

Infrared Small Target Detection Based on Transformer-Based Multi-scale Fusion Attention

Jianchun Zhang^{1,a,*}, Haiyang Yang^{1,b}, Jiangfeng Sun^{1,c}

¹College of Computer Science and Technology, Henan Polytechnic University, Jiaozuo, 454000, China
^azhangjc@hpu.edu.cn, ^b15139185277@163.com, ^csunjiangfeng@hpu.edu.cn
*Corresponding author

Abstract: Infrared weak small target detection is a critical component of infrared target detection and tracking systems, with extensive applications in maritime rescue and military surveillance. However, the task is challenging due to the complex backgrounds and the small size of the targets. Convolutional neural networks (CNNs) are proficient at capturing local details but struggle with processing global context. In contrast, Transformers excel at handling global information but may not perform well with small targets when used alone. Additionally, multiple convolutional layers can lead to the loss of target information. To address these challenges, this paper presents a Transformer-based multi-scale attention network (MATNet). The model integrates Transformer architecture with CNNs to enhance small target features more effectively. It also incorporates a multi-scale pyramid feature fusion module (FPFC) to merge features across various levels and mitigate the loss of features due to multi-layer pooling. Experimental results demonstrate that MATNet achieves superior performance compared to other methods on public datasets.

Keywords: Transformer, feature fusion, infrared small target detection

1. Introduction

Infrared small target detection technology involves identifying small targets amidst cluttered infrared backgrounds, and is known for its robust anti-interference imaging capabilities. It has wide applications in early warning, maritime rescue, precision guidance, and other fields^{[1],[2]}. Although infrared imaging technology has made considerable progress in recent years, the detection process still suffers from issues such as false alarms and low target detection accuracy due to the complex background, noise, and clutter interference^[3]. It is evident that under the influence of these interference factors, infrared target detection technology remains a challenging problem^{[4],[5]}.

The Transformer architecture is known for its strong global feature representation abilities, but its performance may be limited when dealing with infrared small targets that have few features. CNNs, on the other hand, are highly effective at local feature representation. To leverage the strengths of both, a new model called MATNet is proposed. MATNet integrates CNN and Transformer capabilities and utilizes framework akin to that of UNet.

The encoder in MATNet comprises one convolutional module and three ACFT modules. The decoder includes three FPFC modules, one skip connection-dilated convolution module, and three convolutional modules. Each convolutional module consists of two standard convolutional layers.

In the ACFT module, the self-attention mechanism is redesigned by incorporating spatial attention and dilated convolution based on the Transformer architecture. The attention matrix is computed through a dual spatial attention module (DSAM) and integrated with features extracted by spatial attention and convolution operations. The FPFC module generates multi-scale feature representations through the upsampling of feature pyramids, channel attention mechanisms, and feature fusion. By integrating upsampled low-level information with high-level information, the detection performance for small targets is significantly enhanced.

In summary, this paper has several key contributions: 1) The ACFT module, which combines Transformer and CNN, effectively enhances target features by integrating local and global features; 2) The FPFC module greatly improves the detection performance of small targets by combining upsampled low-level feature images with high-level feature images.; 3) Compared to SOTA methods, MATNet exhibits superior and more robust performance in complex backgrounds.

In Section 2, we present the related work. In Section 3, we illustrate the composition architecture of our MATNet. In Section 4, we validated the effectiveness of each module and the entire network of the proposed network through experiments. Finally, in Section 5, we draw the conclusion.

2. Related work

Currently, various techniques for detecting infrared small targets have been proposed, including early traditional model-based approaches such as filtering methods^{[6]-[10]}, local contrast-based methods^{[11]-[15]}, and low-rank-based methods^{[16]-[21]}. However, in complex backgrounds, the detection performance for small targets decreases, indicating a lack of robustness. Vision-based methods are mainly suitable for scenarios where the target is bright and contrasts well with the background. Low-rank-based methods are time-consuming and prone to a high false alarm rate when dealing with infrared images of dark targets. The methods mentioned above rely on prior expert knowledge to extract handcrafted features and are not very efficient in detecting complex scenes.

With the advancement of deep learning, numerous data-driven methods have emerged to address the limitations of model-driven approaches. MDvsFA uses a generative adversarial network to balance missed detections and false alarms^[22]. ACM introduces an asymmetric context modulation fusion module to combine deep and shallow features^[23]. ISNet utilizes Taylor finite differences and bidirectional attention aggregation blocks to precisely detect the shape features of infrared targets^[24]. However, these methods lack the capability to capture global information due to their inherent limitations, which might result in noise in infrared images being detected as targets. Additionally, because the targets in infrared images are small, multiple downsampling operations can easily lead to target loss, affecting the model's detection capability.

Due to the ubiquitous presence of target ambiguity, merely extracting local features is insufficient. Therefore, some have introduced hybrid methods by incorporating Transformers into CNN structures, combining local and global information to achieve better results. IRSTFormer uses a hierarchical Vision Transformer to model long-range dependencies to suppress false alarms, but it does not sufficiently emphasize local details^[25]. IAA Net simply connects the local patch outputs of a simple CNN with the original Transformer, resulting in limited feature extraction, especially in blurred scenes^[26]. Recently, RKformer applied ODE to ISTD tasks, employing the Runge-Kutta method to create coupled CNN-Transformer blocks that enhance infrared small targets and reduce background interference^[27]. However, the Runge-Kutta method is a linear single step method that simply concatenates the two methods, resulting in poor performance due to the lack of deeper fusion.

Therefore, we propose MATNet, which integrates CNN and Transformer to enhance the interaction between local and global features. Additionally, the feature pyramid is used to fuse low-level and high-level feature maps, improving the detection performance for small targets.

3. Method

3.1. Overall Architecture

MATNet integrates CNN and Transformer capabilities and utilizes framework akin to that of UNet. The encoder in MATNet comprises one convolutional module and three ACFT modules. The decoder includes three FPFC modules, one skip connection-dilated convolution module, and three convolutional modules. Each convolutional module consists of two standard convolutional layers.

As illustrated in Figure 1, the skip connection extended convolution module is used as a transition layer between the encoder and decoder, while the pointwise convolution layer at the end of the decoder processes the resulting image after the first few steps to produce the final result.

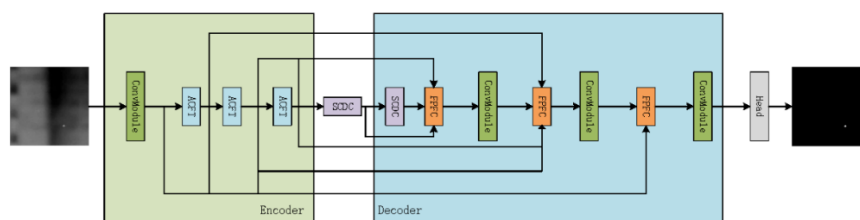


Figure 1: Overall Structure of MATNet

3.2. Attention Convolution Fusion Transformer

The Transformer possesses excellent global feature representation capabilities. However, due to the limited features of infrared small targets, merely having excellent global feature representation does not yield good results. CNN, on the other hand, tends to lose infrared small targets after multiple downsampling operations. Therefore, we incorporate the local feature extraction of CNN and combine it with Transformer, integrating local and global feature extraction. This led to the design of the ACFT module. As illustrated in Figure 2, this module redesigns the self-attention mechanism and consists of convolutional layers, dilated convolutional layers, spatial attention, and fully connected layers. The following section provides a detailed introduction to the ACFT module.

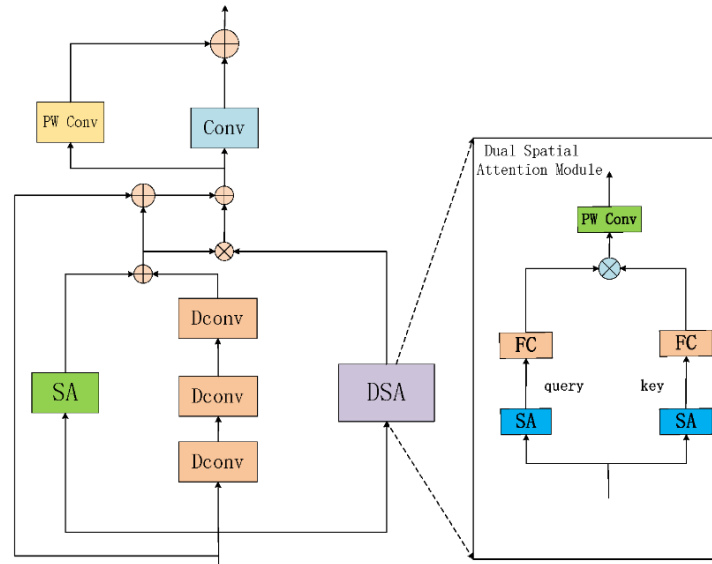


Figure 2: ACFT module diagram

Firstly, we introduce the Dual Spatial Attention (DSA) module, whose main function is to calculate the attention matrix generated by the spatial attention module. As shown in Fig 2, DSA consists of spatial attention, pointwise convolution layers, and fully connected layers. After inputting the features, two new features are obtained after passing through two spatial attention layers. These two features are then reshaped to obtain q and k . Next, both q and k will be computed through a fully connected layer, and then the computed q and k will be subjected to matrix multiplication to obtain the attention matrix. Finally, perform point by point convolution and softmax processing on the calculated attention matrix to obtain the attention matrix we need. Through continuous optimization and learning, the Dual Spatial Attention (DSA) module can more effectively perceive the location of the target.

Since the self-attention structure has been redesigned, it is necessary to calculate the value. In ACFT, the value is obtained through three dilated convolution layers and a spatial attention module. The dilation rates of the three dilated convolution layers are 2, 3, and 2, respectively. Using spatial attention can highlight key regions in an image to extract features from local areas, and the dilated convolution layers, due to their larger receptive fields, can capture more information. This compensates for the lack of detailed feature detection caused by the smaller receptive field of spatial attention. The combination of spatial attention and dilated convolution allows for the extraction of both fine-grained nearby information and distant information, resulting in more comprehensive feature acquisition.

Given an input feature $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$, after processing the input feature through spatial attention and dilated convolution layers, two new features are obtained $\mathbf{I}_{sa} \in \mathbb{R}^{H \times W \times C}$ and $\mathbf{I}_{Dconv} \in \mathbb{R}^{H \times W \times C}$. After adding the two newly computed features, a new feature $\mathbf{v} \in \mathbb{R}^{H \times W \times C}$ is obtained. Finally, multiply the attention matrix obtained from the dual space attention module with the calculated \mathbf{v} to obtain the final output $\mathbf{O}_{attention} \in \mathbb{R}^{H \times W \times C}$. The specific computation process of this operation is shown in Equation (1)-(2).

$$\mathbf{v} = \text{DC}(\mathbf{I}) + \text{SA}(\mathbf{I}) \quad (1)$$

$$\mathbf{O}_{attention} = \mathbf{v} \times \text{attention} \quad (2)$$

Here, $SA(\cdot)$ represents spatial attention, and $DC(\cdot)$ represents dilated convolution. The feature v is obtained after processing through spatial attention and dilated convolution layers. After adopting the self-attention mechanism for modeling, the model inherits the advantages of Transformer and can obtain global information of the target, more effectively collecting the features of infrared small targets, which is helpful for the detection of target features.

Finally, the detected local information and global information are fused into features, and the fused result is calculated and output through a feedforward layer. By combining CNN with Transformer, the model can simultaneously capture local and global information, enhancing the target features and reducing the likelihood of target loss.

3.3. Feature Pyramid Fusion Convolution

After processing the input feature map through the ACFT module, the target features are enhanced and become more distinct. Subsequently, finer information is obtained through the skip connection-dilated convolution. Although some fine features have been acquired, due to the tendency of infrared small target features to disappear, it is necessary to fuse multi-scale features. Therefore, we propose the FPFC module. As shown in Figure 3, the FPFC module consists of an upsampling layer, an average pooling layer, two fully connected layers, a ReLU activation function, a Sigmoid activation function, a feature fusion layer, and a standard convolutional layer.

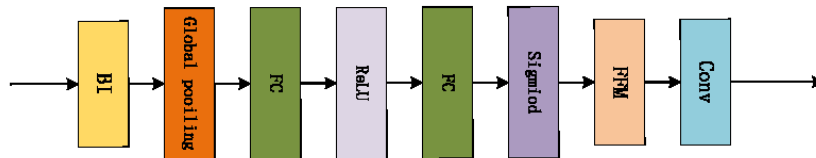


Figure 3: FPFC module diagram

The FPFC module fully leverages the advantages of multi-scale feature fusion. It matches deep features with shallow features through upsampling. Channel attention is used to assign a weight to each channel of the feature map, and these weights are used to adjust the responses of the feature maps. This results in a weighted feature representation that better captures target features and improves model performance. Finally, the weighted features are fused to obtain a more accurate and refined feature map.

4. Experiments

4.1. Experimental setup

Datasets: Our experiments utilized the publicly available IRSTD1k^[24] and NUDT^[31] datasets. We allocated 80% of the images for training and the remaining 20% for testing in each dataset.

Evaluation Metrics: We employed standard metrics for semantic segmentation, including Intersection over Union (IoU), mean IoU (mIoU), and F1 score (F1), to assess our experiments. Their definitions are as follows:

$$IoU = \frac{TP}{TP+FP+FN} \quad (3)$$

$$mIoU = \frac{1}{N} \sum_i IoU(i) \quad (4)$$

$$F1 = \frac{2TP}{2TP+FP+FN} \quad (5)$$

where N represents the total number of samples, and TP , FP , and FN denote the counts of true positive, false positive, and false negative pixels, respectively.

4.2. Implement details and compare results

Implementation Details: We used the AdamW optimizer with a learning rate of 0.0001, a weight decay coefficient of 0.01, and a momentum setting of 0.9. The loss function used is SoftIoU. For hardware, we trained the model using an NVIDIA RTX 4090 GPU. We compared our method with ACM

^[23], AGPCNet ^[29], DNANet ^[30], ABC ^[31], IPI ^[16], RIPT ^[19], and PSTNN ^[28].

Numerical Evaluation: Our proposed MATNet achieves the best results compared to other methods on both the IRSTD1k and NUDT datasets. As shown in Table 1, on the NUDT dataset, our method's IoU is 2.0% higher than the second-best, mIoU is 2.12% higher, and F1 score is improved by 1.03%. Some traditional methods perform poorly due to the complex background of the dataset, where target features are too similar to the background. Deep learning models that use only CNN methods can only detect local features, which may lead to target feature loss in deeper detections, resulting in poor performance. Although ABC combines Transformer and CNN, it only considers the fusion of local and global features without accounting for multi-scale feature fusion, leading to suboptimal performance. Our method integrates local and global features, uses attention to enhance these features, and performs multi-scale feature fusion. From the results in the Table 1, it can be seen that our proposed model can effectively improve performance.

Table 1: IoU (%), mIoU (%), and $F_1 (\times 10^{-2})$ of different SOTA methods on the IRSTD1K and NUDT datasets.

Model	IRSTD1k			NUDT		
	IoU \uparrow	mIoU \uparrow	$F_1\uparrow$	IoU \uparrow	mIoU \uparrow	$F_1\uparrow$
IPI	14.98	34.51	26.05	37.49	48.38	54.53
RIPT	11.33	17.43	20.35	29.17	36.12	45.16
PSTNN	15.93	32.71	27.48	27.72	39.80	43.41
ACM	63.39	60.81	77.59	68.48	69.26	81.29
AGPCNet	68.81	66.18	81.52	88.71	87.48	94.02
DNANet	68.87	67.53	81.57	92.67	92.09	96.20
ABC	72.02	68.81	83.73	92.85	92.45	96.29
MATNet(our)	72.75	69.28	84.19	94.45	94.57	97.32

4.3. Visual result analysis

We demonstrated the visualization results of multiple methods on the NUDT dataset. It is evident from Fig. 4 that MATNet's visualization results are superior to those of other methods. From the first row of the visualizations, it is clear that MATNet demonstrates better target detection capabilities compared to existing algorithms. When encountering noise similar to infrared small targets, other methods show incorrect detections, indicating that our algorithm has better noise suppression capabilities. The second row indicates that, compared to existing algorithms, the proposed network is more accurate in segmentation shapes. After feature enhancement with the ACFT module, our method effectively increases the probability of detecting the target. Following the FPFC module's fusion of multi-scale features, it effectively prevents the loss of very small targets, resulting in more accurate overall detection results.

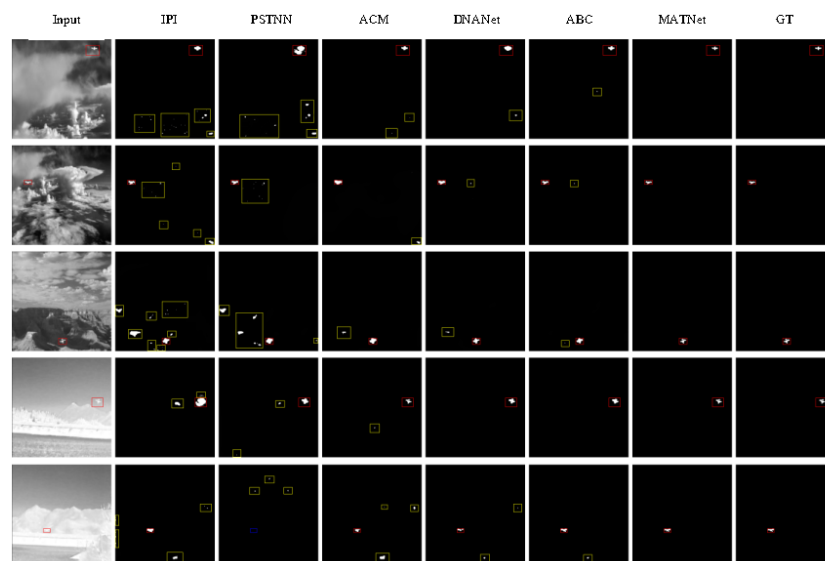


Figure 4: Visualization of various methods on the NUDT dataset. Red boxes indicate correct detections, yellow boxes denote incorrect detections, and blue boxes represent missed detections.

4.4. Ablation Study

We conducted ablation experiments on the model to verify the effectiveness of each module, and used the most basic UNet network as a benchmark network for ablation research on the NUDT dataset.

Impact of the ACFT Module: First, we removed the Dual Spatial Attention (DSA) from the ACFT module. After removing DSA, the model's local feature detection performance decreased, leading to an overall reduction in performance, as shown in Table II. When we removed the dilated convolution from the ACFT module, the model's global feature detection capability decreased. Although there was some improvement compared to the baseline, the performance was still inferior compared to using the complete ACFT module. The last row of Tab 2 indicates that the ACFT module has a significant impact on the model's performance, demonstrating that the ACFT module effectively enhances target features and improves model detection performance.

Table 2: Ablation study of the ACFT module in IoU (%) and mIoU (%), F_1 (10^{-2}).

Model	IoU \uparrow	mIoU \uparrow	$F_1\uparrow$
UNet	88.35	87.43	92.36
UNet+ACFT(no DSA)	90.61	92.76	95.07
UNet+ACFT(no DConv)	89.70	93.03	94.57
UNet+ACFT	93.61	94.19	96.70

Table 3: Ablation study of the FPFC module in IoU (%) and mIoU (%), F_1 (10^{-2}).

Model	IoU \uparrow	mIoU \uparrow	$F_1\uparrow$
UNet	88.35	87.43	92.36
UNet+ACFT	93.61	94.19	96.70
UNet+FPFC	93.03	93.79	95.38
UNet+ACFT+FPFC	94.45	94.57	97.32

Impact of the FPFC Module: we first kept the ACFT module and removed the FPFC module. As shown in Table 3, it can be observed that the performance of the module, compared to the baseline, significantly improves, indicating the effectiveness of the ACFT module. Then, we removed the ACFT module and retained the FPFC module. It is evident that the performance also improves significantly compared to the baseline, demonstrating that the FPFC module effectively enhances the detection performance of the module. The last row of the table shows that with the addition of both the ACFT and FPFC modules, the model's performance reaches its optimal state. This indicates that when both modules are present, the model achieves the best detection performance, with no conflict between them and mutual supplementation.

5. Conclusion

In this paper, we introduced the MATNet model for detecting infrared small targets. This model features the ACFT module, which integrates Transformer and CNN technologies to improve performance in capturing both local and global features. Additionally, we developed the FPFC module to address feature loss of small targets by performing multi-scale pyramid feature fusion, combining deep and shallow target features. Experimental results show that MATNet outperforms existing detection methods on the IRSTD1k and NUDT datasets.

References

- [1] S. Huang, Y. Liu, Y. He, T. Zhang, and Z. Peng, "Structure-adaptive clutter suppression for infrared small target detection: Chain-growth filtering," *Remote Sens.*, vol. 12, no. 1, 2020, Art. no. 47.
- [2] Z. Cao, X. Kong, Q. Zhu, S. Cao, and Z. Peng, "Infrared dim target detection via mode-k1k2 extension tensor tubal rank under complex ocean environment," *ISPRS J. Photogrammetry Remote Sens.*, vol. 181, pp. 167–190, 2021.
- [3] X. Kong, C. Yang, S. Cao, C. Li, and Z. Peng, "Infrared small target detection via nonconvex tensor fibered rank approximation," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–21, 2021.
- [4] X. Wang, Z. Peng, P. Zhang, and Y. He, "Infrared small target detection via nonnegativity-constrained variational mode decomposition," *IEEE Geosci. Remote Sens. Lett.*, vol. 14, no. 10, pp. 1700–1704, Oct. 2017.
- [5] Y. Liu and Z. Peng, "Infrared small target detection based on resampling-guided image model,"

- IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2021.
- [6] Suyog D Deshpande, Meng Hwa Er, Ronda Venkateswarlu, and Philip Chan, “Max-mean and max-median filters for detection of small targets,” in *Signal and Data Processing of Small Targets 1999*. SPIE, 1999, vol. 3809, pp. 74–83.
- [7] K Zhao and X Kong, “Background noise suppression in small targets infrared images and its method discussion,” *Optics and Optoelectronic Technology*, vol. 2, no. 2, pp. 9–12, 2004.
- [8] Fereshteh Seyed Marvasti, Mohammad Reza Mosavi, and Mahdi Nasiri, “Flying small target detection in ir images based on adaptive toggle operator,” *IET Computer Vision*, vol. 12, no. 4, pp. 527–534, 2018.
- [9] TS Anju and NR Nelwin Raj, “Shearlet transform based image denoising using histogram thresholding,” in *ComNet. IEEE*, 2016, pp. 162–166.
- [10] Xiaoyang Wang, Zhenming Peng, Ping Zhang, and Yanmin He, “Infrared small target detection via nonnegativity-constrained variational mode decomposition,” *IEEE GEOSCI REMOTE S*, vol. 14, no. 10, pp. 1700–1704, 2017.
- [11] Andrea Mazzu, Pietro Morerio, Lucio Marcenaro, and Carlo S Regazzoni, “A cognitive control-inspired approach to object tracking,” *IEEE T IMAGE PROCESS*, vol. 25, no. 6, pp. 2697–2711, 2016.
- [12] Yuwen Chen and Yunhong Xin, “An efficient infrared small target detection method based on visual contrast mechanism,” *IEEE GEOSCI REMOTE S*, vol. 13, no. 7, pp. 962–966, 2016.
- [13] He Deng, Xianping Sun, Maili Liu, Chaohui Ye, and Xin Zhou, “Small infrared target detection based on weighted local difference measure,” *IEEE T GEOSCI REMOTE*, vol. 54, no. 7, pp. 4204–4214, 2016.
- [14] He Deng, Xianping Sun, Maili Liu, Chaohui Ye, and Xin Zhou, “Entropy-based window selection for detecting dim and small infrared targets,” *Pattern Recognition*, vol. 61, pp. 66–77, 2017.
- [15] Jinhui Han, Sibang Liu, Gang Qin, Qian Zhao, Honghui Zhang, and Nana Li, “A local contrast method combined with adaptive background estimation for infrared small target detection,” *IEEE GEOSCI REMOTES*, vol. 16, no. 9, pp. 1442–1446, 2019.
- [16] Chenqiang Gao, Deyu Meng, Yi Yang, Yongtao Wang, Xiaofang Zhou, and Alexander G Hauptmann, “Infrared patch-image model for small target detection in a single image,” *IEEE T IMAGE PROCESS*, vol. 22, no. 12, pp. 4996–5009, 2013.
- [17] Xuan Kong, Chunping Yang, Siying Cao, Chaohai Li, and Zhenming Peng, “Infrared small target detection via nonconvex tensor fibered rank approximation,” *IEEE T GEOSCI REMOTE*, vol. 60, pp. 1–21, 2021.
- [18] Hu Zhu, Shiming Liu, Lizhen Deng, Yansheng Li, and Fu Xiao, “Infrared small target detection via low-rank tensor completion with top-hat regularization,” *IEEE T GEOSCI REMOTE*, vol. 58, no. 2, pp. 1004–1016, 2019.
- [19] Yimian Dai and Yiquan Wu, “Reweighted infrared patch-tensor model with both nonlocal and local priors for single-frame small target detection,” *IEEE J-STARS*, vol. 10, no. 8, pp. 3752–3767, 2017.
- [20] Landan Zhang, Lingbing Peng, Tianfang Zhang, Siying Cao, and Zhenming Peng, “Infrared small target detection via non-convex rank approximation minimization joint l_2 , l_1 norm,” *Remote Sensing*, vol. 10, no. 11, pp. 1821, 2018.
- [21] Shunshun Zhong, Haibo Zhou, Xingchao Cui, Xiaobing Cao, Fan Zhang, et al., “Infrared small target detection based on local-image construction and maximum correntropy,” *Measurement*, p. 112662, 2023.
- [22] Huan Wang, Luping Zhou, and Lei Wang, “Miss detection vs. false alarm: Adversarial learning for small object segmentation in infrared images,” in *ICCV*, 2019, pp. 8509–8518.
- [23] Yimian Dai, Yiquan Wu, Fei Zhou, and Kobus Barnard, “Asymmetric contextual modulation for infrared small target detection,” in *WACV*, 2021, pp. 950–959.
- [24] Mingjin Zhang, Rui Zhang, Yuxiang Yang, Haichen Bai, Jing Zhang, and Jie Guo, “Isnet: Shape matters for infrared small target detection,” in *CVPR*, 2022, pp. 877–886.
- [25] Gao Chen, Weihua Wang, and Sirui Tan, “Irstformer: A hierarchical vision transformer for infrared small target detection,” *Remote Sensing*, vol. 14, no. 14, pp. 3258, 2022.
- [26] Kewei Wang, Shuaiyuan Du, Chengxin Liu, and Zhiguo Cao, “Interior attention-aware network for infrared small target detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022.
- [27] Mingjin Zhang, Haichen Bai, Jing Zhang, Rui Zhang, Chaoyue Wang, Jie Guo, and Xinbo Gao, “Rkformer: Runge-kutta transformer with random-connection attention for infrared small target detection,” in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 1730–1738.
- [28] Landan Zhang and Zhenming Peng, “Infrared small target detection based on partial sum of the tensor nuclear norm,” *Remote Sensing*, vol. 11, no. 4, pp. 382, 2019.

- [29] Tianfang Zhang, Siying Cao, Tian Pu, and Zhenming Peng, "Agpcnet: Attention-guided pyramid context networks for infrared small target detection," *arXiv preprint arXiv: 2111.03580*, 2021.
- [30] Boyang Li, Chao Xiao, Longguang Wang, Yingqian Wang, Zaiping Lin, Miao Li, Wei An, and Yulan Guo, "Dense nested attention network for infrared small target detection," *IEEE T IMAGE PROCESS*, 2022.
- [31] Peiwen Pan, Huan Wang, Chenyi Wang, Chang Nie, "ABC: Attention with Bilinea Correlation for Infrared Small Target Detection," *2023 IEEE International Conference on Multimedia and Expo*.