COCO [59], we divide the instances into three kinds of scales in terms of their areas. To be specific, we denote the area of an instance as \mathcal{S} . Then, we set n to n_1 for $\mathcal{S} \leq 32^2$, n_2 for $32^2 < \mathcal{S} \leq 64^2$, and n_3 for $\mathcal{S} > 64^2$. Thus, we can use $SA(n_1, n_2, n_3)$ to represent the scale-adaptive segmentation scheme with certain parameters. The results regarding these two patch segmentation scheme will be reported and analyzed in Sec. V-C.

Once we acquire the masked inputs, we feed them into the victim detector to calculate the FOD. Formally, for a sub-patch, we define the FOD as:

$$\begin{aligned} \text{FOD}_i^j(\mathbf{x}) &= \mathcal{L}(\mathbf{x}) - \mathcal{L}(\mathbf{x}_i^j) \\ &= \mathcal{L}(\mathbf{x}) - \mathcal{L}(\mathbf{M}_i^j \odot \mathbf{x}), \end{aligned} \quad (11)$$

where \mathbf{M}_i^j denotes masking the j -th sub-patch of the i -th instance. Correspondingly, \mathbf{x}_i^j is the input with the j -th sub-patch of the i -th instance is masked, and $\text{FOD}_i^j(\mathbf{x})$ is the first-order difference of the j -th sub-patch belonging to the i -th instance. Associating with the objective function utilized in our TPA (please refer to Sec. IV-C for more detailed introduction), Eq. (11) can be expressed as:

$$(\max \tilde{\mathbf{P}}_i(\mathbf{x}) - \tilde{\mathbf{P}}_i^{\tilde{C}_i(\mathbf{x})}(\mathbf{x}_i^j)) + (1 - \text{IoU}(\tilde{\mathbf{B}}_i(\mathbf{x}), \tilde{\mathbf{B}}_i(\mathbf{x}_i^j))), \quad (12)$$

where $\tilde{\mathbf{P}}_i^{\tilde{C}_i(\mathbf{x})}(\mathbf{x}_i^j)$ denotes the $\tilde{C}_i(\mathbf{x})$ -th entry of $\tilde{\mathbf{P}}_i(\mathbf{x}_i^j)$.

C. Bounding box Drifting Loss

Compared to the image classifier, attacking the object detector is deemed as a more complicated problem, on account of the complex outputs from the detector. Therefore, in addition to attacking the classification branch, destroying

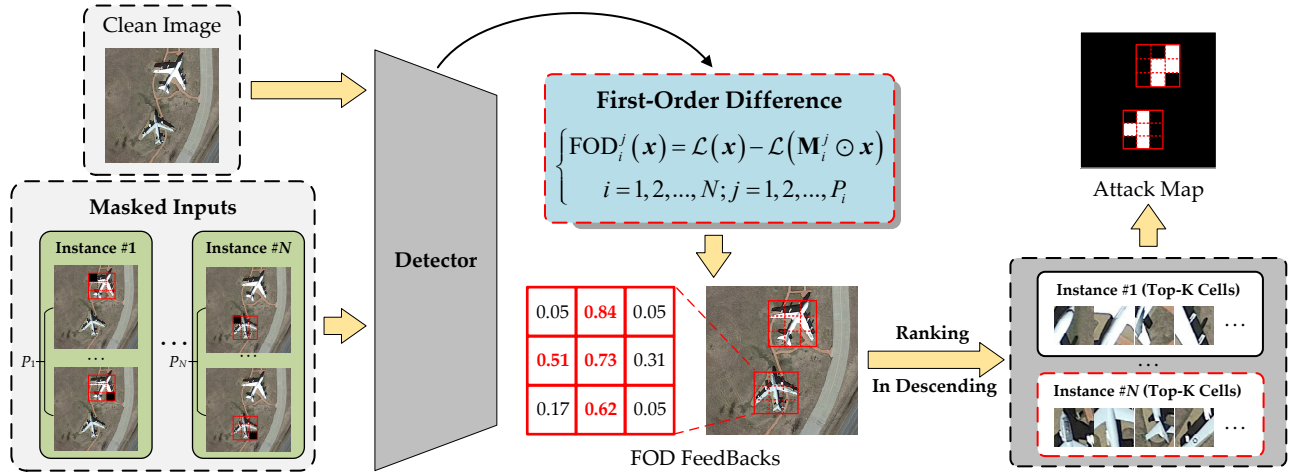


Fig. 5. Illustration of the proposed FOD patch selection scheme.

the regression branch of the detector is another essential step, which could enhance the threats of attacks to a large extent. Existing methods for PAs only focus on the attack against the classification branch. Fortunately, RAP [30] has proposed CBL to break the Bbox regression branch. However, as we analyzed in Sec. I, applying CBL directly to the framework of PAs may be at risk of gradient inundation. In other words, since PAs only attack some specific regions of an image, making all the detected Bboxes covering the entire image seems sets an *impossible* objective. As a result, the gradients of CBL will exist over the entire attack progress. Besides, the norms of the gradient passed from CBL are larger than that passed from the classification loss, on account of the large threshold in CBL. Consequently, the gradient of the classification loss may suffer from being inundated by that of CBL, leading to the stagnation of optimization.

In fact, with the goal of attacking the regression branch, offsetting the detected Bboxes is exact what we want to see. To this end, there exist many solutions for this purpose. One of the most threatening scenario is that there are no intersections between the initial Bbox and detected ones after attacking. Thus, we formulate the attacking on the regression branch as the above situation. That is, drifting the detected Bboxes away from the initial one until there are no overlap between them. When it comes to measuring the overlaps between two Bboxes, a natural idea is to leverage the IoU, a commonly-adopted metric in object detection. Therefore, we formulate our Bbox drifting loss as:

$$\mathcal{L}_{\text{BDL}}(\hat{x}) = \frac{1}{N} \sum_{i=1}^N \max(\text{IoU}(\tilde{B}_i(x), \tilde{B}_i(\hat{x}))), \quad (13)$$

where N denotes the number of instances in the initial results. During the iterations, there may exist a lot of detected Bboxes around the initial one. To this end, Eq. (13) takes the Bbox with the highest IoU between the initial one into calculation. In this way, the detected Bboxes that have no overlaps between the initial one are not taken into consideration. That is, these Bboxes have been attacked successfully.

Another loss function utilized in our TPA is for the attack on the classification branch. Here, we use the loss in RPAttack [35], which is formulated as:

$$\mathcal{L}_{\text{cls}}(\hat{x}) = \frac{1}{k} \sum_{i=1}^k \|\max \tilde{P}_i(\hat{x})\|^2, \quad (14)$$

where k denotes the number of detected results in \hat{x} . Finally, we use Eq. (15) as the total objective function in our TPA.

$$\mathcal{L}(\hat{x}) = \mathcal{L}_{\text{BDL}}(\hat{x}) + \mathcal{L}_{\text{cls}}(\hat{x}). \quad (15)$$

V. EXPERIMENTS

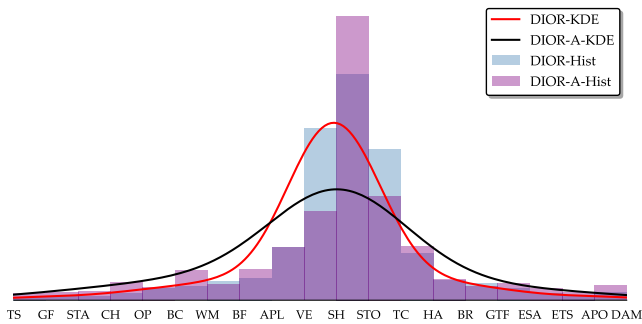
A. Experimental Settings

Datasets. We carry out the evaluations on two widely-adopted benchmarks for object detection in O-RSIs, *i.e.*, DIOR [11] and DOTA [12]. For DIOR dataset, it contains 23463 images in RGB color space, covering 192518 instances of 20 categories. All the images are formatted in a fixed size of 800×800 with a spatial resolution varying across 0.5 to 30 meters. DOTA dataset (we use DOTA-1.0 in this paper) includes 2860 images covering 15 categories, and the size of images in DOTA vary across 800×800 to 4000×4000 . In practical, considering computational burden caused by the large scale of images in DOTA, we split the images into a fixed size, which is set to 800×800 in this paper. Besides, to facilitate the evaluation of the following research regarding adversarial attacks on object detection in O-RSIs, similar to the commonly-adopted protocol [35], [43] in the field of adversarial attacks, we sample 2000 images from the testing subset of DIOR and the validation subset of DOTA, respectively, dubbed DIOR-A and DOTA-A. The class-wise instance distributions of them are exhibited in Fig. 6. Here, we plot both the class-wise instance distribution histograms and their corresponding Kernel Density Estimation (KDE) curves, in which we can see that the sampled datasets share almost the same class-wise distribution with their corresponding parent datasets. Besides, since only 2000 images are utilized for the evaluation, they could reduce the calculation budget to some extent, compared to using the

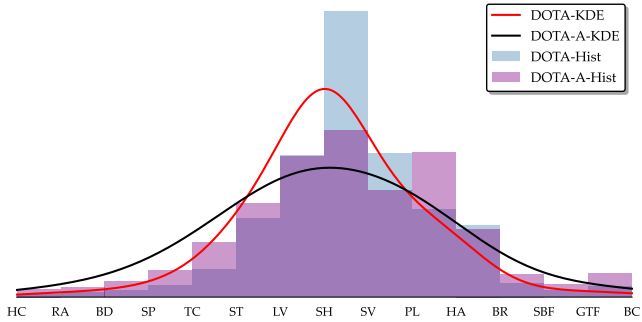
TABLE II
COMPARISON RESULTS ON DIOR-A AND DOTA-A DATASETS. RESULTS ARE SEPARATED BY THE DOUBLE LINE, ABOVE WHICH ARE THE RESULTS ON DIOR-A AND THE REMAINS ARE ON DOTA-A. THE BEST RESULTS ARE SHOWN IN BOLD.

Detector	FR-50		FR-101		FC-50		FC-101		RT-50		RT-101		YOLO-v4	
Method	RPA [35]	Ours	RPA [35]	Ours	RPA [35]	Ours	RPA [35]	Ours	RPA [35]	Ours	RPA [35]	Ours	RPA [35]	Ours
mAP (%)	53.20	47.80	55.80	50.80	-	47.00	-	51.90	-	39.60	-	44.70	50.40	20.30
Recall (%)	69.70	55.80	71.00	58.30	-	61.50	-	66.40	-	58.00	-	60.70	29.60	20.20
ℓ_0 Norm	0.110	0.059	0.120	0.059	-	0.059	-	0.059	-	0.059	-	0.059	0.110	0.059
ℓ_2 Norm	6.840	4.990	6.860	5.160	-	5.640	-	5.790	-	5.170	-	5.310	6.840	4.860
mAP (%)	31.90	26.60	54.00	39.40	-	28.70	-	32.10	-	17.10	-	20.50	34.80	21.70
Recall (%)	49.80	26.60	65.50	48.20	-	47.60	-	50.30	-	37.00	-	40.50	36.40	22.30
ℓ_0 Norm	0.060	0.032	0.060	0.031	-	0.033	-	0.032	-	0.032	-	0.032	0.060	0.032
ℓ_2 Norm	5.412	3.628	5.557	3.756	-	4.304	-	4.411	-	3.944	-	4.005	5.412	3.731

Since RPAttack [35] (RPA) utilizes Faster R-CNN [3] and YOLO-v4 [6] to carry out an ensemble attack, it can not attack RetinaNet and FCOS under the white-box setting. In this case, we only report the results of RPAttack on FC-50, FC-101, and YOLO-v4.



DIOR v.s. DIOR-A w.r.t. Class-Wise Instance Distribution.



DOTA v.s. DOTA-A w.r.t. Class-Wise Instance Distribution.

Fig. 6. Class-Wise instance distributions.

parent datasets for the evaluation. Thus, we will carry out all the following experiments on these sampled datasets.

General settings. We leverage Pytorch framework to implement our method on a single NVIDIA RTX 2080Ti. Here, four kinds of detectors are utilized for the evaluations, where Faster R-CNN [3] (FR), RetinaNet [5] (RT), and FCOS [4] (FC) are trained on MMDetection, and YOLO-v4 [6] are based on DarkNet [7]. For the detectors on MMDetection, we equip them with two backbones including ResNet-50 [8] with FPN [39] and ResNet-101 [8] with FPN [39]. For simplicity, FR-50 denotes Faster R-CNN with ResNet-50+FPN as the backbone, and so on. The testing results on clean images of DIOR-A and DOTA-A are reported in Tab. I. Here, for DIOR-A, we utilize the train-val subsets to train the victim detectors, and training subset is leveraged to train the detectors for DOTA-A. More training details and the sampled datasets

TABLE I
RESULTS ON DIOR-A AND DOTA-A DATASETS.

Detector	FR-50	FR-101	FC-50	FC-101	RT-50	RT-101	YOLO-v4
mAP (%)	88.30	88.60	87.30	87.60	87.30	87.30	89.50
Recall (%)	90.30	90.90	91.30	91.60	92.80	92.80	90.00
mAP (%)	68.70	68.40	65.70	66.80	62.20	64.80	69.70
Recall (%)	77.70	76.10	79.10	80.00	79.50	81.30	76.80

Here, the results above the double lines are those on DIOR-A dataset and the others are those on DOTA-A dataset, the same to Tab. II

are open accessed¹.

Attack settings. Considering the implementation and relativity, we choose RPAttack [35], the state-of-the-art patch attack on object detection in natural images, as our competitor. For the evaluation metrics, we leverage mAP and Recall to measure the strength of different attacks. Besides, we introduce ℓ_2 and ℓ_0 norms to evaluate the visual quality of adversarial examples. For the attack settings, the ℓ_∞ constraint ϵ in Eq. (9) is set to $10/255$ and the number of iterations T is 10. The step size α in Eq. (9) is set to $1/255$. Furthermore, both RPAttack and our TPA attack $\lfloor \frac{n \times n}{2} \rfloor$ patches of an instance, for which we divide to $n \times n$ patches in total. Here, $\lfloor \cdot \rfloor$ represents the floor division. Finally, for the grid segmentation, we use SA(1, 2, 3) as the main scheme. The reason for this choice will be discussed in Sec. V-C.

B. Peer Comparisons

In this subsection, we carry out experiments to validate the advancement of our TPA. The quantitative results are summarized in Tab. II. Surprisingly, the performance of our TPA can surpass the advanced competitor, RPAttack, which leverages the ensemble setting to enhance its threats, where ensemble setting is a powerful attack setting that utilizes more than one victim for the optimization of adversarial examples. Specifically, we summarize the advantages of our TPA as three folds. First, we exceed RPAttack by a large margin in terms of Recall, which indicates that TPA owns the great potential of hiding targets than RPAttack. Second, we can keep the

¹<https://github.com/plpl2019/TPA>

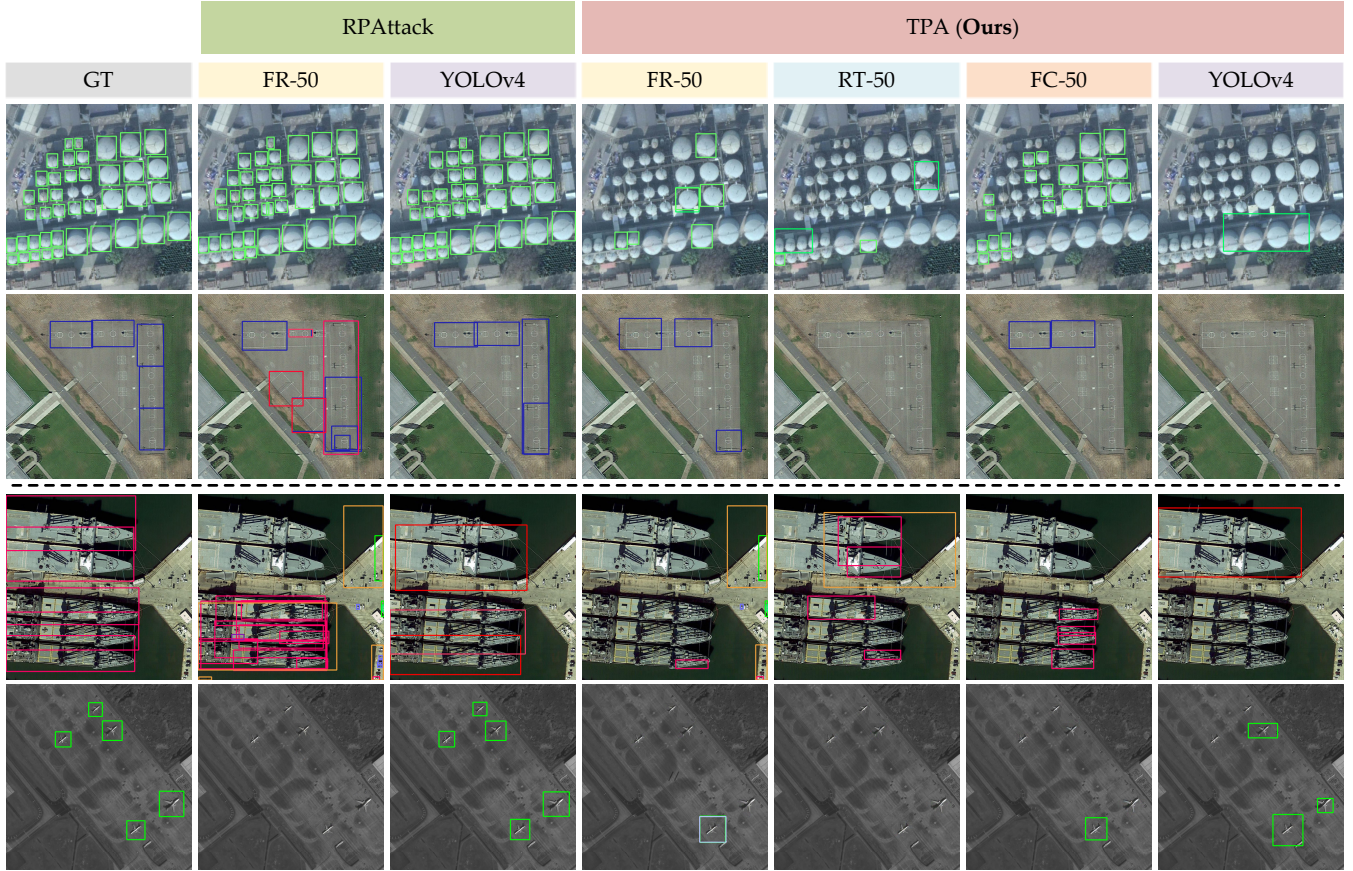


Fig. 7. **Qualitative results w.r.t. the effect of the bounding box drifting loss.** Here, GT for ground truth. Samples extracted from DIOR-A and DOTA-A are separated by the dotted line, where the upper samples belong to DIOR-A and the remains belong to DOTA-A.

highest attack strength while keeping the lowest ℓ_2 norms, which demonstrates the great attack efficiency of our TPA. Finally, different from RPAttack that is restricted by the choice of the victim detector, the proposed TPA is applicable to a variety of detection models.

Later, we also visualize the detected results in Fig. 7. Here, all the targets surrounded by the boxes with the same color belong to the same category, and we use different colors to distinct the categories. As we can see from Fig. 7, compared to RPAttack, TPA can make most targets "invisible". In addition, Fig. 7 releases a significant trend that the invisible ability of RPAttack gets decreased as the density of the targets increases. From this point of view, TPA shows more threats than RPAttack, since the dense instance distribution is a very common phenomenon in O-RSIs.

C. Further Studies

In this subsection, we arrange extended experiments to take a closer look at our TPA. Specifically, we provide three ablation studies regarding the choice of patch selection scheme, the choice of the objective functions for the Bbox regression, and the choice of patch segmentation scheme.

Ablation Study on Patch Selection Scheme. As we have introduced in Sec. I, RPAttack utilizes the gradient feedback to select the most critical regions to be attacked, which

TABLE III
ABLATION RESULTS REGARDING THE PATCH SELECTION SCHEME.

Detector	Method	mAP (%)	Recall (%)	ℓ_0 Norm	ℓ_2 Norm
FR-50	RD	52.80	61.30	0.059	5.050
	GF [35]	50.00	59.40	0.059	4.990
	FOD	47.80	55.80	0.059	4.900
FC-50	RD	52.30	67.50	0.059	5.640
	GF [35]	49.60	65.00	0.059	5.620
	FOD	47.00	61.50	0.059	5.510
RT-50	RD	45.20	64.90	0.059	5.180
	GF [35]	43.10	62.90	0.059	5.180
	FOD	39.60	58.00	0.059	5.170

may suffer from the inconsistency between local and global landscapes, leaving the attack efficiency to be suppressed. To this end, we vary the choice of the patch selection scheme in our TPA to see the effect of our FOD. The results of this part are summarized in Tab. III. Here, GF stands for the gradient feedback patch selection scheme that is leveraged in RPAttack [35] and RD represents selecting the regions in a random manner. Not surprisingly, FOD achieves the best

TABLE IV
ABLATION RESULTS REGARDING THE OBJECTIVE FUNCTION FOR THE BOUNDING BOX REGRESSION.

Detector	Method	mAP (%)	Recall (%)	ℓ_0 Norm	ℓ_2 Norm
FR-50	\mathcal{L}_{cls}	55.20	61.70	0.059	4.800
	$\mathcal{L}_{cls} + \mathcal{L}_{CBL}$	50.90	61.50	0.059	5.340
	$\mathcal{L}_{cls} + \mathcal{L}_{BDL}$	47.80	55.80	0.059	4.990
FC-50	\mathcal{L}_{cls}	55.80	70.06	0.059	5.466
	$\mathcal{L}_{cls} + \mathcal{L}_{CBL}$	52.60	65.30	0.059	5.630
	$\mathcal{L}_{cls} + \mathcal{L}_{BDL}$	47.00	61.50	0.059	5.540
RT-50	\mathcal{L}_{cls}	42.90	61.90	0.059	5.010
	$\mathcal{L}_{cls} + \mathcal{L}_{CBL}$	42.60	61.80	0.059	5.240
	$\mathcal{L}_{cls} + \mathcal{L}_{BDL}$	39.60	58.00	0.059	5.170

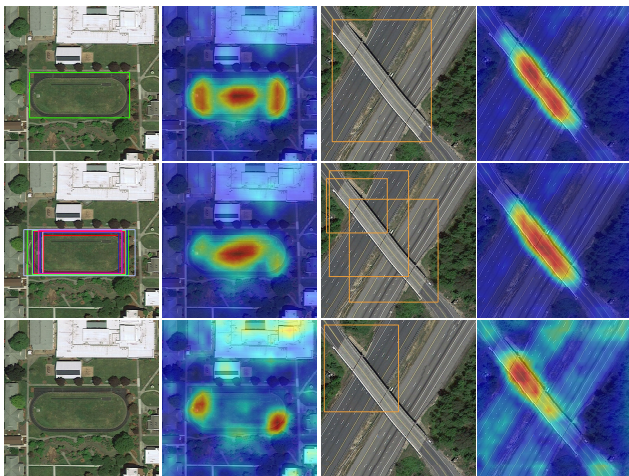


Fig. 8. Visualizations of the detected results (the Singular columns) and the corresponding features (the even columns). Here, the results on clean images are shown in the first row, the results on adversarial examples generated via \mathcal{L}_{cls} only are exhibited in the second row, and the others are the results on adversarial examples generated via $\mathcal{L}_{cls} + \mathcal{L}_{BDL}$.

results in terms of both strength and visual quality, and picking the sub-patches in a random manner exhibits the poorest performance. That is, selecting the sub-patches indeed plays a critical role in final results. Besides, as we have pointed out in Sec. I, FOD selects the most critical sub-patches via imitating the attack scheme in PAs so that to escape from the sub-optimal within the local neighborhood to find the sub-patches with the most attack potential. From this point of view, a relatively constrictive neighborhood could mislead the choice of sub-patches, resulting in poor threats.

Ablation Study on Bbox Regression Loss. Here, we vary the choice of the objective function for Bbox regression to validate the importance of our proposed BDL. Specifically, we first preserve only the objective function for the classification (denoted as \mathcal{L}_{cls} in Tab. IV) to see the importance of attacking the Bbox regression. Then, we set another competitor, *i.e.*, the CBL loss [30]. The results are reported in Tab. IV, where two conclusions can be established. First, attacking only the classification branch shows more inferior threats than taking both Bbox regression into consideration. Meanwhile, both the

TABLE V
ABLATION RESULTS REGARDING THE GRID SEGMENTATION SCHEME.

Detector	Method	mAP (%)	Recall (%)	ℓ_0 Norm	ℓ_2 Norm
FR-50	U(2)	48.50	56.90	0.065	5.180
	U(3)	49.40	57.50	0.057	4.980
	SA(1, 2, 3)	47.80	55.80	0.059	4.990
	SA(2, 3, 4)	46.30	55.20	0.064	5.210
FC-50	U(2)	48.20	63.00	0.065	5.800
	U(3)	49.10	63.50	0.057	5.510
	SA(1, 2, 3)	47.00	61.50	0.059	5.540
	SA(2, 3, 4)	45.20	59.30	0.064	5.840
RT-50	U(2)	40.30	58.20	0.065	5.290
	U(3)	41.30	59.90	0.057	5.140
	SA(1, 2, 3)	39.60	58.00	0.059	5.170
	SA(2, 3, 4)	38.20	55.10	0.064	5.380

introduction of BDL and CBL can sacrifice the ℓ_2 norm, resulting in poorer visual quality than using \mathcal{L}_{cls} only. However, the improvements of the attack strengths are significant compared to the decrease of the visual quality. Thus, introducing the attack on Bbox regression poses more threats than attacking the classification only. Second, compared to CBL, we reach the significant threats with the competitive improvements of the visual quality. We can see that the ℓ_2 norms of CBL are larger than our TPA, which echoes what we discussed in Sec. IV-C, *i.e.*, the problem of gradient inundation in CBL [30] may risk the stagnation of optimization, leading to the decrease of the attack threats. Besides quantitative results, we visualize the effects of our BDL in Fig. 8. Compared to \mathcal{L}_{cls} only, the addition of \mathcal{L}_{BDL} can result in more intense destruction, which can be reflected in the feature maps, where the attention areas are significantly interfered. By contrast, using \mathcal{L}_{cls} only could not cause significant attacks on intermediate representations.

Ablation Study on Grid Segmentation Scheme. Since grid segmentation is the first step in patch selection scheme, and there is no researches regarding the influence of different scheme on final results. To this light, we propose to provide a preliminary experimental exploration in this subsection. Recall that we provided two options for the grid segmentation in Sec. IV-B, *i.e.*, the universal scheme and the scale-adaptive scheme. In this part, we explore the influence on these schemes. Specifically, we set different parameters in $U(n)$ and $SA(n_1, n_2, n_3)$. The results are shown in Tab. V. Generally speaking, the choice of these schemes seems do not play a key role in the final results. When we take a closer look at these results, we can find that with the same visual effect, the scale-adaptive scheme shows more aggressive ability than the universal scheme, while leaving the visual quality get sacrificed slightly. For instance, SA(1, 2, 3) and U(3) achieve almost the same visual effect, but the attack strength of SA(1, 2, 3) is better than U(3). The same case can be found in the comparison between U(2) and SA(2, 3, 4), where SA(2, 3, 4) exhibits more threats than U(2). Thus, con-

sidering the attack strength and visual effect comprehensively, we choose SA(1, 2, 3) as our final choice.

VI. CONCLUSIONS

In this paper, we paid attention to PAs on object detection in O-RSIs and proposed a Threatening PA without the sacrifice of the visual quality, dubbed TPA. Specifically, to address the problem of inconsistency between local and global landscapes in existing patch selection schemes, we proposed to leverage the First-Order Difference of the objective function before and after masking to select the sub-patches to be attacked. Further, considering the problem of gradient inundation when applying existing coordinate-based loss to PAs directly, we designed an IoU-based objective function specific for PAs, dubbed Bounding box Drifting Loss, which pushes the detected bounding boxes far from the initial ones until there are no overlaps between them. Compared to the advanced competitor, the extensive evaluations have witnessed the remarkable effectiveness of our TPA. Moreover, we also replace the key factors of our TPA to see their influence in the final results. These comprehensive explorations also demonstrate the key role of our FOD and BDL. We hope this first attempt can arouse the research interest in further works regarding PAs on object detection in O-RSIs.

REFERENCES

- [1] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 2818–2826. **1, 2**
- [2] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *AAAI*, 2017. **1, 2**
- [3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Proc. Adv. Neural Inform. Process. Syst.*, vol. 28, 2015. **1, 2, 3, 6**
- [4] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *Proc. Int. Conf. Comput. Vis.*, 2019, pp. 9627–9636. **1, 2, 3, 6**
- [5] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proc. Int. Conf. Comput. Vis.*, 2017, pp. 2980–2988. **1, 2, 3, 6**
- [6] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv:2004.10934*, 2020. **1, 2, 3, 6**
- [7] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 779–788. **1, 6**
- [8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2016, pp. 770–778. **1, 2, 6**
- [9] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jul. 2017, pp. 1492–1500. **1**
- [10] G. Cheng, P. Lai, D. Gao, and J. Han, "Class attention network for image recognition," *Sci. China Inf. Sci.*, 2022. **1**
- [11] G. Cheng, J. Wang, K. Li, X. Xie, C. Lang, Y. Yao, and J. Han, "Anchor-free oriented proposal generator for object detection," *IEEE Trans. Geosci. Remote Sens.*, 2021. **1, 2, 5**
- [12] G.-S. Xia, X. Bai, J. Ding, Z. Zhu, S. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang, "Dota: A large-scale dataset for object detection in aerial images," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3974–3983. **1, 2, 5**
- [13] J. Li, Y. Liao, J. Zhang, D. Zeng, and X. Qian, "Semi-supervised degan for optical high-resolution remote sensing image scene classification," *Remote Sens.*, vol. 14, no. 17, p. 4418, 2022. **1**
- [14] B. Niu, Z. Pan, J. Wu, Y. Hu, and B. Lei, "Multi-representation dynamic adaptation network for remote sensing scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–19, 2022. **1**
- [15] L. Pei, G. Cheng, X. Sun, Q. Li, M. Zhang, and S. Miao, "Multi-scale bidirectional feature fusion for one-stage oriented object detection in aerial images," in *IGARSS*, 2021, pp. 2592–2595. **1**
- [16] 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.*, 2023. **1**
- [17] 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.*, 2022. **1**
- [18] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv:1412.6572*, 2014. **1, 2**
- [19] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv:1312.6199*, 2013. **1, 2**
- [20] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *IEEE Symp. Security Privacy*. IEEE, 2017, pp. 39–57. **1, 2**
- [21] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *Proc. Int. Conf. Learn. Represent.*, 2018. **1, 2**
- [22] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 2574–2582. **1, 2**
- [23] F. Tramèr, A. Kurakin, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, "Ensemble adversarial training: Attacks and defenses," in *Proc. Int. Conf. Learn. Represent.*, 2018. **1, 2**
- [24] G. Cheng, X. Sun, K. Li, L. Guo, and J. Han, "Perturbation-seeking generative adversarial networks: A defense framework for remote sensing image scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–11, 2021. **1**
- [25] Y. Xu, T. Bai, W. Yu, S. Chang, P. M. Atkinson, and P. Ghamisi, "Ai security for geoscience and remote sensing: Challenges and future trends," *arXiv:2212.09360*, 2022. **1**
- [26] W. Czaja, N. Fendley, M. Pekala, C. Ratto, and I.-J. Wang, "Adversarial examples in remote sensing," in *Proc. SIGSPATIAL Int. Conf. Adv. Geogr. Inf. Syst.*, 2018, pp. 408–411. **1**
- [27] Y. Xu, B. Du, and L. Zhang, "Assessing the threat of adversarial examples on deep neural networks for remote sensing scene classification: Attacks and defenses," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 2, pp. 1604–1617, 2020. **1**
- [28] C. Xie, J. Wang, Z. Zhang, Y. Zhou, L. Xie, and A. Yuille, "Adversarial examples for semantic segmentation and object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1369–1378. **1, 2**
- [29] S.-T. Chen, C. Cornelius, J. Martin, and D. H. P. Chau, "Shapeshifter: Robust physical adversarial attack on faster r-cnn object detector," in *Proc. ECML PKDD*. Springer, 2019, pp. 52–68. **1, 2**
- [30] Y. Li, D. Tian, X. Bian, and S. Lyu, "Robust adversarial perturbation on deep proposal-based models," in *Proc. Brit. Mach. Vis. Conf.* **1, 2, 3, 5, 8**
- [31] H. Zhang, W. Zhou, and H. Li, "Contextual adversarial attacks for object detection," IEEE, 2020, pp. 1–6. **1, 2**
- [32] O. M. Nezami, A. Chaturvedi, M. Dras, and U. Garain, "Pick-object-attack: Type-specific adversarial attack for object detection," *Comput. Vis. Image Underst.*, vol. 211, p. 103257, 2021. **1, 2**
- [33] S. Wu, T. Dai, and S.-T. Xia, "Dpattack: Diffused patch attacks against universal object detection," *arXiv:2010.11679*, 2020. **1, 2**
- [34] Y. Zhao, H. Yan, and X. Wei, "Object hider: Adversarial patch attack against object detectors," *arXiv:2010.14974*, 2020. **1, 2**
- [35] H. Huang, Y. Wang, Z. Chen, Z. Tang, W. Zhang, and K.-K. Ma, "Rpattack: Refined patch attack on general object detectors," 2021, pp. 1–6. **1, 2, 3, 4, 5, 6, 7**
- [36] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer, "Adversarial patch," *arXiv:1712.09665*, 2017. **1, 2**
- [37] X. Liu, H. Yang, Z. Liu, L. Song, H. Li, and Y. Chen, "Dpatch: An adversarial patch attack on object detectors," *arXiv:1806.02299*, 2018. **1, 3**
- [38] M. Lee and Z. Kolter, "On physical adversarial patches for object detection," *arXiv:1906.11897*, 2019. **1, 3**
- [39] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125. **2, 6**
- [40] Y. Dong, T. Pang, H. Su, and J. Zhu, "Evading defenses to transferable adversarial examples by translation-invariant attacks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 4312–4321. **2**

- [41] C. Xie, Z. Zhang, Y. Zhou, S. Bai, J. Wang, Z. Ren, and A. L. Yuille, "Improving transferability of adversarial examples with input diversity," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 2730–2739. [2](#)
- [42] J. Lin, C. Song, K. He, L. Wang, and J. E. Hopcroft, "Nesterov accelerated gradient and scale invariance for adversarial attacks," in *Proc. Int. Conf. Learn. Represent.*, 2020. [2](#)
- [43] Y. Dong, F. Liao, T. Pang, H. Su, J. Zhu, X. Hu, and J. Li, "Boosting adversarial attacks with momentum," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 9185–9193. [2](#), [5](#)
- [44] J. Chen, M. I. Jordan, and M. J. Wainwright, "Hopskipjumpattack: A query-efficient decision-based attack," in *Proc. IEEE Symp. Secur. Priv. IEEE*, 2020, pp. 1277–1294. [2](#)
- [45] W. Brendel, J. Rauber, and M. Bethge, "Decision-based adversarial attacks: Reliable attacks against black-box machine learning models," *arXiv:1712.04248*, 2017. [2](#)
- [46] X. Sun, G. Cheng, L. Pei, and J. Han, "Query-efficient decision-based attack via sampling distribution reshaping," *Pattern Recognit.*, p. 108728, 2022. [2](#)
- [47] T. Brunner, F. Diehl, M. T. Le, and A. Knoll, "Guessing smart: Biased sampling for efficient black-box adversarial attacks," in *Int. Conf. Comput. Vis.*, 2019, pp. 4958–4966. [2](#)
- [48] X. Sun, G. Cheng, H. Li, L. Pei, and J. Han, "Exploring effective data for surrogate training towards black-box attack," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 15 355–15 364. [2](#)
- [49] M. Zhou, J. Wu, Y. Liu, S. Liu, and C. Zhu, "Dast: Data-free substitute training for adversarial attacks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 234–243. [2](#)
- [50] S. H. Silva and P. Najafirad, "Opportunities and challenges in deep learning adversarial robustness: A survey," *arXiv:2007.00753*, 2020. [2](#)
- [51] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, "Synthesizing robust adversarial examples," in *Proc. Int. Conf. Mach. Learn.* PMLR, 2018, pp. 284–293. [2](#)
- [52] A. Liu, X. Liu, J. Fan, Y. Ma, A. Zhang, H. Xie, and D. Tao, "Perceptual-sensitive gan for generating adversarial patches," in *AAAI*, vol. 33, no. 01, 2019, pp. 1028–1035. [3](#)
- [53] J. Lian, S. Mei, S. Zhang, and M. Ma, "Benchmarking adversarial patch against aerial detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022. [3](#)
- [54] R. den Hollander, A. Adhikari, I. Tolios, M. van Bekkum, A. Bal, S. Hendriks, M. Kruithof, D. Gross, N. Jansen, G. Perez *et al.*, "Adversarial patch camouflage against aerial detection," in *Artificial Intelligence and Machine Learning in Defense Applications II*, vol. 11543. SPIE, 2020, pp. 77–86. [3](#)
- [55] M. Lu, Q. Li, L. Chen, and H. Li, "Scale-adaptive adversarial patch attack for remote sensing image aircraft detection," *Remote Sens.*, vol. 13, no. 20, p. 4078, 2021. [3](#)
- [56] S. Thys, W. Van Ranst, and T. Goedemé, "Fooling automated surveillance cameras: adversarial patches to attack person detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2019, pp. 0–0. [3](#)
- [57] J. Wang, A. Liu, Z. Yin, S. Liu, S. Tang, and X. Liu, "Dual attention suppression attack: Generate adversarial camouflage in physical world," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 8565–8574. [3](#)
- [58] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," in *Artif. Intell. Safety Secur.* Chapman and Hall/CRC, 2018, pp. 99–112. [4](#)
- [59] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2014, pp. 740–755. [4](#)