

# Learning to Reduce Information Bottleneck for Object Detection in Aerial Images

Yuchen Shen, Dong Zhang, Zhihao Song, Xuesong Jiang, Qiaolin Ye, *Member, IEEE*

**Abstract**—Object detection in aerial images is a fundamental research topic in the geoscience and remote sensing domain. However, the advanced approaches on this topic mainly focus on designing the elaborate backbones or head networks but ignore neck networks. In this letter, we first underline the importance of the neck network in object detection from the perspective of information bottleneck. Then, to alleviate the information deficiency problem in the current approaches, we propose a global semantic network (GSNet), which acts as a bridge from the backbone network to the head network in a bidirectional global pattern. Compared to the existing approaches, our model can capture the rich and enhanced image features with less computational costs. Besides, we further propose a feature fusion refinement module (FRM) for different levels of features, which are suffering from the problem of semantic gap in feature fusion. To demonstrate the effectiveness and efficiency of our approach, experiments are carried out on two challenging and representative aerial image datasets (*i.e.*, DOTA and HRSC2016). Experimental results in terms of accuracy and complexity validate the superiority of our method. The code has been open-sourced at [GSNet](#).

**Index Terms**—Information bottleneck, object detection, remote sensing scene, aerial image recognition.

## I. INTRODUCTION

Object detection in aerial images is one of the most fundamental yet challenging research topics in the community of computer vision. This topic aims at recognizing each object with a precise bounding box, which is the foundation of some potential application scenarios. Recently, deep learning based object detection methods have made dramatic progresses [6], [29], [37], thanks to the significant development of deep Convolutional Neural Networks (CNNs) on vision tasks, *e.g.*, semantic segmentation [2], [3], scene classification [4].

However, it is challenging for a standard deep CNNs model to achieve a satisfactory performance, since the ground objects in aerial images usually have the properties of tiny scale,

This work was supported by the Winter Olympic Science and Technology Service Project under Grant DA2020001, the National Science Foundation of China under Grant 62072246, and the Six Talent Peaks Projects of Jiangsu Province. (Corresponding author: Xuesong Jiang and Qiaolin Ye; Yuchen Shen and Zhihao Song contributed equally to this work.)

Y. Shen and Z. Song are with the College of Information Science and Technology, Nanjing Forestry University, Nanjing 210037, Jiangsu, China. E-mail: {shenyuchen, songzhihao}@njfu.edu.cn.

Q. Ye is with the College of Information Science and Technology, Nanjing Forestry University, Nanjing 210037, Jiangsu, China, and also with Key Laboratory Intelligent Information Processing, Nanjing Xiaozhuang University, Nanjing 211171, Jiangsu, China. E-mail: yqlcom@njfu.edu.cn.

D. Zhang is with the Department of Computer Science and Engineering, The Hong Kong University of Science and Technology, Hong Kong, China. E-mail: dongz@ust.hk.

X. Jiang is with the College of Mechanical and Electronic Engineering, Nanjing Forestry University, Nanjing 210037, Jiangsu, China. E-mail: xsjiang@njfu.edu.cn.

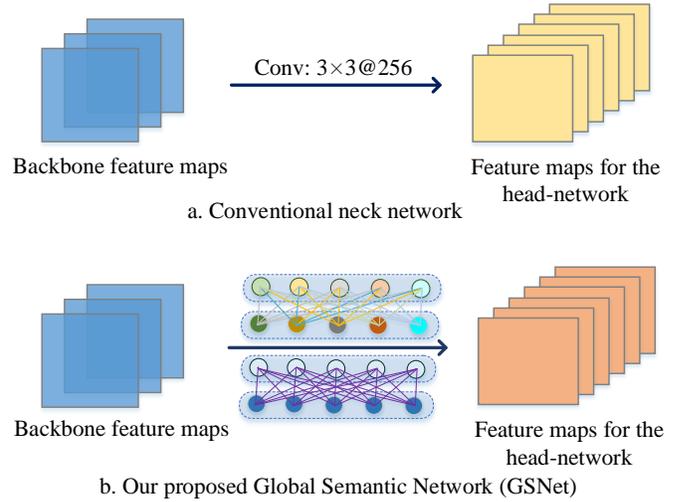


Fig. 1. An illustration for the evolution of the neck network. The conventional neck network (a) is generally based on the coarse accumulation of multiple convolutions. Contrastively, our proposed GSNet (b) can obtain affluent backbone feature cues via a global bilateral scanning operation, thus making the model more suitable for dense prediction tasks.

high density, and intricate background. Especially for some overlapping and vague scenarios, the probability of an awful result is greatly increased [4]. In this paper, we emphasize that if the head network does not change, the reason for the unsatisfactory performance is due to the insufficient feature representation [3], [6]. The information bottleneck mechanism [26], which explores the information flow between elements, can help explain this phenomenon, *i.e.*, the imperfect neck network may cause some task-related information loss [31], [32]. In particular, for the non-discriminative one, the information loss problem is much more significant [1], [26]. To this end, a large number of progressive approaches are proposed to alleviate this problem which mainly start from improving the model information content [31], [33]–[35]. For the first category of methods, *e.g.*, random shifting [31], and dilated convolution [33], receptive fields are expanded by adjusting the down-sampling to acquire global contexts. Methods in the second category, *e.g.*, AugFPN [34], and PANet [35], use multiple convolutional layers to fuse multi-scale features, such that the task-related information can be aggregated.

Although some advanced head networks have also been proposed, the existing methods ignore the fact that the neck network potentially plays a pivotal role [20], [21]. As illustrated in Figure 1 (a), the existing neck networks generally based on the coarse accumulation of multiple convolutions, which have a limited feature aggregation ability. In this letter,

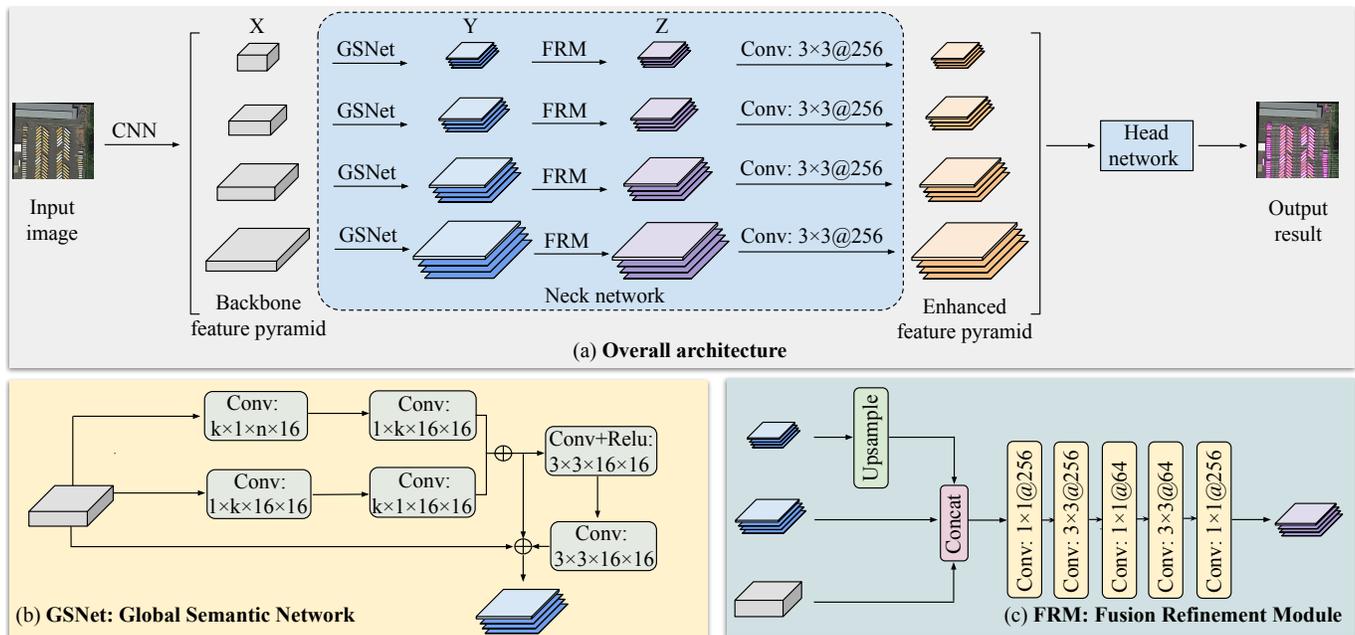


Fig. 2. Our proposed overall architecture, where Global Semantic Network (GSNet) and Fusion Refinement Module (FRM) are implemented on each layer of the backbone feature pyramid. The detailed architectures of GSNet and FRM are shown in (b) and (c), respectively.

we present a Global Semantic Network (GSNet) and a Fusion Refinement Module (FRM), which are based on the feature pyramid network [7]. As shown in Figure 1 (b), GSNet can obtain the rich backbone features via a global bilateral scanning operation, thus making the model more suitable for dense prediction tasks. Besides, FRM is an active module that boosts the model representation by propagating semantic features, such that the problem of semantic gap between features at different scales can be alleviated. To demonstrate the superiority of our model, experiments are implemented on Fast R-CNN [13] and RetinaNet [14] on two commonly used datasets (*i.e.*, DOTA [11] and HRSC2016 [12]) for oriented object detection. Results validate that our GSNet with FRM can achieve the top performance by 79.37% and 74.49% mAP for DOTA, 90.50% and 90.47% for HRSC2016.

Our contributions are summed up as 1) we emphasize that the information bottleneck causes the object detection network to lose information, which has a great performance damage. 2) we propose GSNet and FRM to reduce the information bottleneck by reconstructing the neck network in a bidirectional pattern. 3) our model achieves the competitive 79.37% and 90.50% mAP on two challenging aerial image datasets.

## II. METHODOLOGY

### A. Preliminaries

Considering that the information bottleneck will cause the loss of the effective input information, while the supervised recognition model usually expects to retain the features of the input image as much as possible [1], [26], [31], [34]. To alleviate the conflict between information bottleneck and supervised learning models, we seek to reduce the information bottleneck to minimize feature loss and enhance the network feature representation ability. The overall architecture is shown in Figure 2. Concretely, the ImageNet [15] pre-trained

ResNet [5] is adopted as the backbone following [6]. Based on which, we construct the enhanced feature pyramid via the proposed GSNet and FRM. The prediction results are finally obtained based on the enhanced feature representations with the head network, *e.g.*, Faster R-CNN [13] and RetinaNet [14].

### B. Global Semantic Network (GSNet)

Compared to small convolutional layers, a large convolution kernel convolution can bring large receptive fields [30], [36] with less computational costs, which is empirically beneficial to dense prediction tasks. The trained classical CNNs merely identify small discriminative parts with high response, while the large effective receptive fields help recognize non-discriminative regions by sensing the high-response environment around them. Besides, large kernel enables the detection model to have a tightly connected structure that copes with different transformations. In other words, features generated by convolutions with the same kernel have a stronger spatial correlation and the fully connected layer is not suitable for localization because of its spatially sensitive nature [9].

Motivated by the above observations, we present a GSNet in Figure 2 (b) that explicitly reduces feature loss and improves the model's positioning ability [36]. First, GSNet uses as large convolution kernels as possible or even global convolutions to significantly expand the effective receptive fields. But unlike many classification networks, GSNet does not have the large kernel convolutions of  $k \times k$  directly, which would significantly increase the number of parameters. Instead, our GSNet employs the combined convolutions of  $1 \times k + k \times 1$  and  $k \times 1 + 1 \times k$ . These symmetric and depth-separable combined convolutions [36] incorporate detailed contextual information while decreasing the number of parameters and computational costs, which make it more practical. Besides, GSNet is a fully convolutional network [8] with only linear operations applied

in combined convolutions. The global bilateral scanning operation can be formulated as:

$$M = \text{Conv}(\text{Conv 1D}(X)^T) + \text{Conv}(\text{Conv 1D}(X))^T, \quad (1)$$

where  $X$  as the input are the feature maps extracted from the backbone feature pyramid. Since the localization maps obtained by the recognition network cannot precisely represent the boundary of the target object, we refine the bounding box by modeling the boundary alignment as a residual structure to boost the accuracy. GSNet is introduced into the feature pyramid structure, which is closely linked to the feature maps and trained in an end-to-end manner, making the model more suitable for dense prediction tasks. Formally,

$$Y = M + R(M) + X, \quad (2)$$

$$R(M) = \text{Conv 2D}(\sigma(\text{Conv 2D}(M))), \quad (3)$$

where  $R(\cdot)$  is the residual branch, and  $\sigma$  is the ReLU [10] activation function.

### C. Fusion Refinement Module (FRM)

We proposed a novel FRM in Figure 2 (c). Direct addition is not a reasonable approach for cross-scale fusion since feature maps from the different scales have semantic information gaps. Compared to addition, channel-wise concatenation preserves more feature information, but it also increases the number of model parameters and computation. To this end,  $1 \times 1$  convolutions are adopted at intervals to reduce dimension, alleviating convolution bottlenecks. Besides, a residual branch from the backbone is introduced to inject various spatial context information. The residual structure superimposes depth features on the basis of the original features, realizing the fusion of global and local information. After that, we implement stacked convolutional layers to remove the aliasing effect caused by the interpolation, reducing information loss in the channel and enhancing the feature representational ability. The enhanced feature pyramid contains more higher-level and semantic information. Formally, it can be expressed as

$$Z = f^{1 \times 1} (f^{3 \times 3} (f^{1 \times 1} (f^{3 \times 3} (f^{1 \times 1} ([X_i, Y_i, Y_{i+1}]))))), \quad (4)$$

where  $[\cdot]$  is channel-wise concatenation, and  $X, Y$  are feature maps from the backbone feature pyramid and feature maps processed by GSNet, respectively.  $i$  represents the level of the feature pyramid, which equals 1, 2, 3.  $f^{1 \times 1}$  and  $f^{3 \times 3}$  denote the standard  $1 \times 1$  and  $3 \times 3$  convolutions.

## III. EXPERIMENTS

### A. Datasets and Evaluation Metric

DOTA [11] contains 2,806 aerial images with 15 classes, whose size varies from  $800 \times 800$  to  $4000 \times 4000$ . Figure 3 shows some training samples in this data set. HRSC2016 [12] contains 1061 images with high resolution, whose size ranges from  $300 \times 300$  to  $1500 \times 900$ . For both data sets, we randomly take 3/6 for training, 1/6 and 2/6 for validation and testing, respectively. All images are cropped into  $1024 \times 1024$  patches with a stride of 824. We use mean average precision (mAP) as

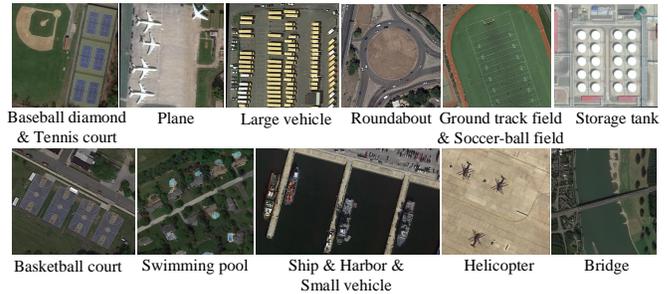


Fig. 3. Some demo training samples of the DOTA dataset [11].

TABLE I  
ABLATION STUDY.

Methods	Backbone	GSNet	FRM	mAP(%)	Params(M)	GFLOPs
Faster R-CNN [13]	Res-101	✗	✗	73.09	74.12	289.19
Faster R-CNN [13]	Res-101	✓	✗	76.61↑3.52	74.61	293.98
Faster R-CNN [13]	Res-101	✗	✓	78.98↑5.89	76.45	325.66
Faster R-CNN [13]	Res-101	✓	✓	<b>79.37</b> ↑6.28	77.91	338.14
RetinaNet [14]	Res-50	✗	✗	68.79	36.42	215.92
RetinaNet [14]	Res-50	✓	✗	70.78↑1.99	37.77	221.54
RetinaNet [14]	Res-50	✗	✓	71.10↑2.31	38.34	226.03
RetinaNet [14]	Res-50	✓	✓	<b>71.61</b> ↑2.82	39.69	231.66

the primary metric. In addition, three commonly used metrics are taken into consideration to verify the model efficiency, which are the GFLOPs, the model Parameters (Params), and the Frames Per Second (FPS).

### B. Experimental Setup

**Baselines.** We deploy two-stage Faster R-CNN [13] and one-stage RetinaNet [14] as baseline models. ResNet101 and ResNet50 (both are pre-trained on the ImageNet [15]) are adopted as backbone networks. FPN is utilized to produce an enhanced neck network. As in [13], [14], the rotated head is developed in RoI-Transformer [6] and RotatedRetinaNet [14] individually. All experimental settings strictly follow as reported in official codes for a fair comparison.

**Training Details.** We use the standard SGD [16] as the optimizer, where the learning rate is initialized to 0.005 and 0.0025 for two baselines. The weight decay and momentum are set to 0.0001 and 0.9, respectively. Models are trained for DOTA and HRSC2016 in by epochs on RTX 3060 with the batch size of 2.

### C. Ablation Study

Our ablation studies aim to validate the effectiveness and efficiency of the proposed modules on different baselines and datasets. For this purpose, we conduct a series of experiments and show some visual comparisons.

**Effectiveness of the proposed modules.** Table I shows result comparisons for effectiveness of the proposed modules. Specifically, we take Faster R-CNN based on ResNet101 as a baseline. It is observed that GSNet and FRM improve

TABLE II  
RESULT COMPARISONS WITH STATE-OF-THE-ART METHODS ON DOTA DATASET [11].

Methods	mAP(%)	PL	BD	BR	GTF	SV	LV	SH	TC	BC	ST	SBF	RA	HA	SP	HC	FPS
<i>two-stage:</i>																	
Gliding Vertex [17]	75.02	89.64	85.00	52.26	77.34	73.01	73.14	86.82	90.74	79.02	86.81	59.55	<b>70.91</b>	72.94	70.86	57.32	10.0
Oriented R-CNN [28]	76.28	88.86	83.48	55.27	76.92	74.27	82.10	87.52	90.90	85.56	85.33	65.51	66.82	74.36	70.15	57.28	<b>15.1</b>
CenterMap OBB [18]	76.03	89.83	84.41	54.60	70.25	77.66	78.32	87.19	90.66	84.89	85.27	56.46	69.23	74.13	71.56	66.06	6.3
RSDet-II [27]	76.34	89.93	84.45	53.77	74.35	71.52	78.31	78.12	<b>91.14</b>	87.35	86.93	65.64	65.17	75.35	79.74	63.31	–
SCRDet++ [19]	76.81	<b>90.05</b>	84.39	55.44	73.99	77.54	71.11	86.05	90.67	87.32	<b>87.08</b>	<b>69.62</b>	68.90	73.74	71.29	65.08	13.0
Faster R-CNN+Ours	<b>79.37</b> <sup>↑2.56</sup>	89.66	<b>86.04</b>	<b>56.25</b>	<b>79.45</b>	<b>79.07</b>	<b>84.29</b>	<b>88.40</b>	90.86	<b>88.10</b>	85.51	65.56	66.01	<b>78.70</b>	<b>79.57</b>	<b>73.02</b>	14.0
<i>one-stage:</i>																	
O <sup>2</sup> -Det [23]	71.04	89.31	82.14	47.33	61.21	71.32	74.03	78.62	90.76	82.23	81.36	60.93	60.17	58.21	66.98	<b>61.03</b>	–
R <sup>3</sup> Det [21]	71.69	89.54	81.99	48.46	62.52	70.48	74.29	77.54	90.80	81.39	83.54	<b>61.97</b>	59.82	65.44	67.46	60.05	14.0
BBAVectors [25]	72.32	88.35	79.96	50.69	62.18	78.43	78.98	<b>87.94</b>	90.85	83.58	84.35	54.13	60.24	65.22	64.28	55.70	–
DRN [24]	73.23	<b>89.71</b>	82.34	47.22	64.10	76.22	74.43	85.84	90.57	<b>86.18</b>	84.89	57.65	61.93	<b>69.30</b>	69.63	58.48	9.8
CFC-Net [22]	73.50	89.08	80.41	<b>52.41</b>	<b>70.02</b>	76.28	78.11	87.21	<b>90.89</b>	84.47	<b>85.64</b>	60.51	61.52	67.82	68.02	50.09	–
RetinaNet+Ours	<b>74.49</b> <sup>↑0.99</sup>	88.92	<b>82.79</b>	51.93	69.53	<b>79.13</b>	<b>79.16</b>	87.26	90.85	82.19	85.12	55.34	<b>66.73</b>	71.28	<b>70.46</b>	56.64	<b>20.0</b>

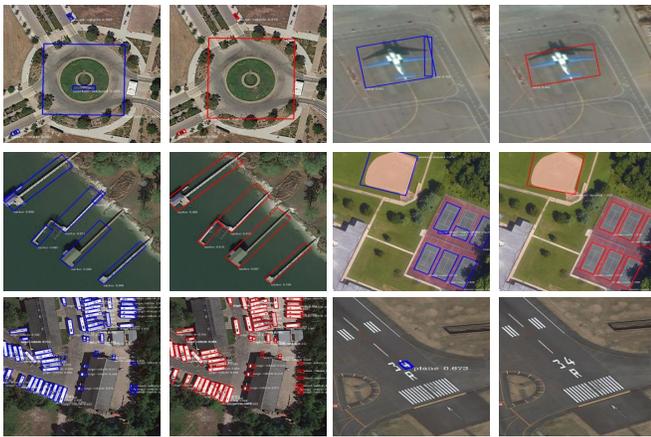


Fig. 4. Visualization of oriented detection results of baseline (blue boxes) and Faster R-CNN + Ours (red boxes) on DOTA dataset [11].

bounding box mAP by 3.52% and 5.89%, respectively. Combining GSNet and FRM, our method achieves 79.37% mAP, which is 6.28% higher than the baseline by a large margin. For model efficiency, we can observe that when GSNet is implemented on the baseline, there are only 0.49M Params and 4.79 GFLOPs. It shows that the combined convolution in GSNet could effectively control model parameters and computational cost.

**Effectiveness on different baselines.** Table I shows the results of our modules deployed to two baselines on DOTA [11]. For RetinaNet, comparing row 9 to row 6, we observe that the proposed modules bring remarkable performance enhancements (*i.e.*, 2.82% mAP). It is because our GSNet and FRM encourage each layer to preserve more features to reduce information bottleneck. This phenomenon is consistent across the HRSC2016 dataset [12]. As we mentioned under Eq. 4, the main reason for computational overheads is the introduction of additional convolutional layers in constructing FRM.

**Visualizations.** Figure 4 shows some visual comparisons of DOTA [11] between the baseline (blue boxes) and Faster R-CNN + Ours (red boxes). Faster R-CNN + Ours refers to our

TABLE III  
THE QUANTITATIVE RESULT COMPARISONS WITH STATE-OF-THE-ART METHODS ON THE TEST SET OF HRSC2016 [12].

Methods	Publication	mAP (%)	FPS
RoI Trans [6]	CVPR 2019	86.20	6.0
Gliding Vertex [17]	TPAMI 2020	88.20	–
R <sup>3</sup> Det [21]	AAAI 2021	89.26	12.0
S <sup>2</sup> A-Net [20]	TGRS 2021	90.17	–
Faster R-CNN + Ours	–	<b>90.50</b> <sup>↑0.33</sup>	14.0
RetinaNet + Ours	–	<b>90.47</b> <sup>↑0.30</sup>	<b>25.1</b>

proposed network based on ResNet-101. We can intuitively observe that our proposed methods have noticeable accuracy improvement in location and boundary, *e.g.*, the roundabout, the harbor, and the plane. From the last two rows, we observe that it also enhances the recall of some small objects.

#### D. Comparisons with State-of-the-arts

In this section, we make result comparisons with the state-of-the-art methods on both DOTA [11] and HRSC2016 [12]. **Results on DOTA [11].** As shown in Table II, compared to the previous best result of 76.81% by SCRDet++ [19] (*i.e.*, the two-stage model) and 73.50% by CFC-Net [22] (*i.e.*, the one-stage model), our GSNet + FRM model ranks the first and improves 2.56% and 0.99% mAP, respectively. Concretely, some hard categories (*e.g.*, the ship, the large vehicle, the harbor, and the helicopter) have notable mAP improvements. These results indicate that our model enhances the feature presentation capabilities. With input image size of 1024×1024, our model achieves 14.0 and 20.0 FPS on two RTX 3060 GPUs, respectively. This observation can validate the efficiency of our proposed method.

**Results on HRSC2016 [12].** The quantitative result comparisons on the test set of HRSC2016 are given in Table III. We can observe that our results are markedly prevail, reaching the top performance (*i.e.*, 90.50% mAP in 14.0 FPS and 90.47% mAP in 25.1 FPS) among all the state-of-the-art methods, which surpass the previous best model by 0.33% and 0.30% mAP with comparable inference speed, individually.

#### IV. CONCLUSIONS

In this letter, we first analyze the existing problems in object detection from the perspective of information bottleneck. Then, we propose a simple GSNet and FRM for enhancing feature representations for the neck network in object detection of aerial images. Extensive experiments on two challenging datasets confirm the superiority of our GSNet + FRM model. The main limitation is that although our network greatly improves the detection rate of small targets, it still fails to detect tiny ones that are difficult to be distinguished by naked eyes. In the future, we will consider adjusting the network structure to overcome this issue and applying our proposed methods to more computer vision tasks, *e.g.*, semantic segmentation, video object detection.

#### REFERENCES

- [1] D. Zhang, Y. Sun, Q. Ye, and J. Tang, "Recursive discriminative subspace learning with l1-norm distance constraint," *IEEE Transactions on Cybernetics*, vol. 50, no. 5, pp. 2138–2151, 2018.
- [2] D. Zhang, H. Zhang, J. Tang, X.-S. Hua, and Q. Sun, "Causal intervention for weakly-supervised semantic segmentation," in *Advances in Neural Information Processing Systems*, 2020.
- [3] D. Zhang, H. Zhang, J. Tang, M. Wang, X. Hua, and Q. Sun, "Self-regulation for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021.
- [4] D. Zhang, N. Li, and Q. Ye, "Positional context aggregation network for remote sensing scene classification," *IEEE Geoscience and Remote Sensing Letters*, 2019.
- [5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [6] J. Ding, N. Xue, Y. Long, G.-S. Xia, and Q. Lu, "Learning roi transformer for oriented object detection in aerial images," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019.
- [7] 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.
- [8] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015.
- [9] Y. Wu, Y. Chen, L. Yuan, Z. Liu, L. Wang, H. Li, and Y. Fu, "Rethinking classification and localization for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2020.
- [10] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, 2011.
- [11] 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.
- [12] Z. Liu, L. Yuan, L. Weng, and Y. Yang, "A high resolution optical satellite image dataset for ship recognition and some new baselines," in *International conference on pattern recognition applications and methods*, 2017.
- [13] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015.
- [14] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- [15] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, 2015.
- [16] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, 2012.
- [17] 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 transactions on pattern analysis and machine intelligence*, 2020.
- [18] J. Wang, W. Yang, H.-C. Li, H. Zhang, and G.-S. Xia, "Learning center probability map for detecting objects in aerial images," *IEEE Transactions on Geoscience and Remote Sensing*, 2020.
- [19] X. Yang, J. Yan, X. Yang, J. Tang, W. Liao, and T. He, "Scrdet++: Detecting small, cluttered and rotated objects via instance-level feature denoising and rotation loss smoothing," *arXiv preprint arXiv:2004.13316*, 2020.
- [20] J. Han, J. Ding, J. Li, and G.-S. Xia, "Align deep features for oriented object detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2021.
- [21] X. Yang, Q. Liu, J. Yan, A. Li, Z. Zhang, and G. Yu, "R3det: Refined single-stage detector with feature refinement for rotating object," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [22] 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 Transactions on Geoscience and Remote Sensing*, 2021.
- [23] H. Wei, Y. Zhang, Z. Chang, H. Li, H. Wang, and X. Sun, "Oriented objects as pairs of middle lines," *ISPRS Journal of Photogrammetry and Remote Sensing*, 2020.
- [24] X. Pan, Y. Ren, K. Sheng, W. Dong, H. Yuan, X. Guo, C. Ma, and C. Xu, "Dynamic refinement network for oriented and densely packed object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2020.
- [25] 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 *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*, 2021.
- [26] A. M. Saxe, Y. Bansal, J. Dapello, M. Advani, A. Kolchinsky, B. D. Tracey, and D. D. Cox, "On the information bottleneck theory of deep learning," *Journal of Statistical Mechanics: Theory and Experiment*, 2019.
- [27] W. Qian, X. Yang, S. Peng, Y. Guo, and J. Yan, "Learning modulated loss for rotated object detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [28] X. Xie, G. Cheng, J. Wang, X. Yao, and J. Han, "Oriented r-cnn for object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021.
- [29] D. Liang, Z. Wei, D. Zhang, Q. Geng, L. Zhang, H. Sun, H. Zhou, M. Wei, and P. Gao, "Learning calibrated-guidance for object detection in aerial images," in *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2021.
- [30] X. Ding, X. Zhang, Y. Zhou, J. Han, G. Ding, and J. Sun, "Scaling up your kernels to 31x31: Revisiting large kernel design in cnns," *arXiv preprint arXiv:2203.06717*, 2022.
- [31] G. Zhao, J. Wang, Z. Zhang *et al.*, "Random shifting for cnn: a solution to reduce information loss in down-sampling layers," in *IJCAI*, 2017.
- [32] Q. Liu, Z. Tan, D. Chen, Q. Chu, X. Dai, Y. Chen, M. Liu, L. Yuan, and N. Yu, "Reduce information loss in transformers for pluralistic image inpainting," *arXiv preprint arXiv:2205.05076*, 2022.
- [33] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European conference on computer vision*, 2018.
- [34] C. Guo, B. Fan, Q. Zhang, S. Xiang, and C. Pan, "Augfpn: Improving multi-scale feature learning for object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [35] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- [36] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun, "Large kernel matters—improve semantic segmentation by global convolutional network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
- [37] 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 journal of photogrammetry and remote sensing*, 2018.