

ENHANCED SINGLE-SHOT DETECTOR FOR SMALL OBJECT DETECTION IN REMOTE SENSING IMAGES

Pourya Shamsolmoali, Masoumeh Zareapoor, Eric Granger, Jocelyn Chanussot, Jie Yang

Abstract—Small-object detection is a challenging problem. In the last few years, the convolution neural networks methods have been achieved considerable progress. However, the current detectors struggle with effective features extraction for small-scale objects. To address this challenge, we propose image pyramid single-shot detector (IPSSD). In IPSSD, single-shot detector is adopted combined with an image pyramid network to extract semantically strong features for generating candidate regions. The proposed network can enhance the small-scale features from a feature pyramid network. We evaluated the performance of the proposed model on two public datasets and the results show the superior performance of our model compared to the other state-of-the-art object detectors.

Index Terms—Object detection, feature pyramid network, remote sensing images.

I. INTRODUCTION

With the rapid progress in remote sensing technology, the analysis of remote sensing imagery (RSI) is a popular field because of its impact in both academic and industry. Object detection in RSI is an essential research field and has been studied and to address the practical issues, several object detection methods [1], [2], [3] have been developed. In the last few years, object detection in natural images has been highly successful due to the great progress of deep learning models, including single-shot multibox detectors (SSDs) [4], models for region-based convolutional neural networks (R-CNNs) [5], [6], [7], feature pyramid network for object detection (FPN) [8], and you only look once (YOLO) models [9], [10]. Recently the object detection methods that are used for natural images have been applied for object detection in RSI [11], [12], [13], [14], [15]. Dong et al. [11] proposed a model based on [7] and adopted transfer learning to reduce the possibility of missing small objects. In [12], the authors proposed a feature capturing network based on FPN to enhance detection accuracy by improving feature representation, and optimizing label assignment. Wang et al. [13] introduced an architecture named feature-merged single-shot detection network (FMSSD), which integrates the information of various sizes by using FPN and different atrous rates to enhance the quality of features. In [14], a model is proposed for small object detection by using low coupling regression and receptive field optimizing layer for

better estimation of Regions of Interests (RoIs). In [15], the authors proposed an architecture to extract semantically strong features in various scales and orientations for better small object detection in RSI. Nevertheless, the problem of small object detection is neglected in the existing models and there is a considerable scope to improve the models' performance. As we earlier discussed, it is challenging to accurately detect small objects that only occupy the region of 10×10 pixels in RSI. In this paper, we propose a new architecture on the basis of SSD to address the above challenges. The main contributions of this paper are as follows.

- We devise a detection pipeline for small objects by integrating an image pyramid network into SSD (IPSSD) to achieve more strong semantic features.
- We propose the rotation pooling layer to cover both the horizontal and oriented region proposals and design a tailored feature fusion model to make the extracted features fuse in a better form.
- We evaluate several recent object detection models in RSIs and the performances are stated.

II. METHODOLOGY

The SSD [4], has shown a promising detection results. In SSD, each prediction layer has different resolutions, where the shallower layers participate in small targets detection and the deeper layers are contributing in large targets prediction. Despite of its high performance, SSD can not detect small objects due to the poor semantic information in earlier layers of the SSD. To address this problem, we enhance the features maps quality by integrating SSD with our propose image pyramid network (IPN) to extract ROIs. Different from the max-pooling layer in the region proposal network (RPN) that only able to cover the horizontal region proposals, the propose rotation pooling layer can handle both the horizontal and oriented region proposals. Moreover, a feature fusion network (FFN) is devised to improve the context information. Fig. 1 illustrates the architecture of IPSSD.

In our propose architecture, the SSD is used as the baseline detector, in which each layer detects a specific scale objects. This implies, the shallower layers estimate small objects, while large objects are estimated by deeper layers. However, due to insufficient semantic information in the shallower layers the SSD can not accurately detect small targets. To solve this problem, we extend the SSD with the IPN to enhance the SSD's performance. As showed in Fig. 1, IPSSD contains two main streams: the standard SSD plus the IPN. For SSD, the backbone is VGG-16 and smaller convolution layers are added

P. Shamsolmoali, M.Zareapoor, J.Yang are with School of Automation, Shanghai Jiao Tong University, Shanghai, China, {pshams, mzarea, jieyang}@sjtu.edu.cn

E. Granger is with Dept. of Systems Engineering, École de technologie supérieure, Université du Québec, Canada, Eric.Granger@etsmtl.ca

J.Chanussot is with Univ. Grenoble Alpes, INRIA, CNRS, Grenoble INP, LJK, Grenoble, France.

This article has supplementary material provided by the authors and color versions of one or more figures available at <https://doi.org/XXXXXX>.

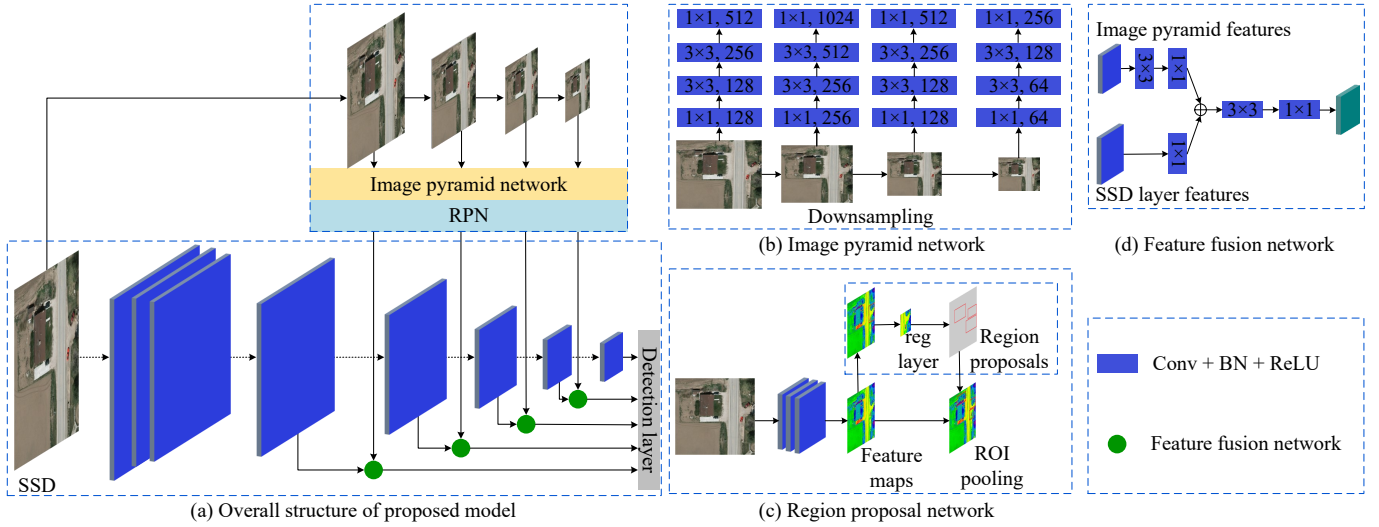


Fig. 1. Pipeline of IPSSD. (a) Network architecture. SSD is adopted with (b) the image pyramid network to extract candidate regions. (c) and (d) shows the RPN and the FFN architecture respectively.

for better feature extraction. In our model, layers of IPN in different scales are integrated into the SSD's layers using FFN.

A. Image pyramid network

The standard FPNs [8] are not computationally efficient and efficient as various scales of each image is processed by a CNN. To handle this problem, we propose an effective model to generate object candidates through RPN in IPN [7]. The network contains a down-scaling process. The IPN as input receives different size images to generate a set of box offsets. Then according to scales of box offset, the module picks a feature map in the optimum size. Firstly from input image X , the model generates multi-scale image $X_p = \{x_1, x_2, \dots, x_n\}$ by down-scaling the input image X , where, n denotes the layers of IPN. To build multi-scale feature maps, the images are processed by the IPN $S_p = \{s_1, s_2, \dots, s_n\}$, where, S_p denotes the features of each layer. The IPN has two 1×1 and two 3×3 Conv layers with different number of channels.

B. Oriented candidate regions network

The standard RPN takes the anchor to create the ROIs. Nevertheless, in RSI, the objects have a tiny scale with various orientations. Indeed, the horizontal candidates created by the standard RPN are not sufficient for challenging objects in the RSIs. To address this problem, we modify the standard RPN as follows: 1) We deleted the last three fully connected and softmax layers; 2) a network called *reg-conv* is added before the convolutional layer *conv*[5-3]; 3) a convolution kernels with a size of $3 \times 3 \times 512$ is employed to generate the 512D feature vector on the categorised feature maps; 4) the generated feature vector is processed by the *pred-score* and *pred-bbox* layers. For the oriented anchor scheme we followed [11] to create ROIs with various orientations and generate more suitable regions for a better small target detection. More specific, the candidate region $H \times W$ is splitted into several sub-regions. Thus, the sub-regions have the equivalent orientation as of the

candidate region and each sub-region has size of $S_h \leftarrow \frac{h}{H}$, $S_w \leftarrow \frac{w}{W}$. In our model the rotation region proposal for each input is defined by (x, y, h, w, θ) , in which (x, y) is the centre of bounding box, (h, w) are the height and width of bounding box (bbox) respectively, and θ denotes the standpoint from the absolute direction of the x -axis to the long side of the oriented bbox with a spatial size S . Therefore, the upper left corner of each sub-region is calculated as:

$$x_0, y_0 \leftarrow x - \frac{h}{2} + uS_h, y - \frac{w}{2} + vS_w \quad (1)$$

in which $u \in \{0, 1, \dots, H-1\}$, $v \in \{0, 1, \dots, W-1\}$ and the rotated coordinate of (x_0, y_0) computed as follows:

$$\begin{aligned} \hat{x} &\leftarrow (x_0 - x)\cos\theta + (y_0 - y)\sin\theta + x \\ \hat{y} &\leftarrow (y_0 - y)\cos\theta + (x_0 - x)\sin\theta + y \end{aligned} \quad (2)$$

C. Feature fusion network

To improve the spatial information, we introduce the FFN to combine features from the IPN layers with the SSD layers (see Fig. 2(d)). In the FFN, first, the output of IPN layer goes through a 3×3 and 1×1 conv layers, however, the output of each SSD layer only goes through a 1×1 conv layer. Then, the features of each IPN layer l_{n-1} and SSD layer l_n are combined via addition. Then, there are a 3×3 and 1×1 Conv layers for detection $d_n = R(\eta(l_{n-1}) \oplus \eta_n(l_n))$, where, $\eta_n(\cdot)$ denotes the processes including a $1 \times 1, 3 \times 3$ Conv, BN layers, and R represents the ReLU activation.

III. EXPERIMENTS

In this section, the performance of our model is evaluated compare to the other approaches on 3 classes of the DOTA [16] and NWPU VHR-10 [11] datasets for small object detection. DOTA and NWPU VHR-10 are two large RSI dataset use for object detection and contains 15 and 10 classes of objects respectively with various scales and orientations. For training and testing phase, the images are divided into the

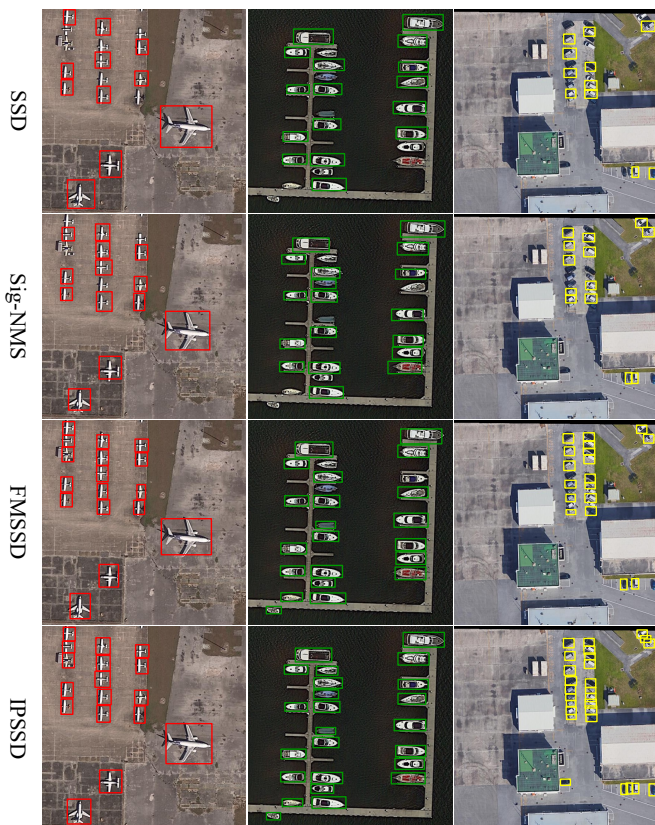


Fig. 2. Qualitative detection results comparison on small objects of DOTA.

600 × 600 pixels patches with an overlap of 100 pixels. Since, the number of images are not enough for training, to increase the number of images in training set, we apply rotation, and rescaling. All our experiments were conducted on a Tesla P-40 GPU. Our model is implemented using Pytorch. The model uses VGG-16 as the backbone which pre-trained on the ImageNet and fine-tuned on RSI dataset. The SSD uses (*Conv4*) and the fully connected layer (*FC*) layers of VGG-16 is transformed to a *Conv* layers. In SSD the last FC layer of the VGG-16 is changed to several small ranges *Conv* layers: [*Conv8*, ..., *Conv11*], with various feature sizes. The layers of SSD are merged with their corresponding layers in IPN by using the FFN. To train our model we adopt Adam as the optimizer and use the initial learning rate of 0.0001 for 30k epochs, and gradually reduce it to 0.00001 for another 20k epochs. we set the batch size to 8 and momentum to 0.9.

A. Model Comparison

To evaluate the performance of the IPSSD for small object detection, several state-of-the-art models are selected for both quantitative and qualitative comparison.

In Tables I and II, we report the detection results of our model compare to the other approaches on three small object categories of the DOTA and NWPU datasets. The detection rate of SSD on the DOTA is 70.72 mAP while processing at 64 FPS. FMSSD [13] achieves the detection rates of 78.06 mAP while processing at 22 FPS. However, IPSSD achieves 79.24 mAP detection rate while processing at 53 FPS. In

TABLE I
COMPARISON OF THE PERFORMANCE FOR SMALL OBJECT DETECTION ON THE TEST SET OF DOTA FOR HBB TASK.

Methods	Plane	SV	Ship	mAP	FPS
SSD [4]	79.64	62.13	70.41	70.72	64
Sig-NMS [11]	86.97	66.19	74.33	75.83	31
FMSSD [13]	89.14	68.32	76.73	78.06	22
IPSSD	89.09	71.39	77.26	79.24	53

TABLE II
COMPARISON OF THE PERFORMANCE FOR SMALL OBJECT DETECTION ON THE TEST SET OF NWPU.

Methods	Plane	Ship	Vehicle	mAP	FPS
SSD [4]	85.16	76.51	62.14	74.59	64
Sig-NMS [11]	90.94	81.03	78.12	83.36	31
FMSSD [13]	99.72	89.90	88.23	92.61	22
IPSSD	99.63	91.07	89.35	93.35	53

Fig. 2 we evaluate the performance of IPSSD compare to the other approaches. As the results show, IPSSD can stably produces precise results. Our model also on the NWPU dataset outperforms the state-of-the-art models. Our detector achieves 93.35% mAP detection rate. This progress is result of the following components.

- 1) By combining the IPN into the SSD, we create an architecture where each image scale is featured and resulted in improving the performance of our detector.
- 2) The FFN can enhance the attention of our proposed model to the whole object parts, which resulted in more accurate small object detection.

Fig. 3 shows the curve plot of mean localization error and confusions with background on the DOTA. As demonstrated, IPSSD has better performance in comparison with the other models.

IV. CONCLUSION

In this paper, we proposed an effective architecture by adopting image pyramid network into SSD to extract more semantic features for small target detection in RSIs. We conducted extensive experiments on two public datasets and the results show that our model performs better than the other state-of-the-art approaches for detecting small objects.

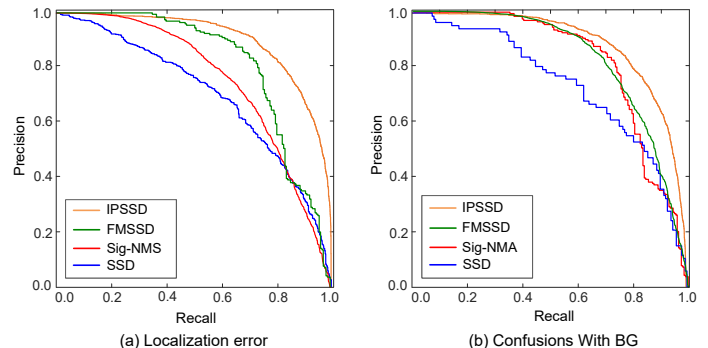


Fig. 3. Performance evaluation.

REFERENCES

- [1] 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.*, 2021.
- [2] P. Shamsolmoali, M. Zareapoor, J. Chanussot, H. Zhou, and J. Yang,
"Multipatch feature pyramid network for weakly supervised object
detection in optical remote sensing images," *IEEE Trans. Geosci. Remote
Sens.*, 2021.
- [3] 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., 2021.
- [4] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, Y. Fu, and A. Berg,
"Ssd: Single shot multibox detector," in *in Proc. ECCV*. Springer, 2016,
pp. 21–37.
- [5] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature
hierarchies for accurate object detection and semantic segmentation,"
in *in Proc. CVPR*, 2014, pp. 580–587.
- [6] R. Girshick, "Fast r-cnn," in *in Proc. ICCV*, 2015, pp. 1440–1448.
- [7] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-
time object detection with region proposal networks," in *Proc. NIPS*,
vol. 201, 2015.
- [8] T. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie,
"Feature pyramid networks for object detection," in *in Proc. CVPR*,
2017, pp. 2117–2125.
- [9] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look
once: Unified, real-time object detection," in *in Proc. CVPR*, 2016, pp.
779–788.
- [10] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement,"
arXiv preprint arXiv:1804.02767, 2018.
- [11] 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.
- [12] 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.*, 2021.
- [13] 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, 2019.
- [14] Y. Yuan and Y. Zhang, "Olcnn: An optimized low coupling network for
small objects detection," *IEEE Geosci. Remote Sens. Lett.*, 2021.
- [15] P. Shamsolmoali, M. Zareapoor, J. Chanussot, H. Zhou, and J. Yang,
"Rotation equivariant feature image pyramid network for object detec-
tion in optical remote sensing imagery," *IEEE Trans. Geosci. Remote
Sens.*, 2021.
- [16] 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 *in Proc. CVPR*, 2018, pp. 3974–3983.