

# Multimodal Object Detection of UAV Remote Sensing Based on Joint Representation Optimization and Specific Information Enhancement

Jinpeng Wang<sup>1</sup>, Congan Xu<sup>1</sup>, Chunhui Zhao, Long Gao, Junfeng Wu, Yiming Yan<sup>2</sup>, *Member, IEEE*, Shou Feng<sup>3</sup>, *Member, IEEE*, and Nan Su<sup>4</sup>, *Member, IEEE*

**Abstract**—With the development of Earth observation technology, it becomes easier and easier to acquire multimodal image data at the same time. To improve the performance of a multimodal remote-sensing detection algorithm, a new fusion feature optimization detection network is proposed. The method is designed to solve the problem of performance degradation caused by the unreliability of single-modal data in multimodal remote-sensing data. The key to obtain high-quality fusion features from multimodal data with interference is to suppress single-modal redundant features and fully integrate multimodal features. The proposed method mainly includes two improvements. First, a novel joint expression optimization module is designed to enhance the target features and suppress the redundant and interference features that affect the fusion effect. In addition, we propose a novel specific information enhancement module to further enhance the discriminative feature information of targets within each modal image. Experiments on the DroneVehicle dataset show that our proposed method is state of the art on this dataset.

**Index Terms**—Joint expression optimization module (JEOM), multimodal object detection, specific information enhancement module (SIEM).

## I. INTRODUCTION

THE object detection technology of Earth observation data is widely used in military and civilian fields, such as

Manuscript received 24 January 2024; revised 5 February 2024; accepted 22 February 2024. Date of publication 6 March 2024; date of current version 8 July 2024. This work was supported in part by the National Natural Science Foundation of China under Grant 62271159, Grant 62071136, Grant 62002083, Grant 61971153, and Grant 62371153, in part by the Excellent Youth Foundation of Heilongjiang Province of China under Grant YQ2022F002, and in part by the Fundamental Research Funds for the Central Universities under Grant 3072022QBZ0805. (*Corresponding authors: Congan Xu; Nan Su.*)

Jinpeng Wang is with the College of Information and Communication Engineering, Harbin Engineering University, Harbin 150001, China, and also with Naval Aeronautical University, Yantai 264000, China (e-mail: wangjinpeng9521@hrbeu.edu.cn).

Congan Xu, Long Gao, and Junfeng Wu are with Naval Aeronautical University, Yantai 264000, China (e-mail: xcatougao@163.com; 1017730430@qq.com; patrickwu0609@163.com).

Chunhui Zhao, Yiming Yan, Shou Feng, and Nan Su are with the College of Information and Communication Engineering, Harbin Engineering University, Harbin 150001, China, and also with the Key Laboratory of Advanced Marine Communication and Information Technology, Ministry of Industry and Information Technology, Harbin 150001, China (e-mail: zhaochunhui@hrbeu.edu.cn; yanyiming@hrbeu.edu.cn; fengshou@hrbeu.edu.cn; sunan08@hrbeu.edu.cn).

Digital Object Identifier 10.1109/JSTARS.2024.3373816

intrusion warning, aerospace, etc. [1], [2]. The optical image has the texture and detail information of the target, and the infrared image can provide the temperature information of the target. These two types of target information are complementary. Currently, how to make full use of multimodal data has gradually become a new research hotspot [3]. However, in the case that there may be modal interference in multimodal remote-sensing data, how to obtain high-quality fusion features and give full play to the complementary advantages of multimodal information is a major challenge for fusion detection technology.

Optical remote-sensing images have the advantages of easy access and high resolution. Many researchers use optical remote-sensing images for object detection [4], [5]. However, in some challenging visual scenes, such as low illumination, smoke interference, etc., relying solely on optical images for detection often fails. Infrared remote-sensing images can obtain temperature information of the target without relying on visual factors in the environment. Some researchers have focused on using infrared images for object detection [6]. Due to the low resolution of infrared images, it is difficult to detect difficult targets such as low contrast, small scale, and lack of texture. In the face of some highly difficult suspected targets, even the human eye is difficult to judge. If the temperature information of the infrared image and the details and colors of the optical image can be used at the same time, the detection performance can be greatly improved. Therefore, it is worth exploring how to solve the inherent limitations of single-modal data by using multimodal complementary information to improve the detection performance [7], [8].

At present, deep learning technology is widely used in various fields [9], and it is also the focus of research in the field of multimodal target detection. Fig. 1 shows different fusion strategies in multimodal object detection algorithms in detail, which are image-level fusion, feature-level fusion, and decision-level fusion. In this figure, the red part is responsible for feature extraction, the green part represents the fusion step, and the blue part points to the object detection. Many studies have shown that implementing multimodal feature fusion in the middle layer of the network can usually obtain better multimodal detection results [10]. Nowadays, multimodal object detection methods

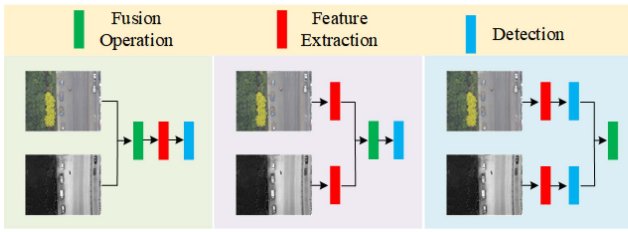


Fig. 1. Diagram of three types of fusion schemes. (a) Detection by fusion image. (b) Detection by fusion feature. (c) Fusion both detection results.

based on feature fusion have been widely concerned and become the mainstream trend.

To use optical and infrared remote-sensing images for object detection, some researchers carry out direct weighted summation of each modal branching feature. Li et al. [11] introduce illumination information to guide the fusion of multimodal image features. They use a simple CNN and prediction head to evaluate the illumination of RS optical images, and according to the evaluation results, the importance of optical image branch features is obtained, and finally, the weighted summation of each mode feature is carried out to obtain the fusion feature. Guan et al. [12] went a step further by using a deeper hierarchy to evaluate the illuminance information and then used the evaluation results to guide multimodal feature fusion. However, the branch learning efficiency of the modal weight calculation of this fusion method is low, which affects the detection performance.

In order to obtain the fusion features adaptively, Tu et al. [13] extracted the fusion feature information by layer-by-layer joining and convolution of multimodal features, which improved the overall efficiency of the algorithm. Zhang et al. [14] went further, using a more complex concatenation convolutional structure to obtain fusion features of multiple subspaces simultaneously, and finally, grouping these features along the channel dimension to obtain fusion features for prediction. This approach directly combines the features of the two branches and ignores the possible interference in the modal information.

Recently, attention structures have been widely used in various architectures because of their excellent performance. Meng and Liu [15] use a residual attention structure to perform self-attention operations on convolution-acquired fusion features to highlight their useful components. Zhao et al. [16] use transformer architecture to perform self-attention operations on multimodal features and perform attention weighting in multiple subspaces. Wang et al. [17] applied the attention structure to feature fusion networks combined with reliability weighting operations for high-level semantic information. This method combines the high efficiency of multimodal detection based on feature fusion with the robustness based on reliability weighting and has achieved great success on a remote-sensing dataset. However, in the aforementioned methods, the attention structure is only applied to the feature fusion method while the backbone network still uses the general structure for feature extraction.

In summary, multimodal object detection still faces many challenges due to the significant modal differences between different modal data. The simple weighted fusion method struggles to fully aggregate the multimodal feature information.

The fusion method based on concatenation convolution fails to consider how to suppress feature information that hinders fusion. Although the fusion methods based on attention structure design have shown excellent performance, they neglect to improve the single-modal feature extraction ability of the backbone network, ultimately reducing the efficiency of multimodal feature fusion. In addition, optical images can introduce interference in the feature fusion process, especially under low-illumination conditions. The discriminative information is the information that can distinguish the target from the background, such as the temperature difference between the target and the background in the infrared image, and the color difference in the optical image. The simple feature fusion methods may inadvertently introduce interference information and reduce the overall detection performance.

The main contributions of this article are as follows.

- 1) To address the aforementioned challenges, this article proposes a novel two-branch multimodal detection fusion feature optimization detection network (FFODNet). The FFODNet aims to adaptively fuse target feature information from multimodal remote-sensing images and achieve high-performance detection. Specifically, the method consists of two key improvements: 1) the backbone network and 2) the feature fusion module.
- 2) To fully integrate multimodal features and suppress interference information that is unfavorable to fusion, we propose a joint expression optimization module (JEOM) based on cross-concern. The JEOM is designed to adaptively extract high-quality multimodal joint expressions of objects of interest from remote-sensing data with uncertain primary and secondary states.
- 3) To enhance the discriminative feature information of the target in the single-modal image, we propose a new specific information enhancement module (SIEM) in the two-branch backbone network. The SIEM is designed to suppress irrelevant background feature information and further improve the efficiency of the subsequent feature fusion operation.

The rest of this article is organized as follows: In Section II, we describe the network structure and methods in detail. Section III gives the details of our work and experimental results and related comparisons to verify the effectiveness of our method. Finally, Section IV concludes this article.

## II. PROPOSED METHOD

The overall architecture of the proposed detection method is illustrated in Fig. 2. Since the infrared image and optical image have a similar data format, we utilize an isomorphic backbone network to extract features from the multimodal images. To ensure that the image features of each modality have the same dimension within the network, we expand the single-channel infrared image to three channels by duplicating the same value across all channels. The feature extraction process begins with a double-branch structure, which extracts features from the multimodal images. Subsequently, the JEOM utilizes these features to suppress information that is not conducive to fusion, thereby

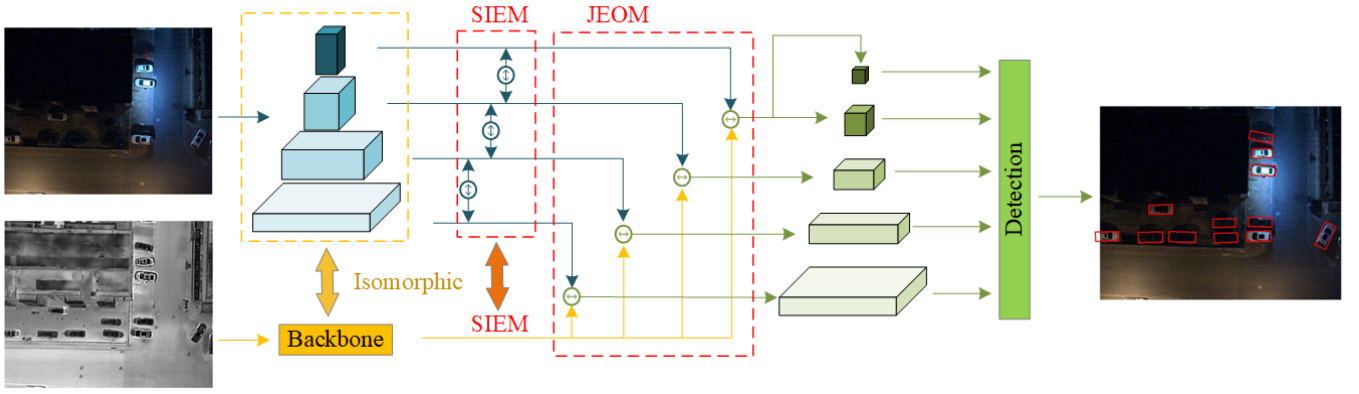


Fig. 2. Proposed multimodal object detection method mainly includes backbone, SIEM, JEOM, and detection.

extracting high-quality joint feature representations. Simultaneously, the SIEM is employed to enhance the extracted features, thereby improving the discriminative features of each individual modality and further enhancing overall performance. Finally, the fused features are passed to the detection head, where the detection results are obtained. This detection head leverages the fused features to identify and localize the target objects in the scene.

The proposed architecture combines multimodal feature fusion and single-modal feature enhancement, allowing the network to effectively capture useful complementary information in multimodal remote-sensing images and improve the performance of detection tasks. In addition, the isomorphic backbone network enables the extraction of consistent and compatible features from different modalities, facilitating the fusion process and enhancing the overall performance of the detection method.

The specifics of the JEOM and SIEM are elaborated upon in Section II-A and Section II-B, respectively.

#### A. Joint Expression Optimization Module

To obtain high-quality multimodal fusion features, we propose the JEOM based on cross attention to address possible interference in remote-sensing data. This module incorporates both single-modal features and multimodal joint features to perform attention operations. The objective of these operations is to enhance useful information while suppressing redundant information that is not beneficial for fusion. Networks generally exhibit better performance when they have access to more useful information. By incorporating cross-attention mechanisms within the JEOM, our proposed method effectively focuses on important features and enhances their representation in the fusion process. This allows the network to leverage the most relevant and discriminative information from both single-modal and multimodal features, thereby leading to improved overall performance.

The query tensor is used to search for and enhance useful features in each single mode. If the query tensor can search for the target feature more accurately, then it can also be regarded as learning a lot of useful information. To further improve the fusion efficiency, an additional step is designed during the fusion process. This involves adding the query vector, which is

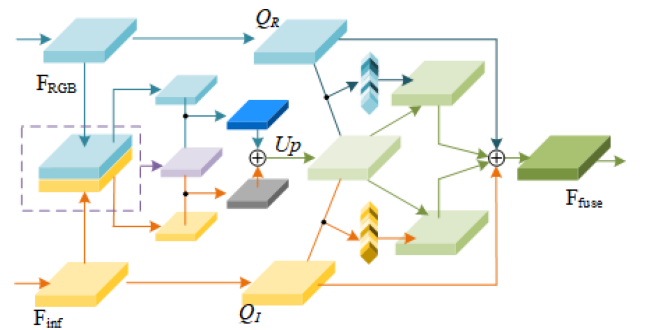


Fig. 3. Schematic diagram of JEOM.

calculated from each single-modal feature, to the fusion feature. The structure of this enhanced fusion feature is depicted in Fig. 3.

As shown in Fig. 3, the optical and infrared features are input to the fusion module and mapped as query tensors, respectively. These query tensors are then combined with the simple fusion multimodal features for attention operations. Finally, the query tensor and the enhanced fusion feature are added together to obtain the final fusion feature. The entire process leverages the attention mechanism to suppress redundant information, thereby obtaining a high-quality fusion feature expression.

By incorporating this enhanced fusion feature into the network architecture, aiming to further emphasize the useful information contained within the single-modal features, thereby increasing the proportion of useful information in the final fusion feature.

In addition, an attention-enhancing structure similar to the transformer strategy is designed, which has a strong ability to acquire high-quality fusion features [18]. In this structure, in order to carry out adaptive information fusion, we use the form of convolution instead of tensor product operation. First, the module fuses the multimodal input features, obtains the preliminary fusion features, and regards them as key tensors and value tensors. Second, the query tensor is calculated using each modal feature. Then, two query tensors and key tensors are used to calculate the weight vector, respectively, and the value tensor is weighted by the weight vector. Finally, the enhanced fusion features are added to the query tensor to obtain the final

fusion features. Finally, the enhanced fusion feature is added to the query tensor to obtain the final fusion feature. The formula of the overall calculation process is expressed as follows:

$$\text{JEOM} (F_i^{\text{RGB}}, F_i^{\text{Inf}}) = \tilde{F}_i^{\text{Inf}} + F_i^{\text{RGB}} + Q_I + Q_R \quad (1)$$

where  $\tilde{F}_i^{\text{Inf}}$  and  $F_i^{\text{RGB}}$  refer to the enhanced feature tensor, respectively.  $Q_R$  and  $Q_I$  refer to query tensors calculated from RGB image features and infrared image features, respectively

$$Q_R, Q_I = \text{Conv}_{1 \times 1} (F_i^{\text{RGB}}), \text{Conv}_{1 \times 1} (F_i^{\text{Inf}}). \quad (2)$$

$\text{Conv}_{1 \times 1}$  represents a convolution operation with kernel size 1 that does not change the dimension. Take the feature calculation of optical image as an example, the formula is as follows:

$$F_i^{\text{RGB}} = \text{Tr} (Q = Q_R, K = F_i^f, V = F_i^f). \quad (3)$$

$\text{Tr}$  is a cross-attention enhancement operation similar to the transformer strategy. It uses the query tensor and the key tensor to calculate the weight vector, and weights the value tensor as follows:

$$\text{Tr}(Q, K, V) = W(Q, K) \cdot V. \quad (4)$$

The dot multiplication in the formula refers to the multiplication of values along the channel dimension after broadcasting, so as to achieve the purpose of feature selection

$$W(Q, K) = \text{Pool}_{\text{average}}[\text{CBL}_{1 \times 1}(Q, K)]_{h,w} \quad (5)$$

$$\text{CBL}_1(X, Y) = L\_ReLU\{\text{Bn}\{\text{Conv}_1[\text{Cat}(X, Y)c]\}\}. \quad (6)$$

The previous formula represents the use of a weight vector to enhance the channel dimension of the feature graph where  $\text{Pool}_{\text{average}}$  refers to the global averaging pooling operation of feature tensors along the width and height directions.  $\text{CBL}_{1 \times 1}$  refers to the concatenation convolution operation.  $\text{Bn}$  and  $L\_ReLU$  are the batch normalized operation and the activation operation using the leaky ReLU function, respectively. In the previous equation,  $F_i^f$  refers to the preliminary fusion feature, and the calculation formula of it is as follows:

$$F_i^f = \text{Upsample} (\tilde{F}_i^r + \tilde{F}_i^i). \quad (7)$$

The primary fusion features are obtained by adding and applying operations on the calculated  $\tilde{F}_i^r$  and  $\tilde{F}_i^i$

$$\tilde{F}_i^r = \text{Conv}_{1 \times 1} (\text{Cat} (F_i^r, F_i^{\text{add}})) \quad (8)$$

$$\tilde{F}_i^i = \text{Conv}_{1 \times 1} (\text{Cat} (F_i^i, F_i^{\text{add}})). \quad (9)$$

To fully capture the discriminative information from each modality, we concatenate the optical feature and infrared feature and perform convolutions with the fusion feature separately. This process enables the network to focus not only on the combined features but also on individual modal features, thus providing essential information for subsequent fusion operations. The resulting fusion feature map is denoted that primarily emphasizes the optical modal features as  $\tilde{F}_i^r$ .  $F_i^{\text{add}}$  in the formula refers to the multimodal features of the initial fusion. This step enhances the discriminative capability of the network and contributes to

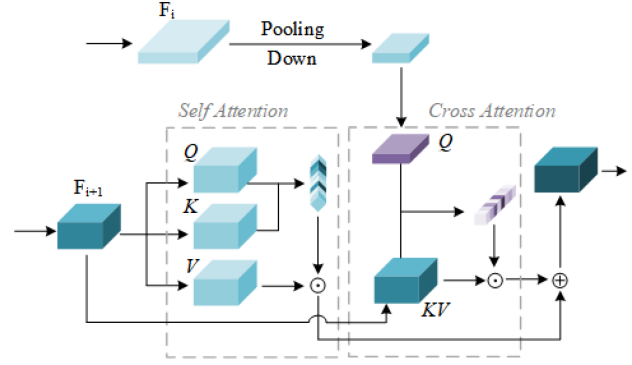


Fig. 4. Schematic diagram of SIEM.

improved multimodal object detection performance

$$F_i^{\text{add}} = \text{Conv}_{3 \times 3}^{s=2} (F_i^{\text{RGB}} + F_i^{\text{Inf}}) \quad (10)$$

$$F_i^r, F_i^i = \text{Conv}_{3 \times 3}^{s=2} (F_i^{\text{RGB}}), \text{Conv}_{3 \times 3}^{s=2} (F_i^{\text{Inf}}). \quad (11)$$

$\text{Conv}_{3 \times 3}^{s=2}$  in the formula represents a convolution operation with kernel size 3 and step size 2. In order to obtain the context information of different modal features, we carry out a downsampling operation on the features involved in the calculation when calculating the initial fusion features, so as to obtain useful information conducive to fusion in a larger area.

## B. Specific Information Enhancement Module

To enhance the feature extraction capability of single-modal remote-sensing images, we propose a self-attention SIEM module to further improve the detection performance. This module leverages multiscale feature information to enhance the discriminative features of the target, thereby improving the network's attention to target information during the feature fusion stage. We adopt the attention-enhancing structure proposed in this article within the SIEM. Given that the input features consist of both deep and shallow features that exhibit strong correlations, the operation methods of the query tensor, key tensor, and value tensor within the attention network in this section differ from those introduced in the previous section. These modifications are necessary to effectively capture and enhance the discriminative information within the multiscale feature representation. By incorporating the self-attention SIEM module, we enable the network to dynamically emphasize target-related information and enhance the discriminative power of the fused features.

The structure diagram of SIEM is shown in Fig. 4. SIEM utilizes shallower features to enhance attention toward deeper features. In many feature extraction networks, downsampling operations are performed on feature maps to reduce computation. However, this downsampling process can result in the loss of certain low-level spatial location information. SIEM addresses this issue by enhancing the discriminative features of the target and reducing the loss of spatial position information caused by downsampling during feature extraction.

In SIEM, the input shallow feature is scaled, and the module performs separate mapping operations on the deep feature and the scaled shallow feature. Specifically, key features and query



features are mapped, and their interaction generates an attention vector. This attention vector is then utilized to enhance the discriminative feature information within the deep feature.

By incorporating SIEM into the network structure, we leverage richer spatial information from the shallower features to enhance the deeper features. Simultaneously, the deeper features undergo self-attentional operations to enhance useful target feature information. The useful information component in the feature is enhanced to improve the efficiency of the subsequent feature fusion network.

The overall calculation process is as follows:

$$\text{SIEM}(F_{i-1}, F_i) = \text{Att}_{\text{self}}(F_i) + \text{Att}_{\text{cro}}(F_{i-1}, F_i). \quad (12)$$

$F_i$  and  $F_{i+1}$  represent shallow features and deep features in the process of feature extraction. The output of this module is the sum of the arithmetic result of the deep feature self-attention enhancement and the weighted feature of the shallow feature. The formula for self-attention enhancement is as follows:

$$\text{Att}_{\text{self}}(F_i) = \text{Tr}(Q_s(F_i), K_s(F_i), F_i) \quad (13)$$

$$Q_s(F_i), K_s(F_i) = \text{Conv}_{1 \times 1}(F_i), \text{Conv}_{1 \times 1}(F_i). \quad (14)$$

In the cross-attention operation, deep and shallow features need to be aligned in wide and high dimensions. In order to retain more low-level spatial information, we only use one convolutional layer for dimensional alignment. Although the method also downsamples shallow features, compared with the backbone network, the single-layer convolution operation can retain more spatial information rather than extract more semantic information. In addition, the single-layer convolution structure can also establish additional residual paths and improve the training efficiency of the method. The cross-attention formula is as follows:

$$\text{Att}_{\text{cro}}(F_{i-1}, F_i) = \text{Tr}(Q_c(F_{i-1}), K_c(F_i), F_i) \quad (15)$$

$$Q_c(F_{i-1}), K_c(F_i) = \text{Conv}_{3 \times 3}^{s=2}(F_{i-1}), \text{Conv}_{1 \times 1}(F_i). \quad (16)$$

$Q_c$  and  $K_c$  in the formula refer to the computational structure of the query tensor and the bond tensor in this part. The operation process of the function  $\text{Tr}$  in the formula is the same as the formula (4).

In the overarching process, we input the multimodal image data into the network and initiate distinct feature extraction procedures using a specialized isomorphic backbone network tailored for each modality. This step enables us to capture and emphasize unique characteristics inherent to each type of data. Following this initial extraction, we propose the SIEM to dynamically amplify the discriminative information present within the features of each individual modal. This enhancement process occurs independently within dedicated branches for each modal.

Building upon this enhancement, we activate the JEOM, which orchestrates the fusion of multimodal features at the same hierarchical level. The objective here is to ensure a harmonious integration of information across modalities, promoting a synergistic representation. Within this fusion process, careful attention is given to utilizing each modal feature map in order to selectively suppress redundant information embedded within

TABLE I  
VOLUME DIAGRAM TABLE FOR EACH DATASET

Dataset	Classes	Train set	Test set	Total
DroneVehicle	5	17 990	1 469	19 459
FLIR-aligned	3	4 129	1 013	5 142

the initial fusion feature map. This approach aims to distill and preserve only the most pertinent details, generating fusion features of elevated quality.

In the culmination of this sophisticated process, we deploy resulting high-quality fusion features for critical object detection tasks. Our comprehensive approach encompasses distinct feature extraction, individual modal enhancement, multimodal fusion, and information refinement, all working collectively to ensure the robust and effective performance of our network in discerning and identifying targets within multimodal image data. By taking advantage of the different strengths of each mode, our network extracts target discrimination features from multimodal images that capture both multimodal shared information and single-modal specific information. The individual modal enhancement module refines the features of each modality, boosting their discriminative power and facilitating more accurate object detection. The multimodal fusion process effectively integrates the enhanced features from different modalities, enabling the network to exploit the complementary information and achieve a more comprehensive understanding of the scene.

### III. EXPERIMENT

#### A. Experimental Datasets

The DroneVehicle dataset consists of 19459 pairs of RGB-infrared images, classified as vehicles, captured by camera-equipped drones [19]. Regional scenes are divided into urban roads, residential areas, and highways. The lighting conditions were night and day. We used 17990 RGB-infrared image pairs for training and 1469 pairs for validation. The overall dataset contains the following five categories of objectives: 1) car, 2) freight car, 3) truck, 4) bus, and 5) van. Some images have the problem of low contrast or low illumination, which will cause the network to be interfered by certain modal data in the feature fusion stage, which requires the feature balancing performance of the network. According to the interference degree of modal data, the dataset is divided into the following two parts: 1) the weak interference subset and 2) the strong interference subset. The subsets are shown in Fig. 5. It is important to note that we have only divided the test set.

To ensure the consistency of the image scale between the two modalities in FLIR, we conducted experiments on the ‘‘aligned’’ version [20]. The ‘‘aligned’’ FLIR contains 5142 RGB-infrared image pairs, of which we used 4129 pairs for training and 1013 pairs for testing. It covers different urban street scenes and includes three object categories: 1) bicycle, 2) car, and 3) person.

The training, testing, and overall data volumes for each dataset are shown in Table I.

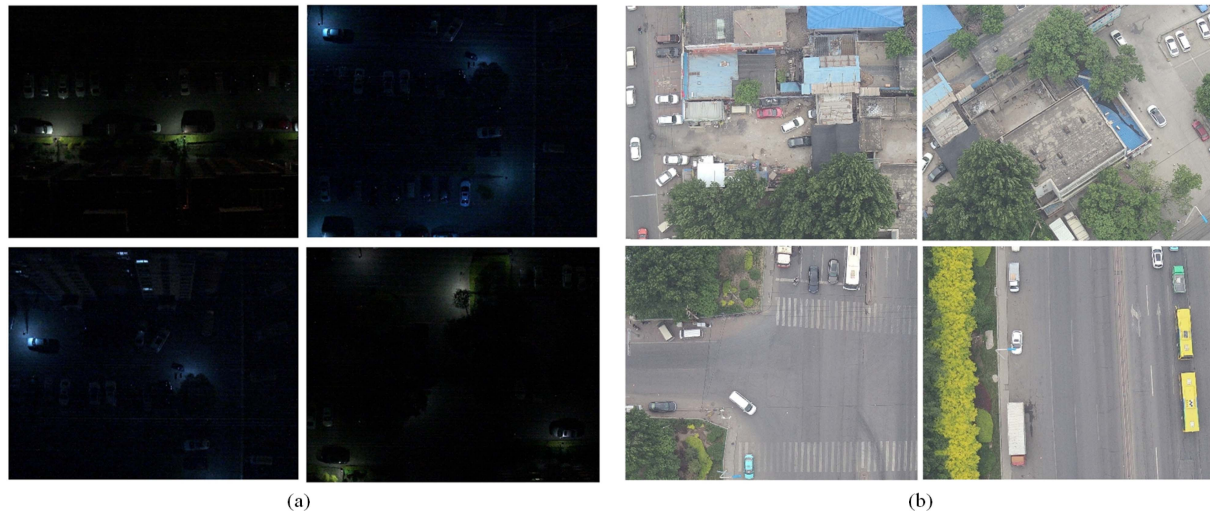


Fig. 5. Schematic samples of the weak interference subset and the strong interference subset are shown in the figure. (a) Schematic data of the strong interference subset. (b) Schematic data of the weak interference subset.

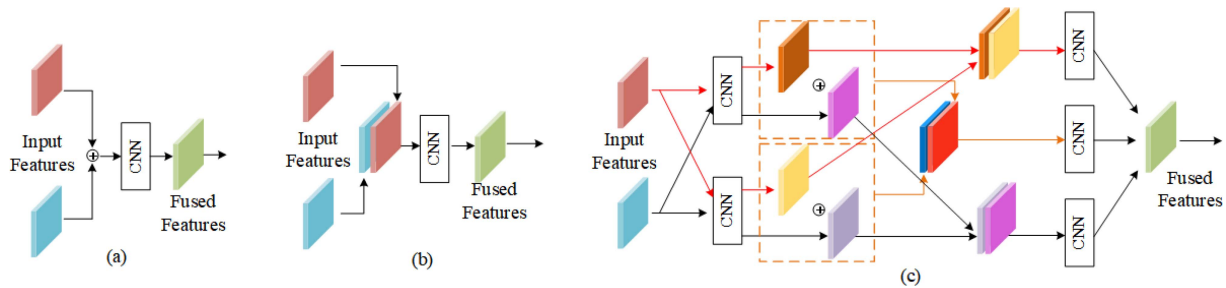


Fig. 6. Comparison of three fusion structure diagrams. (a) Direct addition of the feature map. (b) Convolution after concatenation of feature graphs. (c) Multiple interleaved concatenation and convolution of feature graphs.

### B. Implementation Details

We executed all experiments using PyTorch on a machine equipped with a GeForce RTX 3090 GPU. The optimization process employed the stochastic gradient descent algorithm, with an initial learning rate set at 0.003, an attenuation weight of 0.0001, and a momentum value of 0.9. To quantitatively assess the performance of multimodal object detection, we employed the conventional evaluation metric known as average mean precision (mAP).

### C. Performance Evaluation on DroneVehicle Dataset

Our base structure is a two-branch object detection network without SIEM and JEOM in Fig. 2. The baseline of the proposed method is improved by the single-branch Faster-RCNN [21]. To explore a better feature fusion method, we conducted some relative comparative experiments. Hong et al. [22] listed a number of multimodal feature fusion methods. The authors in [14] and [23] used pointwise addition and Concat-Conv, respectively, to carry out multimodal feature fusion, and achieved certain results. We reproduced the most complex fusion method in their paper and compared it with direct addition and concatenation convolution, two simple fusion methods, their structures are shown in Fig. 6.

These three fusion methods are called pointwise addition, Concat-Conv, and Cross-Concat-Conv in turn. Each of the methods in the table uses ResNet50 for feature extraction of single-modal images and only uses different structures for feature fusion. Finally, the fused features are used for object detection. By comparing the performance differences of various fusion methods, experiments show that the proposed fusion method is superior to the earlier fusion methods.

Table II shows that a more complex feature fusion network can obtain better multimodal fusion features, thus improving the performance of the object detection method. To improve the learning efficiency of the network and make a more explicit performance comparison, we refer to the detection network using direct addition operation in the fusion part as the baseline.

Table II shows the comparison between the proposed method and the current object detection method with multimodal feature fusion capability. RISNet and UA-CMDet are good fusion object detection algorithms at present [17], [19], both using a mixture of feature-level fusion and decision-level fusion strategies, but their performance is still inferior to our proposed method.

For a better comparison, the detection structure of the method in Fig. 6(a) is used as our baseline. Our improved mAP improves by about 16% compared to baseline results and is significantly higher than the single-modal object detection



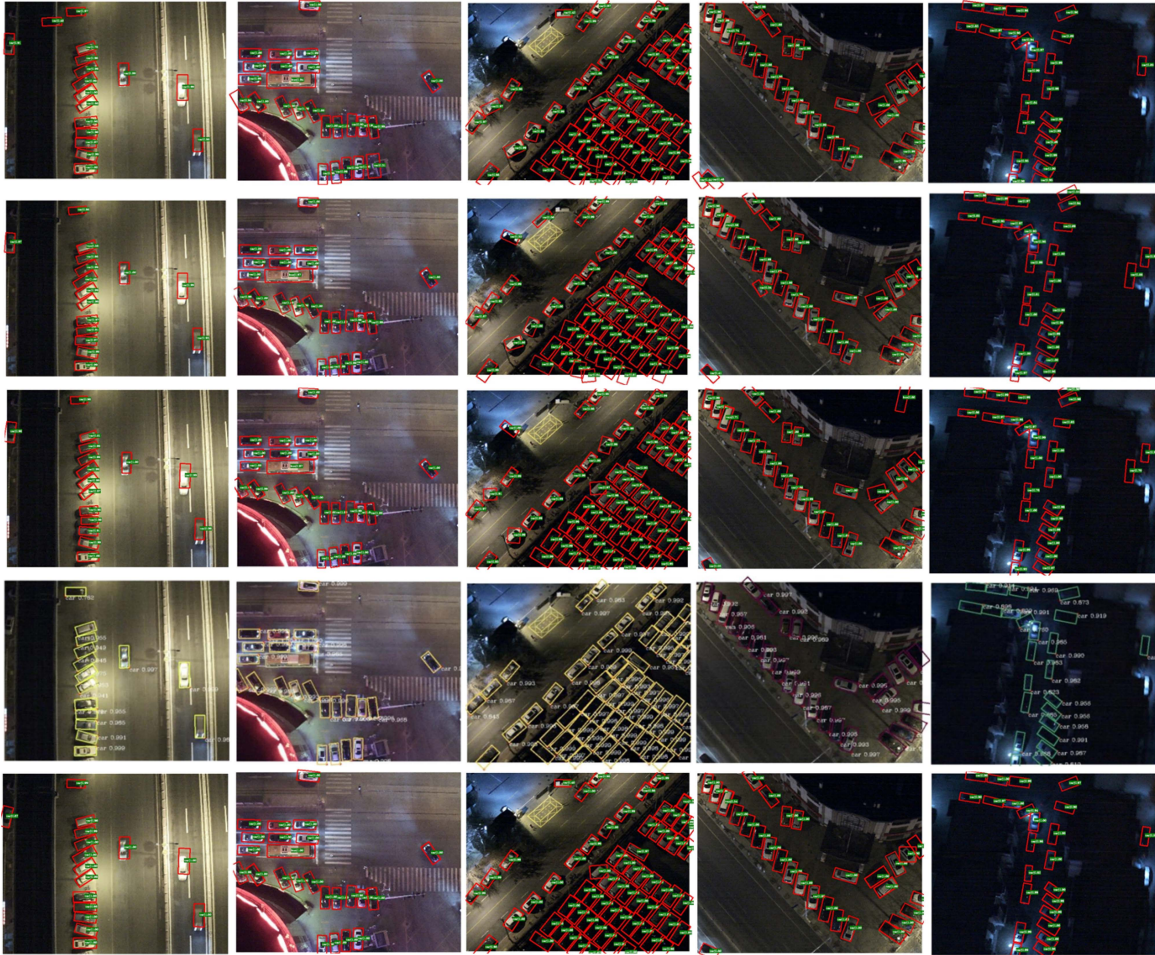


Fig. 7. Visualization of detection results of each algorithm. From top to bottom are the test results of pointwise addition method, Concat-Conv, Cross-Concat-Conv, UA-CMDet, and FFODNet.

TABLE II  
EXPERIMENTAL PERFORMANCE OF EACH OBJECT DETECTION METHOD ON DRONEVEHICLE DATASET AS WELL AS THE MODAL IMAGES IT USES

Method	Modality	mAP
Faster R-CNN [21]	R\I	43.94\52.63%
RoITransformer [24]	R\I	47.91\59.15%
ReDet [25]	R\I	51.04\60.54%
Gliding Vertex [26]	R\I	52.48\62.89%
Pointwise addition [23]	R+I	60.82%
Concat-Conv [14]	R+I	61.63%
Cross-Concat-Conv [22]	R+I	64.24%
UA-CMDet [19]	R+I	64.01%
RISNet [17]	R+I	66.40%
AR-CNN [27]	R+I	71.58%
FFODNet (ours)	R+I	<b>76.93%</b>

The optimal detection results are shown in bold.

method. Experiments show that the proposed method can extract high-quality fusion feature information from infrared-optical image pairs, and its performance is state of the art.

The subjective detection results of each algorithm are presented in Fig. 7. The first line represents the fusion method of pointwise addition, the second line represents the fusion method

TABLE III  
EXPERIMENTAL PERFORMANCE OF EACH OBJECT DETECTION METHOD ON DRONEVEHICLE DATASET AS WELL AS THE MODAL IMAGES IT USES

Method	mAP in subset (a)	mAP in subset (b)
Pointwise addition [23]	59.8%	65.5%
Concat-Conv [14]	59.4%	72.9%
Cross-Concat-Conv [22]	64.2%	77.3%
FFODNet (ours)	<b>75.2%</b>	<b>83.2%</b>

The optimal detection results are shown in bold.

of Concat-Conv, and the third line represents the fusion method of Cross-Concat-Conv. The fourth line displays the detection result of UA-CMDet, and the image is sourced from the original paper. The final line demonstrates the test results of the proposed method. The proposed method demonstrates excellent detection effectiveness, performing well on dense targets and targets with obstructed edges

The quantitative analysis of each fusion method on the two subsets is compared in Table III. To evaluate the robustness of the algorithm in difficult scenarios, the test set of the original dataset is divided into two test subsets according to the difficulty of the scenario. Experiments show that the proposed method is robust in a low-illumination environment.



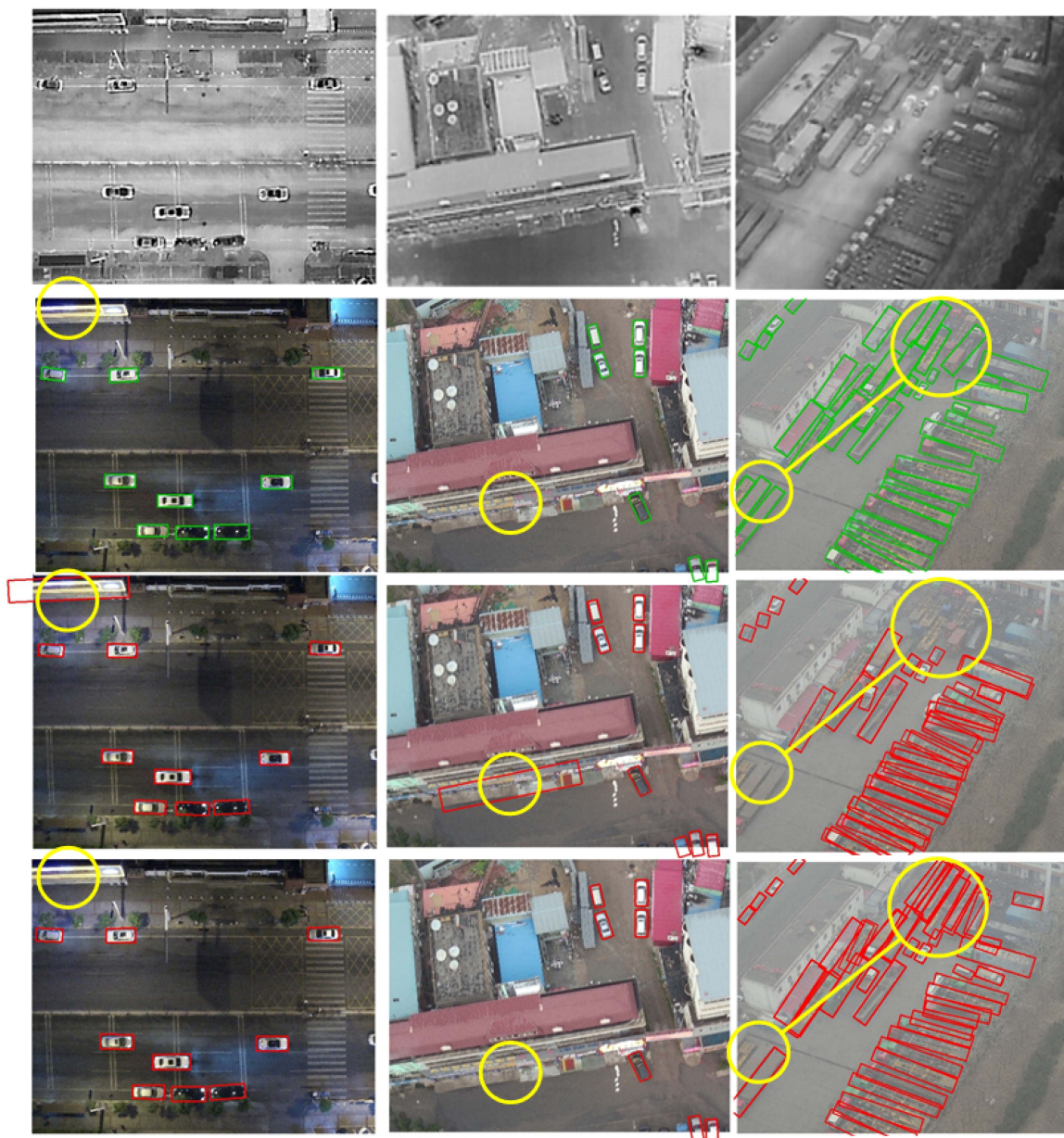


Fig. 8. Subjective diagram of test results. The rows from top to bottom are labeled optical images, baseline detection results, and the results of the FFODNet.

As shown in Fig. 8, the true value is shown in green on the first row of the RGB image. The second and third rows show the detection results for the baseline and the FFODNet, respectively. To enhance visual clarity, we highlight the objects with our approach above the baseline in yellow.

Visualizations of some of the detection results in the strong interference subset are depicted in Fig. 9. In the first column of the presented data, the infrared data exhibits a low-contrast phenomenon while the optical image suffers from cloud interference. The feature fusion network needs to address the low-illumination interference caused by optical images in the second and fourth columns of data. In the third column of data, both cloud and low-illumination problems are observed in the optical images. These phenomena indicate that a certain modality of data is not always reliable in the fusion object

detection task. Consequently, it is crucial for the algorithm to initially treat all modal data as equally important during feature fusion and adaptively suppress the interference introduced by modal data throughout the fusion process.

The experimental results unequivocally validate the effectiveness and feasibility of our algorithm, particularly when applied to the challenging strong interference subset. The algorithm demonstrates excellent detection performance, surpassing expectations and showcasing its potential for real-world applications. The exceptional performance of our algorithm can be attributed to its innovative fusion strategy, adaptive weighting mechanism, and end-to-end training approach. These key components enable the algorithm to effectively handle interference from different modalities, prioritize relevant information, and optimize the feature fusion process for accurate object detection.



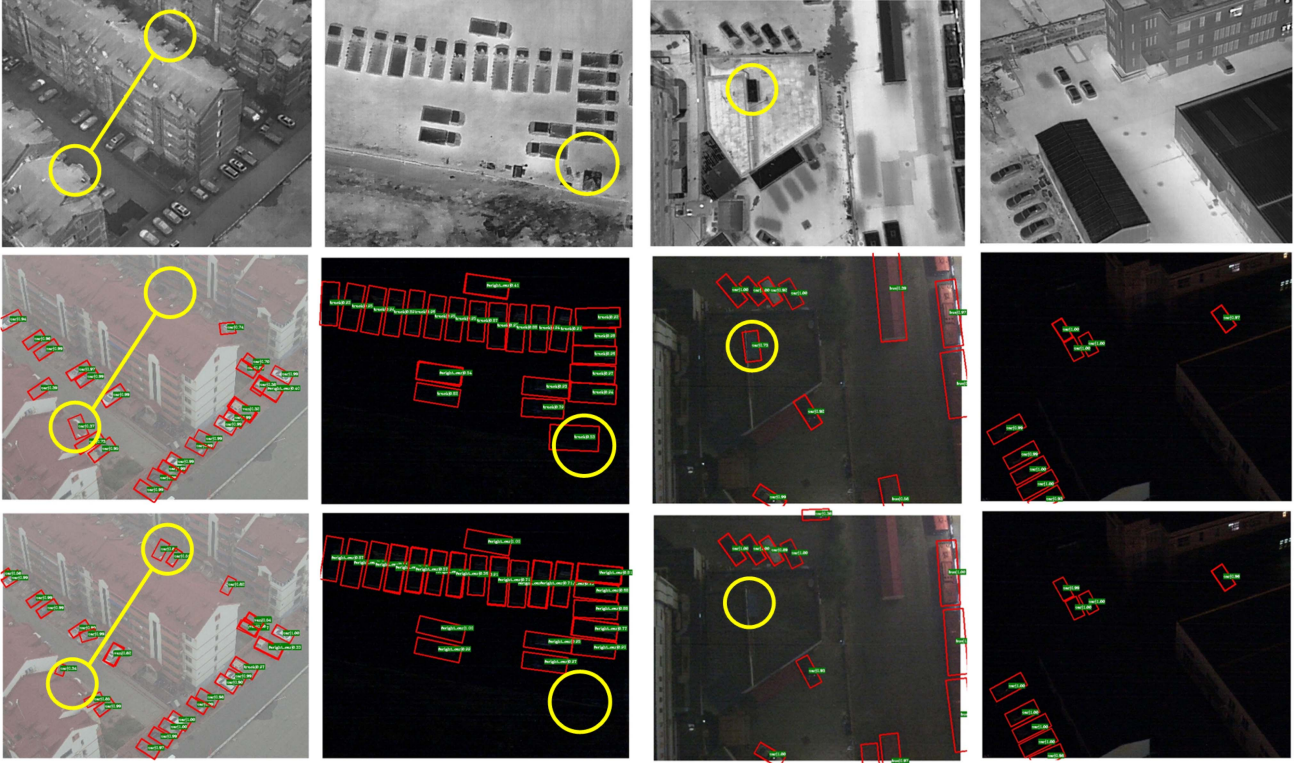


Fig. 9. Figure of data sample detection results with modal data interference. The first is the infrared image data, the second is the detection result of the baseline, and the third line is the detection result of the proposed method.

TABLE IV  
OUR PROPOSED METHOD WAS COMPARED WITH BASELINE ABLATION EXPERIMENTS

Method	car	freight-car	truck	bus	van	mAP
baseline	89.70%	35.90%	49.00%	88.30%	41.20%	60.82%
base+S	90.10%	52.60%	64.20%	88.20%	54.50%	69.92%
base+J	90.30%	62.30%	71.70%	89.40%	60.70%	74.88%
base+J+S	<b>90.40%</b>	<b>68.40%</b>	<b>72.60%</b>	89.20%	<b>64.10%</b>	<b>76.93%</b>

The optimal detection results are shown in bold.

The algorithm's ability to adaptively suppress interference and exploit the complementary characteristics of multimodal data contributes to its outstanding performance.

#### D. Ablation Experiment on DroneVehicle Dataset

To verify the effectiveness of the module proposed in this article, we conducted ablation experiments on JEOM and SIEM. Table IV shows the experimental results of the ablation experiments. The baseline in Table IV refers to the method pointwise addition in Table II.

J and S in Table IV, respectively, refer to the JEOM and SIEM in Section II of this article. The fusion module is compared with the simple fusion method. Furthermore, SIEM is used to improve the ability of the network to extract the discriminant features of the target, so as to improve the efficiency of feature fusion. In the proposed method, the detection efficiency is improved through the synergistic effect of JEOM and SIEM. Experiments show that JEOM can optimize the detection performance by

TABLE V  
EXPERIMENTAL PERFORMANCE OF EACH OBJECT DETECTION METHOD ON FLIR DATASET, AS WELL AS THE MODAL IMAGES IT USES

Method	Modality	mAP
Faster R-CNN [21]	R\I	63.60\75.30%
HalfwayFusion [28]	R+I	71.17%
DALFusion [29]	R+I	72.11%
CFR [28]	R+I	72.39%
GAFF [30]	R+I	73.80%
YOLO-MS [31]	R+I	75.20%
MFF-YOLOv5 [15]	R+I	78.20%
UA-CMDet [19]	R+I	78.60%
FFODNet(ours)	R+I	78.30%

improving the efficiency of feature fusion. On this basis, SIEM can further improve the overall network performance. On the basis of optimizing the feature fusion structure, it is meaningful to enhance the ability of each feature extraction branch.

#### E. Performance Evaluation on FLIR Dataset

To further validate the effectiveness of our proposed method, we conducted comparison experiments with other state-of-the-art methods on the FLIR-aligned dataset.

Table V presents the results of several advanced object detection methods that possess multimodal fusion capabilities. As observed from Table V, our approach outperforms the other methods, establishing itself as the leading method for object detection on the FLIR dataset. These results further demonstrate

the superior performance and effectiveness of our proposed method in multimodal object detection tasks.

As observed in Table III, using only infrared modal images can achieve a higher degree of precision, mainly because the infrared images in the FLIR dataset provide a better view compared to the optical images. Inadequate fusion methods may introduce interference information that hinders fusion, ultimately reducing the detection performance of the network. Notably, the YOLO-MS is a recently advanced multimodal fusion object detection method. However, in our experiments, we have achieved higher performance compared to the method. Furthermore, while our approach performs equally well as UA-CMDet on the FLIR dataset, it outperforms UA-CMDet on the drone dataset. These results highlight the superior performance and effectiveness of our approach in both FLIR and drone datasets.

#### IV. CONCLUSION

In this article, a novel multimodal detection FFODNet is proposed, which adaptively fuses the target feature information of multimodal remote-sensing images to achieve high-performance detection. It includes the improvement of the backbone network and fusion module. In order to obtain high-quality fusion features by enhancing object-specific features and suppressing redundant information that may hinder fusion, a new JEOM is proposed. Based on this, a new SIEM is designed to suppress irrelevant background feature information and further improve the efficiency of subsequent feature fusion operations. Experimental results show that our proposed method outperforms existing state-of-the-art methods on the DroneVehicle dataset.

#### REFERENCES

- [1] W. Lu et al., "A CNN-transformer hybrid model based on CSWin transformer for UAV image object detection," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 16, pp. 1211–1231, Jan. 2023, doi: [10.1109/JSTARS.2023.3234161](https://doi.org/10.1109/JSTARS.2023.3234161).
- [2] Y. Han, W. Meng, and W. Tang, "Capsule-inferenced object detection for remote sensing images," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 16, pp. 5260–5270, Apr. 2023, doi: [10.1109/JSTARS.2023.3266794](https://doi.org/10.1109/JSTARS.2023.3266794).
- [3] J. Xue, D. He, M. Liu, and Q. Shi, "Dual network structure with interweaved global-local feature hierarchy for transformer-based object detection in remote sensing image," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 15, pp. 6856–6866, Aug. 2022, doi: [10.1109/JSTARS.2022.3198577](https://doi.org/10.1109/JSTARS.2022.3198577).
- [4] S. Tian, L. Kang, X. Xing, J. Tian, C. Fan, and Y. Zhang, "A relation-augmented embedded graph attention network for remote sensing object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–18, May 2021, Art. no. 1000718, doi: [10.1109/TGRS.2021.3073269](https://doi.org/10.1109/TGRS.2021.3073269).
- [5] Y. Tang et al., "An object fine-grained change detection method based on frequency decoupling interaction for high-resolution remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 62, pp. 1–13, Nov. 2024, Art. no. 5600213, doi: [10.1109/TGRS.2023.3337816](https://doi.org/10.1109/TGRS.2023.3337816).
- [6] D. Zhou and X. Wang, "Robust infrared small target detection using a novel four-leaf model," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 17, pp. 1462–1469, Nov. 2023, doi: [10.1109/JSTARS.2023.3337996](https://doi.org/10.1109/JSTARS.2023.3337996).
- [7] Y. Lin, Y. Tu, Z. Dou, L. Chen, and S. Mao, "Contour Stella image and deep learning for signal recognition in the physical layer," *IEEE Trans. Cogn. Commun. Netw.*, vol. 7, no. 1, pp. 34–46, Mar. 2021.
- [8] Q. Wang et al., "Multispectral point cloud superpoint segmentation," *Sci. China Technol. Sci.*, vol. 67, pp. 1270–1281, 2024.
- [9] Y. Tu, Y. Lin, C. Hou, and S. Mao, "Complex-valued networks for automatic modulation classification," *IEEE Trans. Veh. Technol.*, vol. 69, no. 9, pp. 10085–10089, Sep. 2020.
- [10] Q. Wang, C. Yin, H. Song, T. Shen, and Y. Gu, "UTFNet: Uncertainty-guided trustworthy fusion network for RGB-thermal semantic segmentation," *IEEE Geosci. Remote Sens. Lett.*, vol. 20, pp. 1–5, Oct. 2023, Art. no. 7001205, doi: [10.1109/LGRS.2023.3322452](https://doi.org/10.1109/LGRS.2023.3322452).
- [11] C. Li, D. Song, R. Tong, and M. Tang, "Illumination-aware faster R-CNN for robust multispectral pedestrian detection," *Pattern Recognit.*, vol. 85, pp. 161–171, 2019.
- [12] D. Guan, Y. Cao, J. Yang, Y. Cao, and M. Y. Yang, "Fusion of multispectral data through illumination-aware deep neural networks for pedestrian detection," *Inf. Fusion*, vol. 50, pp. 148–157, 2019.
- [13] Z. Tu, Y. Ma, Z. Li, C. Li, J. Xu, and Y. Liu, "RGBT salient object detection: A large-scale dataset and benchmark," *IEEE Trans. Multimedia*, vol. 25, pp. 4163–4176, May 2023, doi: [10.1109/TMM.2022.3171688](https://doi.org/10.1109/TMM.2022.3171688).
- [14] J. Zhang, J. Lei, W. Xie, Z. Fang, Y. Li, and Q. Du, "SuperYOLO: Super resolution assisted object detection in multimodal remote sensing imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, 2023, Art. no. 5605415.
- [15] S. Meng and Y. Liu, "Multimodal feature fusion YOLOv5 for RGB-T object detection," in *Proc. China Automat. Congr.*, 2022, pp. 2333–2338.
- [16] C. Zhao, H. Liu, N. Su, and Y. Yan, "TFTN: A transformer-based fusion tracking framework of hyperspectral and RGB," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–15, Oct. 2022, Art. no. 5542515, doi: [10.1109/TGRS.2022.3215816](https://doi.org/10.1109/TGRS.2022.3215816).
- [17] Q. Wang, Y. Chi, T. Shen, J. Song, Z. Zhang, and Y. Zhu, "Improving RGB-infrared object detection by reducing cross-modality redundancy," *Remote Sens.*, vol. 14, no. 9, pp. 526–530, 2022.
- [18] C. Zhao, H. Liu, N. Su, C. Xu, Y. Yan, and S. Feng, "TMTNet: A transformer-based multimodality information transfer network for hyperspectral object tracking," *Remote Sens.*, vol. 15, no. 4, 2023, Art. no. 1107.
- [19] Y. Sun, B. Cao, P. Zhu, and Q. Hu, "Drone-based RGB-infrared cross-modality vehicle detection via uncertainty-aware learning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 32, no. 10, pp. 6700–6713, Oct. 2022.
- [20] H. Zhang, É. Fromont, S. Lefèvre, and B. Avignon, "Multispectral fusion for object detection with cyclic fuse-and-refine blocks," in *Proc. IEEE Int. Conf. Image Process.*, 2020, pp. 276–280. [Online]. Available: <https://api.semanticscholar.org/CorpusID:221970539>
- [21] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, Jun. 2017.
- [22] D. Hong et al., "More diverse means better: Multimodal deep learning meets remote-sensing imagery classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 5, pp. 4340–4354, May 2021.
- [23] M. Sharma et al., "YOLORs: Object detection in multimodal remote sensing imagery," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 14, pp. 1497–1508, Nov. 2021, doi: [10.1109/JSTARS.2020.3041316](https://doi.org/10.1109/JSTARS.2020.3041316).
- [24] J. Ding, N. Xue, Y. Long, G.-S. Xia, and Q. Lu, "Learning ROI transformer for oriented object detection in aerial images," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 2844–2853.
- [25] J. Han, J. Ding, N. Xue, and G. Xia, "ReDet: A rotation-equivariant detector for aerial object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 2785–2794.
- [26] Y. Xu et al., "Gliding vertex on the horizontal bounding box for multi-oriented object detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 4, pp. 1452–1459, Apr. 2021.
- [27] L. Zhang, X. Zhu, X. Chen, X. Yang, Z. Lei, and Z. Liu, "Weakly aligned cross-modal learning for multispectral pedestrian detection," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 5126–5136.
- [28] H. Zhang, E. Fromont, S. Lefevre, and B. Avignon, "Multispectral fusion for object detection with cyclic fuse-and-refine blocks," in *Proc. IEEE Int. Conf. Image Process.*, 2020, pp. 276–280.
- [29] H. Zhang, E. Fromont, S. Lefevre, and B. Avignon, "Deep active learning from multispectral data through cross-modality prediction inconsistency," in *Proc. IEEE Int. Conf. Image Process.*, 2021, pp. 449–453.
- [30] H. ZHANG, É. Fromont, S. Lefèvre, B. Avignon, and U. de Rennes, "Guided attentive feature fusion for multispectral pedestrian detection," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 2021, pp. 72–80.
- [31] Y. Xie, L. Zhang, X. Yu, and W. Xie, "YOLO-MS: Multispectral object detection via feature interaction and self-attention guided fusion," *IEEE Trans. Cogn. Develop. Syst.*, vol. 15, no. 4, pp. 2132–2143, Dec. 2023.