# RemoteDet-Mamba: A Hybrid Mamba-CNN Network for Multi-modal Object Detection in Remote Sensing Images

Kejun Ren*, Xin Wu*, Lianming Xu† and Li Wang*

*School of Computer Science (National Pilot Software Engineering School),
Beijing University of Posts and Telecommunications, Beijing, China, Beijing, China
Email:{kejun.ren,xin.wu,liwang}@bupt.edu.cn
† School of Electronic Engineering, Beijing University of Posts and Telecommunications, Beijing, China
Email:xulianming@bupt.edu.cn

*Abstract*—Unmanned aerial vehicle (UAV) remote sensing is widely applied in fields such as emergency response, owing to its advantages of rapid information acquisition and low cost. However, due to the effects of shooting distance and imaging mechanisms, the objects in the images present challenges such as small size, dense distribution, and low inter-class differentiation. To this end, we propose a multimodal remote sensing detection network that employs a quad-directional selective scanning fusion strategy called RemoteDet-Mamba. RemoteDet-Mamba simultaneously facilitates the learning of single-modal local features and the integration of patch-level global features across modalities, enhancing the distinguishability for small objects and utilizing local information to improve discrimination between different classes. Additionally, the use of Mamba's serial processing significantly increases detection speed. Experimental results on the DroneVehicle dataset demonstrate the effectiveness of RemoteDet-Mamba, which achieves superior detection accuracy compared to state-of-the-art methods while maintaining computational efficiency and parameter count.

*Index Terms*—Unmanned aerial vehicle, remote sensing, multimodal, object detection, mamba.

## I. INTRODUCTION

Unmanned aerial vehicle (UAV) remote sensing, with its high flexibility and low operational costs, has become a vital complement to traditional satellite remote sensing [1]. It is widely used in environmental monitoring [2], disaster response, and urban planning [3]. As shown in [4], UAVs can be equipped with various modal sensors, such as infrared, visible light, and LiDAR, to achieve all-weather monitoring of ground objects. Due to the differences in spectral range, resolution, and imaging conditions among different modal sensor data, there is an urgent need for efficient fusion methods to achieve precise, effective, and fast ground object detection.

Most existing multimodal fusion methods focus on research within the frameworks of convolutional neural networks (CNNs) and transformers [5], [6]. For example, TSFADet [7] designed a CNN-based alignment framework to address the issue of weak alignment in paired image modalities. However, CNN-based methods struggle to effectively learn global context information due to the limitations of their receptive fields. Although Transformers excel at capturing long-range dependencies and global information, their high computational load limits their application in remote sensing object detection. Consequently, researchers often employ a combination of CNN and Transformer methods to address these issues, e.g., C²Former [8] utilizes a CNN as the backbone network to extract image features.

Then, it employs the Transformer's CrossAttention module to achieve feature alignment and complementarity between modalities. However, C²Former cannot avoid the Transformer's computational complexity. Furthermore, to alleviate the computational burden, these methods typically map features to a lower dimension when computing global attention, which inevitably leads to some loss of information and affects the fusion effectiveness.

Recently, the Mamba [9] architecture has enabled networks to efficiently capture global contextual information of objects with linear computational complexity. Various enhanced versions based on the Mamba architecture have been explosively successful in the field of computer vision, including Vision Mamba [10], VMamba [11], MIM-ISTD [12], and Mamba-YOLO [13]. Recently, research related to Mamba has also found widespread application in the multimodal domain. The Fusion-Mamba [14] incorporates a Fusion-Mamba block (FMB) module that maps cross-modal features into hidden state spaces for interaction, reducing discrepancies between cross-modal features and enhancing the consistency of the fused feature representation. Sigma [15] is a Siamese Mamba semantic segmentation network. It employs Mamba for feature extraction and integrates a Mamba fusion mechanism to select critical information across different modalities. However, these methods mainly incorporate some CNN structures with original visual Mamba blocks via the query-key-value (QKV) mechanism. These methods introduce redundant parameters and fail to leverage the advantages of the Mamba architecture, limiting the effectiveness of the fusion.

To this end, we propose RemoteDet-Mamba, a novel framework specifically designed for multimodal UAV object detection in remote sensing images. RemoteDet-Mamba mainly comprises a Siamese CNN network and a Cross-modal Fusion Mamba (CFM) module. In this framework, the Siamese CNN encoder extracts multi-scale local information. The CFM module is designed based on Mamba's selective scanning 2D mechanism (SS2D), which has been introduced for the first time in UAV multimodal object detection in remote sensing images. Specifically, the CFM module employs the SS2D to perform four-directional scans of the extracted multi-scale features. This linear scanning strategy decouples dense detection objects, enabling selective feature fusion. The process exhibits linear time complexity and captures long-range dependencies at the patch level. The main contributions of this paper are specified as follows:

- We propose a UAV multimodal object detection framework in remote sensing images called RemoteDet-Mamba. This framework learns single-modality local features and facilitates multimodal patch-level global feature fusion, which enhances the distinguishability of small objects and improves inter-class differentiation.
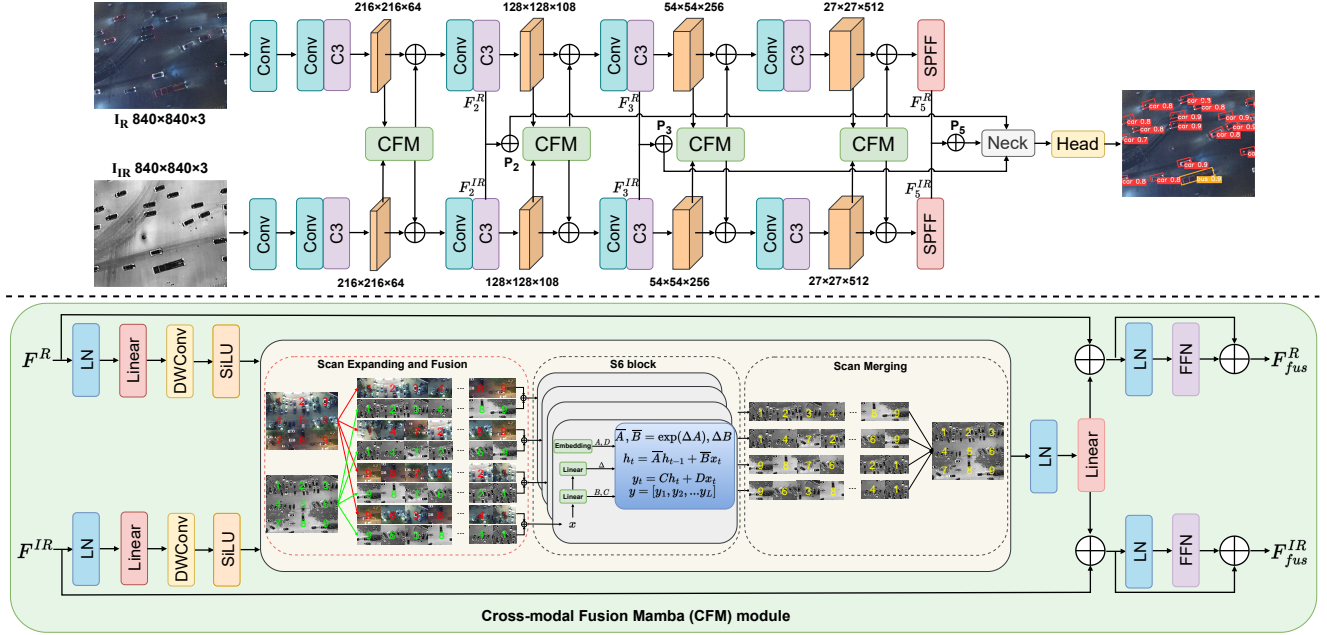- We designed a CFM module based on SS2D to perform four-

Fig. 1: The architecture of the RemoteDet-Mamba method. The top portion is the outline of the RemoteDet-Mamba. The bottom section provides a detailed view of our CFM module.

directional scans of the extracted multi-scale features at the patch level. This module decouples densely distributed small objects and facilitates the extraction of global information.

- Experimental results on the RGB-T drone remote sensing object detection dataset (Dronevehicle) demonstrate that our method achieves state-of-the-art performance while maintaining low computational and parameter costs.

## II. PROPOSED METHOD

Fig.1 gives the architecture of the proposed multimodal UAV object detection framework, which consists of a Siamese CNN network and cross-model fusion mamba (CFM) module. Specifically, the Siamese CNN network effectively extracts multi-scale features from two multimodal images, e.g., RGB and TIR. The CFM module is located between the two modality feature extraction networks [16] and achieves deep patch-level global feature fusion between the two modalities. The final fusion feature comprises the integrated outputs from the CFM module, the visible light branch network, and the infrared branch network.

### A. Siamese CNN network

Given two modalities as an example, the input image is defined as $I_s$, where $s$=1,2 corresponds to the two different modalities. These images are processed through convolutional blocks to extract their multi-scale features.

$$F_i^S = \begin{cases} \text{Conv}(I_s), & i = 0 \\ \text{Conv}(\text{C3}(F_{i-1}^s)), & i = 1, 2, 3, 4 \\ \text{SPFF}(F_{i-1}^s), & i = 5, \end{cases} \quad (1)$$

where $F_i^s$ denotes the features extracted from the $i$th layer, $F_i^1$ and $F_i^2$ represent the multi-scale features of two different modalities, respectively. For multi-scale feature extraction, the inputs to the Neck part consist of three components: the fused outputs from the 2nd, 3rd, and 5th layers of the Siamese CNN network, denoted as $P_i$.

$$P_i = F_i^1 + F_i^2, \quad i = 2, 3, 5. \quad (2)$$

After the Neck performs multi-scale feature fusion on $P_i$, it is passed to the Head, generating the final remote sensing detection results.

### B. Cross-modal Fusion Mamba (CFM) module

The design of the CFM module is based on the idea of the Mamba network. The section below the dashed line in Fig.1 illustrates the fusion process of the CFM module. Firstly, a set of multi-scale features $F^1$ and $F^2$ generated by Section II-A are sent to a LayerNorm layer to normalize the input features.

$$\overline{F}^i = \text{Linear}(\text{LN}(F^i)), F^i = \{F^1, F^2\}, \quad (3)$$

where $\text{LN}(\cdot)$ denotes the LayerNorm operation, and $\text{Linear}(\cdot)$ represents the projection operation using a linear transformation. Then, depthwise convolution is applied to the features of these two modalities to facilitate inter-channel communication.

$$f^i = \text{SiLU}(\text{DWConv}(\overline{F}^i)), \overline{F}^i = \{\overline{F}^1, \overline{F}^2\}, \quad (4)$$

where $\text{DWConv}(\cdot)$ denotes the depthwise convolution operation, and $\text{SiLU}(\cdot)$ refers to the SiLU activation function. The Mamba architecture employs a unique selective scanning mechanism that adjusts parameters based on the input data, enabling selective scanning fusion processing of the two modal features. Similar to Vmamba, the features are flattened along four directions to generate 1D sequences, each of length $HW \times C$. Afterward, deep feature fusion at the patch level is achieved through element-wise addition.

$$\begin{aligned} f_i^1 &= \text{flatten}_i(f^1), \\ f_i^2 &= \text{flatten}_i(f^2), \quad i = 1, 2, 3, 4 \\ f_i^{\text{FUS}} &= f_i^1 + f_i^2. \end{aligned} \quad (5)$$

where flatten$(\cdot)$ represents the scanning operation along the i-th direction. These one-dimensional fused sequences are individually processed by the S6 block for feature extraction, yielding four outputs, denoted as $y_1$, $y_2$, $y_3$, and $y_4$, respectively.

$$y_i = \text{S6}_i(f_i^{\text{FUS}}), \quad i = 1, 2, 3, 4, \tag{6}$$

where $S6_i$ denotes the i-th S6 block. the outputs of the S6 blocks are unfolded and recombined to generate new feature maps, denoted as $Y^{\text{FUS}}$. These fused features are then projected back to the size of the original input feature space.

$$Y^{\text{FUS}} = \sum_{i=1}^{4} \text{unflatten}_i(y_i), \tag{7}$$

where unflatten$(\cdot)$ signifies the operation of reconstructing a one-dimensional fused sequence along the i-th direction into a two-dimensional feature map. The resulting $Y^{\text{FUS}}$, along with the original inputs from the two modalities, is then processed through residual connections to yield the complementary features $\hat{F}^R$ and $\hat{F}^{IR}$.

$$\hat{F}^i = F^i + Y^{\text{FUS}}, \quad i \in \{R, IR\}. \tag{8}$$

Finally, the complementary features from each modality are processed separately through layer normalization operations and a feed-forward neural network (FFN). Then, we merge the complementary features into the multiscale features to enhance feature representation by the addition operation.

$$F_{\text{fus}}^i = \hat{F}^i + \text{FFN}(\text{LN}(\hat{F}^i)), \quad i \in \{R, \text{IR}\}, \tag{9}$$

$$\text{FFN}(\text{LN}(\hat{F}^i)) = \text{GELU}(\text{LN}(\hat{F}^i)W_1 + b_1)W_2 + b_2, \tag{10}$$

where GELU(Gaussian Error Linear Units) is the nonlinear activation function. The parameters $W_1$, $W_2$, $b_1$, and $b_2$ represent the weight matrices and bias vectors, respectively, for the linear transformations within the FFN.

### C. Loss Function

The total loss function $L_{total}$ of the proposed detection framework is composed of three parts: boundary regression loss $L_{box}$, confidence loss $L_{obj}$, classification loss $L_{cls}$. It is defined as follows:

$$L_{total} = L_{box} + L_{obj} + L_{cls}. \tag{11}$$

We employed the standard Horizontal Bounding Box (HBB) loss, as it effectively separates angular information from boundary parameter data. The Complete Intersection over Union (CIoU) was selected as the loss function $L_{box}$. For $L_{cls}$, we used Smooth Binary Cross Entropy (Smooth BCE) to enhance numerical stability. The same Smooth BCE was applied to compute the confidence loss $L_{obj}$, with CIoU added between horizontal edges to accelerate training.

## III. EXPERIMENTS AND DISCUSSIONS

### A. Data Description

A drone-based RGB-Infrared vehicle detection dataset, termed DroneVehicle [17] [1], is used for quantitative and qualitative analysis of the proposed RemoteDet-Mamba. The DroneVehicle dataset is a large-scale UAV aerial remote sensing dataset for vehicle detection and counting tasks. The DroneVehicle dataset encompasses 28,439 pairs of visible light and infrared images, totaling 953,087 detection instances, with a broad spectrum of lighting conditions ranging

[1]https://github.com/VisDrone/DroneVehicle.

from daytime to nighttime. This dataset captures diverse scenarios, including urban streets, highways, parking lots, and residential areas. It features detailed annotations with oriented bounding boxes for five vehicle categories: car, truck, bus, van, and freight car. In the experiment, DroneVehicle is divided into training, validation, and test sets, comprising $17,990$, $1,469$, and $8,980$ image pairs, respectively.

### B. Experimental Setup

In the experiment, CSPDarkNet53 is used as the backbone network for the Siamese CNN network, and data augmentation is performed using random color transformations and image reconstruction. The optimization of model parameters is the stochastic gradient descent (SGD), with an initial learning rate set at 1e-2, subsequently reduced to 2e-3. The proposed RemoteDet-Mamba architecture is implemented using the PyTorch framework on a single NVIDIA GeForce RTX 3090, with Ubuntu 18.04 and CUDA version 12.0.

### C. Evaluation Metrics

The experiment employed the commonly used object detection evaluation metric, mean Average Precision (mAP), for quantitative analysis. Here, mAP denotes the average of the average precision values computed at an Intersection over Union (IoU) threshold of 0.5.

$$mAP = \frac{\sum_{i=1}^{k} \text{AP}_i}{k}, \tag{12}$$

where AP= $\int P(R)\,dR$ represents the Area Under the Precision-Recall (P-R) Curve. K represents the number of detection categories.

### D. Experiments and Discussions

**Ablation analysis:** Table I gives the ablation analysis of the DroneVehicle dataset, including unimodal and multi-modal with different fusion strategies. In the DroneVehicle dataset, visible light and thermal infrared ground truth (GT) boxes, denoted as RGB and TIR GT, are independent. Due to the label misalignment between the two modal ground truth boxes, we retain the larger area ground truth box for training and testing in multi-modal fusion, referred to as Fusion GT. This approach facilitates an indirect alignment at the image level.

Fig.2 shows the ground truth boxes under different GT forms for training and testing. It shows that the thermal infrared ground truth boxes provide more accurate object labeling across varying lighting conditions than visible light and fused ground truth boxes. Similarly, TABLE I presents the quantitative results for different ground truth scenarios. When using only TIR or RGB images as inputs, the mean Average Precision (mAP) values are 0.757 and 0.794, respectively. With multi-modal inputs, employing addition fusion, bidirectional scanning fusion, and the proposed fusion strategy, the mAP values achieved are 0.808, 0.811, and 0.818, respectively. This indicates that the proposed fusion method improves performance by 2.4% over the thermal infrared unimodal approach, demonstrating that this strategy achieves a more comprehensive fusion of multi-modal data at the patch level.

**Performance Analysis:** In the experiment, night commonly used methods for object detection in remote sensing images were selected for both quantitative and qualitative comparisons. These methods include RetianNet [18], Refined Rotation RetinaNet($R^3$Det) [19], Single-Shot Alignment Network($S^2$ANet) [20], Faster R-CNN [21], Region of Interests Transformer(RoITransformer) [22], Rotation-equivariant Detector(ReDet) [23], Uncertainty-Aware Cross-Modality Vehicle Detection(UA-CMDet) [17], Multimodal Knowledge Distillation(MKD) [24], Two-Stream Feature Alignment Detector(TSFADet)

| GT Form | Fusion Strategy | | | | mAP (%) | mAP 0.5:0.95 (%) |
|---|---|---|---|---|---|---|
| | No fusion | Add | Bid-scanning | CFM | | |
| RGB GT | ✓ | | | | 75.7 | 50.3 |
| | | ✓ | | | 80.7 | 57.5 |
| | | | ✓ | | 80.8 | 58.2 |
| | | | | ✓ | 81.1 | 58.3 |
| Fusion GT | ✓(RGB) | | | | 67.7 | 40.5 |
| | ✓(TIR) | | | | 76.1 | 53.2 |
| | | ✓ | | | 78.3 | 55.8 |
| | | | ✓ | | 78.3 | 56.4 |
| | | | | ✓ | 78.7 | 56.1 |
| TIR GT | ✓ | | | | 79.4 | 55.9 |
| | | ✓ | | | 80.8 | 57.7 |
| | | | ✓ | | 81.1 | 58.5 |
| | | | | ✓ | **81.8** | **58.9** |

**Note:** No fusion refers to single modality, Bid-Scanning refers to bidirectional scanning fusion, mAP represents the mAP value at a threshold of 0.5, and mAP 0.5:0.95 refers to the mAP within the IoU threshold of 0.5 to 0.95.

TABLE II: Performance comparison of different methods on the DroneVehicle dataset

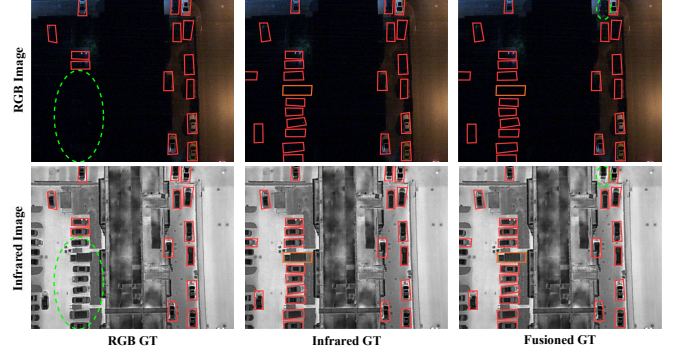| Method | Modality | Car | Trunk | FreightCar | Bus | Van | mAP(%) | Speed(fps) | Size(MB) |
|---|---|---|---|---|---|---|---|---|---|
| RetianNet [18] | | 78.5 | 34.4 | 24.1 | 69.8 | 28.8 | 47.1 | 14.53 | 218 |
| R³Det [19] | | 80.3 | 56.1 | 42.7 | 80.2 | 44.4 | 60.8 | - | - |
| S²ANet [20] | | 80.0 | 54.2 | 42.2 | 84.9 | 43.8 | 61.0 | - | - |
| Faster R-CNN [21] | RGB | 79.0 | 49.0 | 37.2 | 77.0 | 37.0 | 55.9 | 13.18 | 232 |
| RoITransformer [22] | | 61.6 | 55.1 | 42.3 | 85.5 | 44.8 | 61.6 | 11.25 | 233 |
| ReDet [23] | | 69.48 | 47.87 | 31.46 | 77.37 | 29.03 | 51.04 | 9.11 | 125 |
| RetianNet [18] | | 88.8 | 35.4 | 39.5 | 76.5 | 32.1 | 54.5 | 14.53 | 218 |
| R³Det [19] | | 89.5 | 48.3 | 16.6 | 87.1 | 39.9 | 62.3 | - | - |
| S²ANet [20] | | 89.9 | 54.5 | 55.8 | 88.9 | 48.4 | 67.5 | - | - |
| Faster R-CNN [21] | TIR | 89.4 | 53.5 | 48.3 | 87.0 | 42.6 | 64.2 | 13.18 | 232 |
| RoITransformer [22] | | 89.6 | 51.0 | 53.4 | 88.9 | 44.5 | 65.5 | 11.25 | 233 |
| ReDet [23] | | 89.47 | 53.95 | 42.82 | 79.89 | 34.39 | 60.54 | 9.11 | 125 |
| UA-CMDet [17] | | 87.51 | 60.70 | 46.80 | 87.08 | 37.95 | 64.01 | 9.12 | 234 |
| MKD [24] | | 93.49 | 62.48 | 52.73 | 91.93 | 44.50 | 69.03 | 42.39 | 242 |
| TSFADet [7] | RGB +TIR | 90.0 | 69.2 | 65.5 | 89.7 | 55.2 | 73.9 | 18.6 | 104.7 |
| C²Former [8] | | 90.2 | 68.3 | 64.4 | 89.8 | 58.5 | 74.2 | - | 118.47 |
| DMM [25] | | 90.4 | 79.8 | 68.2 | 89.9 | 68.6 | 79.4 | - | 87.97 |
| Ours | | 98.2 | 81.2 | 67.9 | 95.7 | 65.1 | 81.8 | 24.01 | 71.34 |



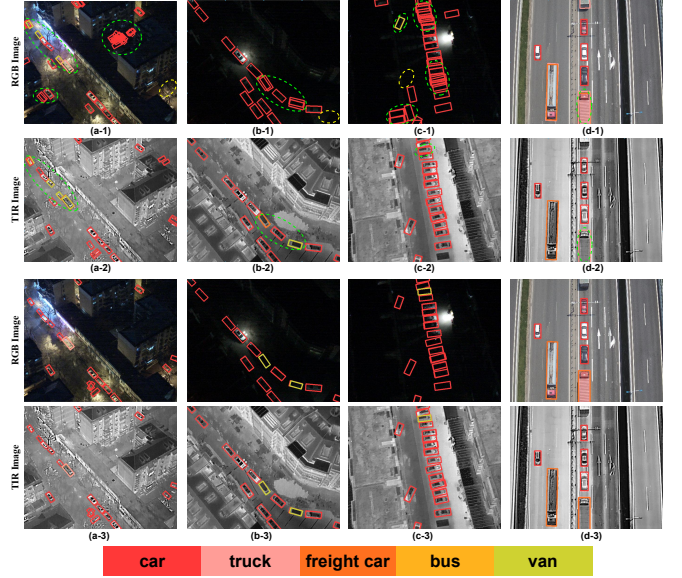Fig. 2: The visual ground truth boxes under different GT forms.



Fig. 3: The visual detection results of (a-1)-(d-1) is single RGB modality, (a-2)-(d-2) is single TIR modality, and (a-3)-(d-3) is the proposed RemoteDet-Mamba on the DroneVehicle dataset.

[7], Calibrated and Complementary Transformer(C²Former) [8], Disparity-guided Multispectral Mamba(DMM) [25]. Table II lists quantitative comparison results with existing state-of-the-art methods. Although these methods have achieved a certain level of accuracy in detection, their parameter quantity and tracking speed limit their application in the field of remote sensing. It shows that the proposed RemoteDet-Mamba achieves a mAP of 81.8% with 71.34M parameter quantity and 24.01 Frames Per Second (FPS) detection efficiency, marking a 1.2% improvement over the suboptimal DMM a [25] method.

Fig. 3 shows the detection results for unimodal and the proposed multimodal emoteDet-Mamba. The first and second rows give the detection results using visible light and infrared images, respectively. Notably, objects within the green and yellow dashed circles in Fig. 3 show that: 1) There are missed detections in the visible light images due to low-light conditions, while the infrared images exhibit false detections due to the absence of color and detailed features. 2) The densely distributed objects in remote sensing images often lead to numerous redundant prediction boxes, resulting in false detection outcomes, as depicted in (b-1)-(b-3) and (c-1)-(c-3) in Fig. 3. The third and fourth rows show the detection results of the proposed RemoteDet-Mamba utilizing two modality images. The enhanced precision of detection boxes and improved accuracy in classification demonstrate that the Mamba fusion method effectively addresses the challenge of densely distributed objects in remote sensing by employing selective feature fusion, which significantly improves detection accuracy.

## IV. CONCLUSIONS

In this paper, we propose RemoteDet-Mamba, a novel framework for multimodal UAV object detection in remote sensing images. RemoteDet-Mamba integrates the complementary strengths of CNNs and the Mamba architecture. It obtains the multi-scale local feature representation capabilities of CNNs along with the linear complexity and global receptive field provided by the Mamba architecture. For multimodal fusion, we developed the CFM module based on the Mamba architecture, which utilizes a selective scanning mechanism to execute comprehensive global scans from four directions. This method effectively isolates densely packed detection objects, facilitating refined and targeted feature fusion. Experiments on public datasets demonstrate the effectiveness of the RemoteDet-Mamba framework for multimodal UAV remote sensing object detection. This combination of detection accuracy and efficiency underscores the utility of RemoteDet-Mamba in advanced remote sensing applications.

REFERENCES

[1] D. Hong, B. Zhang, X. Li, Y. Li, C. Li, J. Yao, N. Yokoya, H. Li, P. Ghamisi, X. Jia *et al.*, "Spectralgpt: Spectral remote sensing foundation model," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, no. 8, pp. 5227–5244, 2024.

[2] X. Wu, W. Li, D. Hong, R. Tao, and Q. Du, "Deep learning for unmanned aerial vehicle-based object detection and tracking: A survey," *IEEE Geoscience and Remote Sensing Magazine*, vol. 10, no. 1, pp. 91–124, 2021.

[3] X. Wu, W. Li, D. Hong, J. Tian, R. Tao, and Q. Du, "Vehicle detection of multi-source remote sensing data using active fine-tuning network," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 167, pp. 39–53, 2020.

[4] X. Wu, D. Hong, and J. Chanussot, "Convolutional neural networks for multimodal remote sensing data classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–10, 2021.

[5] D. Hong, B. Zhang, H. Li, Y. Li, J. Yao, C. Li, M. Werner, J. Chanussot, A. Zipf, and X. X. Zhu, "Cross-city matters: A multimodal remote sensing benchmark dataset for cross-city semantic segmentation using high-resolution domain adaptation networks," *Remote Sensing of Environment*, vol. 299, p. 113856, 2023.

[6] C. Li, B. Zhang, D. Hong, J. Zhou, G. Vivone, S. Li, and J. Chanussot, "Casformer: Cascaded transformers for fusion-aware computational hyperspectral imaging," *Information Fusion*, vol. 108, p. 102408, 2024.

[7] M. Yuan, Y. Wang, and X. Wei, "Translation, scale and rotation: crossmodal alignment meets rgb-infrared vehicle detection," in *European Conference on Computer Vision*. Springer, 2022, pp. 509–525.

[8] M. Yuan and X. Wei, "C2former: Calibrated and complementary transformer for rgb-infrared object detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2024.

[9] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," *arXiv preprint arXiv:2312.00752*, 2023.

[10] L. Zhu, B. Liao, Q. Zhang, X. Wang, W. Liu, and X. Wang, "Vision mamba: Efficient visual representation learning with bidirectional state space model," *arXiv preprint arXiv:2401.09417*, 2024.

[11] Y. Liu, Y. Tian, Y. Zhao, H. Yu, L. Xie, Y. Wang, Q. Ye, and Y. Liu, "Vmamba: Visual state space model," 2024. [Online]. Available: https://arxiv.org/abs/2401.10166

[12] T. Chen, Z. Tan, T. Gong, Q. Chu, Y. Wu, B. Liu, J. Ye, and N. Yu, "Mim-istd: Mamba-in-mamba for efficient infrared small target detection," *arXiv preprint arXiv:2403.02148*, 2024.

[13] Z. Wang, C. Li, H. Xu, and X. Zhu, "Mamba yolo: Ssms-based yolo for object detection," *arXiv preprint arXiv:2406.05835*, 2024.

[14] W. Dong, H. Zhu, S. Lin, X. Luo, Y. Shen, X. Liu, J. Zhang, G. Guo, and B. Zhang, "Fusion-mamba for cross-modality object detection," *arXiv preprint arXiv:2404.09146*, 2024.

[15] Z. Wan, Y. Wang, S. Yong, P. Zhang, S. Stepputtis, K. Sycara, and Y. Xie, "Sigma: Siamese mamba network for multi-modal semantic segmentation," *arXiv preprint arXiv:2404.04256*, 2024.

[16] C. Li, B. Zhang, D. Hong, X. Jia, A. Plaza, and J. Chanussot, "Learning disentangled priors for hyperspectral anomaly detection: A coupling model-driven and data-driven paradigm," *IEEE Transactions on Neural Networks and Learning Systems*, 2024.

[17] Y. Sun, B. Cao, P. Zhu, and Q. Hu, "Drone-based rgb-infrared crossmodality vehicle detection via uncertainty-aware learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 10, pp. 6700–6713, 2022.

[18] 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, pp. 2980–2988.

[19] X. Yang, J. Yan, Z. Feng, and T. He, "R3det: Refined single-stage detector with feature refinement for rotating object," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 4, 2021, pp. 3163–3171.

[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*, vol. 60, pp. 1–11, 2021.

[21] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1137–1149, 2016.

[22] 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/CVF conference on computer vision and pattern recognition*, 2019, pp. 2849–2858.

[23] J. Han, J. Ding, N. Xue, and G.-S. Xia, "Redet: A rotation-equivariant detector for aerial object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 2786–2795.

[24] Z. Huang, W. Li, and R. Tao, "Multimodal knowledge distillation for arbitrary-oriented object detection in aerial images," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[25] M. Zhou, T. Li, C. Qiao, D. Xie, G. Wang, N. Ruan, L. Mei, and Y. Yang, "Dmm: Disparity-guided multispectral mamba for oriented object detection in remote sensing," *arXiv preprint arXiv:2407.08132*, 2024.