YOLO-Ant: A Lightweight Detector via Depthwise Separable Convolutional and Large Kernel Design for Antenna Interference Source Detection

Xiaoyu Tang, Member, IEEE, Xingming Chen, Jintao Cheng,

Jin Wu, Member, IEEE, Rui Fan, Senior Member, IEEE, Chengxi Zhang, Member, IEEE, Zebo Zhou

Abstract—In the era of 5G communication, removing interference sources that affect communication is a resource-intensive task. The rapid development of computer vision has enabled unmanned aerial vehicles to perform various high-altitude detection tasks. Because the field of object detection for antenna interference sources has not been fully explored, this industry lacks dedicated learning samples and detection models for this specific task. In this article, an antenna dataset is created to address important antenna interference source detection issues and serves as the basis for subsequent research. We introduce YOLO-Ant, a lightweight CNN and transformer hybrid detector specifically designed for antenna interference source detection. Specifically, we initially formulated a lightweight design for the network depth and width, ensuring that subsequent investigations were conducted within a lightweight framework. Then, we propose a DSLK-Block module based on depthwise separable convolution and large convolution kernels to enhance the network's feature extraction ability, effectively improving small object detection. To address challenges such as complex backgrounds and large interclass differences in antenna detection, we construct DSLKVit-Block, a powerful feature extraction module that combines DSLK-Block and transformer structures. Considering both its lightweight design and accuracy, our method not only achieves optimal performance on the antenna dataset but also yields competitive results on public datasets.

Index Terms-YOLO-Ant, Antenna interference source de-

This research was supported by the National Natural Science Foundation of China under Grant 62001173, the Project of Special Funds for the Cultivation of Guangdong College Students' Scientific and Technological Innovation ("Climbing Program" Special Funds) under Grant pdjh2022a0131 and pdjh2023b0141, the National Natural Science Foundation of China under Grant 42074038, the Fundamental Research Funds for the Central Universities and Xiaomi Young Talents Program.

Corresponding author: Xiaoyu Tang. E-mail address: tangxy@scnu.edu.cn. Xiaoyu Tang, Xingming Chen and Jintao Cheng is with the School of Electronic and Information Engineering, Faculty of Engineering, South China Normal University, Foshan, Guangdong 528225, China, and also with the School of Physics, South China Normal University, Guangzhou, Guangdong 510000, China.

Jin Wu is with the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong.

Rui Fan is with the College of Electronics & Information Engineering, Shanghai Research Institute for Intelligent Autonomous Systems, the State Key Laboratory of Intelligent Autonomous Systems, and Frontiers Science Center for Intelligent Autonomous Systems, Tongji University, Shanghai 201804, China (e-mail: rui.fan@ieee.org).

Chengxi Zhang is with the School of Internet of Things Engineering, Jiangnan University, Wuxi, 214122, China.

Zebo Zhou is with the School of Aeronautics & Astronautics University of Electronic Science and Technology of China and Aircraft swarm intelligent sensing and cooperative control Key Laboratory of Sichuan Province, Chengdu 610097, China, and also with the National Laboratory on Adaptive Optics, Chengdu 610209, China.

tection, Small object detection, Lightweight, CNN-transformer fusion.

I. INTRODUCTION

TO ensure high-quality communication in people's work and daily lives, various wireless devices operate in different frequency bands. 5G communication is of particular note due to its introduction of new frequency bands into everyday communication. However, due to the presence of numerous private wireless signals that have not undergone spectrum allocation by communication regulatory authorities, the 5G communication network has accumulated a considerable number of sources of interference. If individuals operate in the same geographical areas and occupy similar or adjacent frequency bands as these interference signals in their everyday communication, this will result in a significant deterioration in communication quality, as shown in Fig. 1. Regular remediation of radio interference sources is vital for communication departments to alleviate this situation. The identification of signal interference sources necessitates monitoring personnel to visually inspect areas where communication quality is compromised due to the presence of suspicious antennas elevated at high altitudes, constituting a time-consuming and labor-intensive task. In light of the mature advancements in unmanned aerial vehicle (UAV) cruising technology and object detection techniques within computer vision, unmanned drones have become viable alternatives for handling complex and challenging detection tasks previously performed by humans. For example, [1] [2] [3] noted that object detection tasks in deep learning combined with UAVs have been useful in production and other areas. The success of these approaches has demonstrated the feasibility of utilizing UAVs for object detection tasks aimed at interference source antennas. However, due to the nascent stage of this detection task within the current domain of object detection, the creation of a suitable antenna dataset and the exploration of appropriate object detection methodologies are of paramount importance.

Convolutional architectures are the basis for most object detection frameworks in industrial scenarios and rely on the development of efficient convolutional neural networks in deep learning. When addressing various tasks and technical challenges, corresponding enhancements to these architectures are necessary. The antenna interference source object detection task presents three main challenges. The first issue pertains to the lightweight nature and low computational complexity of a detection model; consequently, object detectors can be deployed on lightweight computing devices, enabling realtime object detection via UAVs. Previous research, exemplified by GhostNet [4] and EfficientNet [5], has focused on designing lightweight networks as potential backbones for different detection models to achieve an overall lightweight solution. However, these networks are susceptible to information feature loss. The second difficulty in the antenna interference source object detection task lies in the differences arising from the different angles and heights from which the UAV captures the antennas. These variations result in a nonuniform distribution of target sizes within the images, most of which are extremely small in size. Additionally, there is a significant interclass dissimilarity issue, wherein antennas of the same type exhibit markedly different morphologies in different images. To address these issues, researchers have explored two aspects: multiscale feature learning and attention mechanisms [6] [7] [8] [9]. However, despite improvements in small object detection accuracy, these methods encounter challenges related to the model's weak generalizability and robustness. As a result, the overall detection accuracy for the targets was compromised. Moreover, the computational complexity of such models is higher. The third difficulty is complex target backgrounds, which cause serious false and missed detections. Given that antennas are commonly installed on tall buildings or fenced balconies in practical scenarios, the resulting complex and mutually obscuring environment between the target and the background significantly hinders detection. Researchers have suggested using attention or selfattention mechanisms to address this difficulty. In [10] [11] [12], a squeeze-and-excitation(SE) attention module was proposed, or a self-attention structure was used to build the whole network for object detection. The advantage of these models lies in their ability to effectively capture the spatial relationship between the target and the background. This capability significantly enhances object detection performance on complex backgrounds. However, these mechanisms tend to consume considerable computational resources and memory, which is not consistent with the original lightweight design intention. Additionally, networks built solely on self-attention mechanisms also suffer from long training times and poor detection accuracy for small objects.

In response to the aforementioned limitations, we propose YOLO-Ant, a lightweight one-stage detector designed for detecting antenna interference sources with small targets and complex backgrounds. Initially, we analyze the scale and number of channels in each feature layer of the model; subsequently, we design the network's width and depth to ensure that the entire detection process is performed within a lightweight framework. Our design considerations aim to balance detection accuracy with the reduction of model parameters and computational complexity. To address the issues of small target size and large interclass variation, we implement an efficient feature extraction module based on depthwise separable convolution, DSLK-Block, which is applied to each feature layer in the model. This method effectively enhances



Fig. 1: The process of 5G communication in the CBN-U-H5H-0713 area is shown in the figure. Two antenna interference source signals appear in it. The gNB (gNodeB) denotes a 5G base station. The UE (User Equipment) denotes the terminal equipment that users use to access the wireless network.

the network's feature learning and fusion capabilities, leading to a significant improvement in detection accuracy for all types of targets, particularly small targets. Additionally, this approach contributes to reducing the model's overall weight. Finally, to address the problem of complex backgrounds, YOLO-Ant uses an innovative CNN and transformer hybrid structure to act on the neck of the model. This process enables us to fully utilize both local and global feature learning to address the challenges posed by complex backgrounds while still accounting for small object detection. This approach significantly improves all the detection accuracy indicators while only slightly increasing the model's number of parameters and computational complexity. To demonstrate the model's generalizability and robustness, we also tested YOLO-Ant on public datasets and achieved highly competitive results. In conclusion, the main contributions of this paper are as follows:

- In response to the lack of learning samples for antenna object detection schemes, we conducted image acquisition and manual annotation of the three most common types of antennas encountered in real-world interference source investigation tasks. This dataset is pioneering and establishes the foundation for subsequent work.
- 2) We initially pruned YOLOv5-s [13], obtaining a lightweight detection framework. Within this framework, (i) a lightweight plug-and-play module based on depthwise separable convolution combined with large convolutional kernels was proposed to effectively improve the feature extraction and detection capabilities of the network for small targets; (ii) the innovative use of a transformer module to construct the neck structure of the detection model improved the detection capability of the network without increasing the model's parameter count or computational complexity, effectively solving the problem of dealing with complex backgrounds.
- 3) Our proposed method achieves state-of-the-art (SOTA)

performance on the antenna dataset, striking a balance between lightweight design and detection accuracy. Moreover, YOLO-Ant yields competitive results on public datasets such as COCO, validating its robustness and superior performance. Source code is released in https://github.com/SCNU-RISLAB/YOLO-Ant.

The remainder of this paper is structured as follows. Section II presents related work, briefly introducing the improvement points of the model proposed in this paper and the related work involved, including CNN network development, the emergence of the transformer detection model and the crossover development between the CNN and transformer. Section III discusses the proposed lightweight detection framework based on the YOLOv5-s improvement, including pruning the baseline model, the design of the DSLK-Block module, and the neck structure built based on DSLKVit-Block. The experimental results are given in Section IV, and Section V concludes the paper.

II. RELATED WORK

A. CNN (convolutional neural network)-based object detection

The development of object detection in the computer vision field has been greatly influenced by CNN-based methods. Traditional approaches using hand-designed features and classifiers have been shown to be inadequate, leading to the dominance of CNN-based methods. The initial CNN model, LeNet-5 [14], was limited by computational resources and model size. However, with advancements in computational power and larger datasets, deeper and more complex CNN models, such as AlexNet [15], VGGNet [16], GoogLeNet [17], and ResNet [18], have emerged. These models have improved network accuracy, reduced parameters, and addressed network degradation issues, laying a solid foundation for 2D object detection.

Two distinct methods have emerged from the convoluted development of 2D object detection: two-stage and one-stage detectors. Two-stage detectors, such as R-CNN [19] and Fast R-CNN [20], generate candidate frames using algorithms and perform classification and regression on each candidate frame. Faster R-CNN [21] introduces the region proposal network (RPN) for candidate frame generation. In contrast, one-stage detectors, such as YOLO [22] and SSD [23], perform classification and regression directly on each location in the input image. YOLOv2 [24] and YOLOv3 [25] improved detection accuracy through methods such as multiscale prediction, batch normalization, and feature pyramid networks (FPNs). SSD introduces multiscale detection using multiple-scale feature maps, while RetinaNet [26] focuses on addressing the category imbalance problem. For the aforementioned model, one-stage detectors are more suitable for real-time detection tasks on UAVs than two-stage detectors are because they do not require additional networks or algorithms for fine-tuning. However, to compensate for the deficiency in accuracy resulting from the pursuit of detection speed, improvements need to be made to the backbone and neck of the one-stage detector by developing various efficient feature extraction modules or structures. The backbone and neck are the basic components of object detection models. The backbone is a CNN trained on image classification datasets such as ImageNet [27], in which the input image is transformed into a high-dimensional feature representation. The neck module further processes the feature map, changing the scale and resolution to extract different levels of feature information. Numerous object detection models, such as NAS-FPN [28], EfficientDet [29], YOLOv4 [30], and YOLOv7 [31], have been developed based on these concepts, incorporating various improvements and techniques to enhance accuracy and performance. However, these general models are often designed with modules that consider various common tasks, exhibiting generalizability but not effectively addressing specific challenges in particular scenarios. For instance, there are several challenges, such as small object detection and complex backgrounds, in our task. Therefore, making taskspecific modifications is crucial when contemplating different tasks.

B. Developing an attention mechanism in the CV domain

Attention mechanisms, initially utilized in natural language processing, have gained significant traction in computer vision [32], [33], particularly in the field of object detection. Attention mechanisms such as channel attention, spatial attention, and their combinations have been introduced [34] [35] [36]. They effectively utilize global and local information in feature maps, improving feature representation and attention weighting, thereby enhancing model accuracy and efficiency. However, for these conventional attention mechanisms, a fixed window size or other constraints are typically employed to regulate the correlation between each position and others. In contrast, self-attention mechanisms can extract information from different positions in the information sequence more flexibly, enabling the extraction of global information. This flexibility has contributed to the widespread application of transformer [37] models based on self-attention mechanisms, including in various domains such as computer vision. For example, the Vision Transformer (ViT) [38] splits images into patches for self-attention computations. The swin transformer [39] improves local information processing by using a window-based partitioning approach. Detection with transformers (DETR) [11] adopts a global self-attention mechanism, allowing each position to obtain contextual information from the entire image. Naturally, transformers incur substantial computational costs and training time, posing challenges for model convergence. To address these challenges, researchers have introduced lightweight transformer object detectors, including MobileViT [40] and EdgeViT [41]. Moreover, innovative approaches such as conditional DETR [42] and DN-DETR [43] have been developed to address the crucial issue of slow training convergence. However, due to their simplified design, the majority of current lightweight transformer structures are applicable only to classification tasks involving small-sized image inputs and are not suitable for detection tasks. The proposed detection methods aimed at addressing slow convergence have made transformer models more complex. Therefore, achieving a balance between the lightweight

nature of transformer models and detection accuracy remains a crucial research scope in the current field of computer vision.

C. Combination of CNN and Transformer

In object detection, CNNs and transformers have distinct applications and advantages. CNNs are known for their strong image feature extraction abilities, ability to perform multichannel processing, and ability to learn spatial correlations. However, CNN-based models have limitations in handling objects of different sizes and proportions due to fixed window sizes and strides. On the other hand, transformers exhibit excellent performance in capturing long-range dependencies within input sequences without prior knowledge, albeit at a slower speed and requiring substantial amounts of training data. Evidently, the amalgamation of CNNs and transformers offers complementarity across various dimensions, and researchers have already delved into numerous methodologies to explore this synergy.

The pioneering DETR model replaces fully connected and convolutional layers with transformers while using ResNet as the feature extractor, improving accuracy and efficiency. Huawei's CMTBlock combines depthwise separable convolution and the transformer's multihead self-attention module for local and global information fusion. The CMT model [44] stacks the CMTBlock in a hybrid CNN-transformer structure. The Conformer [45] adopts a dual-network structure, where the CNN branch enhances local perception of the transformer branch. The mobile-former [46] features parallel CNN and transformer modules with bidirectional bridges, leveraging MobileNet [47] for local processing and the transformer for global interaction. However, networks or models employing such hybrid structures face challenges in effectively balancing accuracy and lightweight design. For instance, detectors such as DETR, lacking FPN structures, exhibit suboptimal performance in small object detection. While the CMT and Conformer networks have proven effective in classification tasks, their application to downstream tasks such as object detection deviates from the realm of lightweight design. In contrast to the aforementioned models, which concatenate both structures, an alternative approach involves making transformerstyle improvements directly on the CNN network. ConvNeXt [48] implements novel architectures and optimization strategies similar to those of transformers, achieving competitive results without attention structures. RepLKNet [49] employs large convolutional kernels to widen the receptive field, thus emulating the transformer-like capability for global feature extraction. By investigating the computational principles of transformers, ACMix [50] maps their operation process onto convolutional operators, thereby combining them with traditional convolution operations to construct a novel CNN architecture. Parc-Net [51] introduces circular convolution for global information extraction within a pure convolutional structure. Although these innovative networks may not achieve SOTA performance, their greater significance lies in exploring the factors contributing to the success of transformers from a CNN perspective, providing inspiration for subsequent research endeavors. The fusion of transformers and CNNs offers a flexible and diverse range of integration methods. Future research should strive to deepen the understanding of their interactions to improve design and optimization.

D. Object Detection of Antenna Interference Sources

Regularly monitoring and mitigating antenna interference sources has become one of the most critical tasks in the wireless communication field. In the past, detecting antenna interference sources mainly relied on traditional techniques such as spectrum analysis, signal recognition and positioning. However, these methods have many limitations. For example, when detection personnel identify a radio interference signal through a spectrum analyzer, they can determine only the approximate direction of the interference source based on the strength of the received signal and cannot accurately determine its position.

The rapid advancement of deep learning and computer vision has facilitated the successful application of object detection-assisted tasks in various industries. Examples include defect detection in industrial settings, pest/weed detection in agriculture, and vehicle and pedestrian detection in transportation [52] [53] [54] [55] [56]. These solutions provide effective ideas for our antenna interference source detection task. When investigators confirm the approximate direction of the interference source antenna through a signal receiver and spectrum analyzer, they can use drones with cameras and related object detection algorithms to replace manual accurate positioning work. Unfortunately, the field of antenna interference source detection based on object detection tasks has largely not been explored. Due to the lack of learning samples and models for related antenna interference source detection, existing detection methods are not suitable for antenna detection. Therefore, it is urgent and meaningful to create a professional dataset and train a model suitable for this detection task to address the difficulty of locating antenna interference sources in the wireless communication field.

III. PROPOSED DETECTION FRAMEWORK

A. Overall model structure

The overall idea for the network(Fig. 2) lies in the combination of a CNN and transformer, both the inductive bias ability of the convolutional operation and the ability of the transformer to extract global information, while also meeting the needs of a lightweight model with low computational complexity. YOLO-Ant adopts DSLKNet, which is composed of DSLK-Blocks, as the backbone for downsampling and feature extraction in images. In DSLKNet, four DSLK-Layers employ convolutional kernels of varying sizes to sequentially extract rich features from different receptive fields of the image. To address the challenge of detecting small objects, we incorporate the neck structures of the FPN and PAN for multiscale feature learning. On the neck component, we conducted pruning based on YOLOv5-s (detailed data provided in Section IV. EXPERIMENT). In comparison to the baseline model, the pruned neck model features an increased number of module stacks and a reduced number of channels in each module. This structural modification effectively alleviates redundancy



Fig. 2: The YOLO-Ant model can be roughly divided into three parts: the backbone consisting of DSLKNet, the FPN+PAN structure consisting of DSLK-Layer and DSLKVit-Block forming the neck.

in feature extraction, resulting in an overall model that is not only more lightweight but also attains higher detection accuracy. In this lightweight framework, we introduce the DSLKVit-Block, which is a combination of transformer and convolutional modules. Even though the transformer module has a larger parameter count and computational complexity than does the convolutional module, the final YOLO-Ant model remains lighter than the baseline model. Overall, the integration of CNN and transformer in the model is manifested as follows: (i) a transformer-like pure CNN structure is proposed using depthwise separable convolution and a large convolution kernel, effectively expanding the perceptual field of the convolution operation and extracting the target context information; (ii) the FPN structure of the pure CNN and the PAN structure designed with the transformer module are combined in a parallel structure to complement each other and thus improve the feature processing capability of the model. The following two subsections provide a detailed description of the working principles of the DSLK-Block and DSLKVit-Block within the YOLO-Ant.

B. More efficient feature extraction module, DSLK-Block

The DSLK-Block structure built with depthwise separable convolution was introduced using large convolutional kernels. The design of this structure is based on several starting points. (1) Depthwise separable convolution is used instead of conventional convolution operations to achieve model lightweighting. (2) Large convolutional kernels are used to increase the receptive fields to extract a greater amount of contextual information. Models such as RepLKNet have shown that pure convolutional networks can achieve performance comparable to that of transformer-style networks in this way. (3) Inspired by the ConvNeXt approach, DSLKBlock uses fewer normalization and activation functions, replacing the rectified linear unit (ReLU) with a Gaussian error linear unit (GELU).

However, unlike other feature extraction blocks, DSLK-Block has three places that change correspondingly with the network location of the DSLK-Layer to balance the relationships between parameter volume, computational complexity, and accuracy. (1) The size of the large convolutional kernel in the backbone changes according to the location of the DSLKBlock. The rationale behind this design primarily stems from several considerations. First, if all DSLK-Layers adopt excessively large convolutional kernels, the model will be greatly burdened in terms of parameters and detection speed. Second, small objects typically have pixels in the range of 32×32 , roughly equivalent to 1/40 of the original image size. When the network input size is set to 640×640 , the corresponding size of the small objects is approximately 16×16 . To ensure that the early layers of the network can sufficiently extract feature information from small objects and prevent the loss of small object details caused by large convolutions, we employ 5×5 and 9×9 convolutional kernels in the first two layers. Finally, as the downsampling rate increases, the feature map sizes decrease. We employ larger convolutional kernels to handle larger-sized objects, while these larger kernels also provide more comprehensive contextual information for small feature maps (derived from the relationships between adequately extracted small objects and their surrounding environments). In the final model backbone, the sizes of the large convolutional kernels are set to 5×5 , 9×9 , 13×13 , and 27×27 , thereby achieving efficient detection, particularly for objects of different sizes, especially small objects. (2) For the DSLK-Layer within the model backbone, to prevent the potential loss of important information as the convolutional kernel size increases, we incorporated a parallel pathway using 3×3 convolutional



Fig. 3: DSLK-Block & DSLK-Layer

kernels within the DSLK-Block. The output from the small kernel pathway is fused with that of the large kernel pathway using an addition operation. To ensure model lightweightness and expedite model convergence, within the neck structure, all the large convolutional kernels within the DSLK-Block were resized to 3×3 dimensions, while the small convolutional kernel pathway was modified to follow a conventional shortcut form. (3) Depthwise separable convolution first uses depthwise convolution to convolve each feature point within the same channel and then extracts information between different channels of the same feature point through pointwise convolution. After decomposing conventional convolution operations into these two steps, the computational cost of the convolution operation is effectively reduced. The pointwise convolution of the DSLK-Block adopts a variable factor to control the number of channels. To further balance the relationships between the model parameters, computational complexity, and accuracy, the DSLK-Block uses different pointwise convolution depths at different positions in the network and achieves good results. The final structures of DSLK-Block and DSLK-Layer are shown in Fig. 3. The DSLK-Block is part of the DSLK-Layer, as illustrated on the right-hand side of Fig. 3, and forms each CNN feature processing layer within the backbone and neck of the model. Assuming that the input feature map is denoted as $X_i \in R^{C_i \times H_i \times W_i}$ (where C, H and W denote the number of channels, spatial height and width, respectively), its workflow can be represented as follows:

$$F_{DSLK-Layer}(X_i) = f_{CBS}(1, C_i, f_{CBS}(1, C_i/2, X_i))$$

$$\otimes F_{DSLK-Block}(f_{CBS}(1, C_i/2, X_i))$$
(1)

where \otimes represents the Concat operation and f_{CBS} represents a module that sequentially undergoes convolution, normalization, and activation functions, which can be expressed using

the following formula:

$$f_{CBS}(k, c, X_i) = \rho(\sigma(f_{conv}(k, c, 1, X_i)))$$
(2)

where $\rho(x)$ represents batch normalization and $\sigma(x)$ represents the SiLU activation function. where $f_{conv}(k, c, g, X_i)$ represents the convolution operation, k is the kernel size, c is the output channel number, and g represents the number of groups (g=1 in regular convolution, and $g = C_i$ in depthwise separable convolution).

We represent the workflow of the DSLK-Block using formula $F_{DSLK-Block}$:

$$F_{DSLK-Block}(X_i) = X_i + f_{pw}(e, f_{dw}(K_L, X_i) + f_{dw}(K_S, X_i))$$
(3)

where K_L and K_S represent the large and small convolution kernels used in the two depthwise convolution paths respectively (K_L takes values of 3, 5, 9, 13, and 27 in the model. K_S takes a value of 0 or 3. If K_L =0, then the path becomes a normal shortcut operation). where f_{dw} represents the substitution of conventional convolutions with depthwise convolution in f_{CBS} , as expressed by the following formula:

$$f_{dw}(K, X_i) = \rho\left(\sigma\left(f_{conv}(K, C_i, C_i, X_i)\right)\right) \tag{4}$$

while f_{pw} represents the pointwise convolution block, as expressed by the following formula:

$$f_{pw}(e, X_i) = \rho(f_{conv}(1, 1, C_i, \sigma(f_{conv}(1, 1, e \times C_i, X_i))))$$

$$(5)$$

where e represents the variable expansion coefficient, which is used to control the channel expansion and scaling factor in the pointwise convolution process.

C. DSLKVit-Neck structure for efficient feature fusion

The efficiency of transformer models relies heavily on their global attention mechanism, which differs from convolution



Fig. 4: DSLKVit-Block

operations in that information between feature points is calculated only within the size of the convolution kernel. Instead, transformer models consider the interactions between each feature point and all other points in the feature map. Using the transformer's self-attention mechanism to enhance the contextual information extraction of feature points is useful for dealing with problems such as target shape differences caused by multiangle shots from UAVs and interference from complex environmental backgrounds in antenna detection tasks. However, this design comes at the cost of consuming a significant amount of computational resources, making it unacceptable for lightweight models. Therefore, finding methods to efficiently utilize transformer structures while conserving resources remains a significant challenge.

In MobileViT, the transformer calculates information only between feature points at the same position in each patch of the feature map. In EdgeViT, a convolution operation is used to aggregate all the information within a local window of size k \times k. Then, the transformer calculates the information between all the feature points that contain the aggregated information. Due to the large amount of data redundancy inherent in image data, the difference in information between adjacent pixels is often not significant, allowing this computational savings to occur. Motivated by models such as MobileViT and EdgeViT, we introduce an innovative DSLKVit-Block module, which incorporates both CNN and transformer architectures, as shown in Fig. 4. To reduce the computational complexity of the entire model, DSLKVit-Block calculates information only between "representative feature point" within small regions of the feature map instead of computing the mutual information between every feature point.

The DSLKVit-Block initially conducts feature extraction on localized regions of the feature map through convolution operations within the DSLK-Block, resulting in a set of "representative feature points". Each of these feature points represents the aggregated features from their respective regions. The new feature map, composed of all these "representative feature points", has reduced dimensions, conserving resources for subsequent self-attention computations. Assuming that the input feature map is denoted as $X_i \in R^{C_i \times H_i \times W_i}$ (where C, H and W denote the number of channels, spatial height and width, respectively), we denote the process of aggregating local information in the DSLKViT-Block as follows using the formula f_{Local} :

$$f_{Local}(X_i) = f_{LA}(sr, \rho(f_{PE}(X_i) + F_{DSLK-Block}(X_i)))$$
(6)

where $\rho(x)$ represents the NormLayer operation, f_{PE} represents the positional encoding achieved through convolution operations and $f_{LA}(sr, x)$ represents the operation of aggregating local information, which can be achieved using pooling layers. The sr parameter indicates the local range covered by the "representative feature points" and is equal to $sr \times sr$.

Subsequently, the feature map undergoes multihead attention calculations, enabling the interactions among the "representative feature points" that aggregate information within each region, thereby acquiring rich contextual information between various regions on the original feature map. Finally, the model employs deconvolution operations to map the "representative feature points" back to their respective corresponding regions and enhances the network's expressive capabilities through a feed-forward network (FFN). The process of obtaining global information can be expressed using the following formula:

$$f_{Global}(sr, X_i) = X_i + \rho(f_{LD}(sr, MHSA(f_{Local}(X_i)))$$
(7)

where $f_{LD}(sr, x)$ represents the operation of mapping representative feature points back to their original regions, which can be accomplished via deconvolution. *MHSA* represents the multihead self-attention computation:

$$MHSA(X) = MultiHead(Q, K, V)$$
$$= Concat(head_1, head_2, ..., head_h)W^O$$
(8)

where $head_i = Attention(QW_i^Q, KW_i^K, VW_i^V), W_i^Q \in R^{d_{model} \times d_k}, W_i^K \in R^{d_{model} \times d_k}, W_i^V \in R^{d_{model} \times d_v}$, and $W^O \in R^{hd_v \times d_{model}}$. In this work we employ h = 4 parallel attention heads. For each of these we use $d_k = d_v = d_{model}/4 = C$. For single headed attention, we compute the attention function on a set of queries simultaneously, packed together into a marix Q. The keys and values are also packed together into matrices K and V. We compute the matrix of outputs as follows:

Attention
$$(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$
 (9)

where d_k is a scaling factor. We can obtain $Q \in R^{N \times C}, K \in R^{N \times C}$, and $V \in R^{N \times C}$ through a linear layer by input feature $X \in R^{C \times H \times W}$ respectively, where $N = H \times W$.

Ultimately, the entire workflow of DSLKVit-Block can be represented as follows:

$$F_{DSLKVit}(sr, X_i) = FFN(f_{Global}(sr, (f_{Local}(X_i))))$$
(10)

$$FFN(X_i) = X_i + \rho(MLP(X_i)) \tag{11}$$

where the MLP function can be implemented using two fully connected layers.

After finalizing the foundational DSLKVit-Block module, we proceeded to modify the model's neck structure at the macro level, aiming to harness the complementary advantages of transformers and CNNs to their fullest extent. We observed that the neck structure of YOLOv5-s adopts an FPN with a PAN architecture. In this configuration, the FPN's topdown structure merges high-level feature maps with low-level feature maps to propagate semantic information from higher levels, enhancing the representational capacity of the lowerlevel feature maps. Subsequently, the connection of the PAN, a bottom-up structure, aims to assist high-level feature maps in acquiring richer positional information from shallow-level feature maps. The PAN and FPN also engage in lateral connections to fuse feature maps of the same dimensions, enriching the information content. Therefore, the use of the transformer to construct PAN is logically sound. During the bottom-up process, the DSLKVit-Block gradually extracts finer features from large-sized feature maps to small-sized feature maps. Moreover, lateral connections are leveraged to incorporate the convolution results of the DSLK-Block to compensate for the insensitivity of the transformer structure to positional information. Furthermore, due to the computational complexity of the transformer being $O(N^2)$, applying it to feature maps of size 80×80 would results in a prohibitively high computational cost. Therefore, we deployed the DSLKVit-Block only in the smallest two feature layers (40×40 and 20×20).

The DSLKVit-Neck structure constructed from the above design incorporates the complementary concepts of CNNs and transformers at both the microlevel (within DSLKVit-Block) and macrolevel (the entire model's neck structure).

IV. EXPERIMENT

A. Experimental Description

Dataset:

1) Antenna interference source dataset: In the antenna interference source detection task proposed in this article, the methods and datasets in the relevant field are still lacking. Therefore, our first priority is to create the first dataset specifically for antenna interference source detection. In the daily interference source investigation and cleaning activities of the communication bureau, we identified three common antenna interference sources that significantly impact communication: the Yagi antenna, plate log antenna, and patch antenna, as shown in Fig. 5 a-c. The three antennas have different shapes, bringing multiple challenges to the detection task. For example, the Yagi antenna is easily confused with a complex fence background when it is installed on a balcony, as shown in Fig. 5d. The plate log antenna has a flat and wide shape that can be recognized from the front well. However, the angle from the side greatly reduces the recognition rate, and different angles create problems



Fig. 5: Antenna Interference Source Dataset

with significant interclass differences, as shown in Fig. 5 e. The patch antenna is the most difficult to recognize because it is merely a monochromatic rectangular object accompanied by a lengthy wire, resulting in a lack of intricate features available for the model to learn from. Additionally, when the patch antenna is wrapped around a pole and placed high up, the object captured by aerial photography is small and difficult to distinguish in the image, as shown in Fig. 5f. Moreover, we simulate different angles, distances, and lighting effects to enrich the dataset when drones take pictures of the antennas. Based on the collected images of the three types of antenna interference sources, we used common data augmentation techniques such as rotation, flipping, and image color space changes to expand the original images. Finally, the dataset was divided into training images (1777) and validation images (449) at an 8:2 ratio and labeled by professionals.

- 2) COCO: The Common Objects in Context (COCO) dataset is widely used in the computer vision community for benchmarking object detection and segmentation algorithms. It is designed for use in object detection, segmentation, and captioning. The dataset contains more than 330,000 images with more than 2.5 million object instances labeled across 80 different object categories. Additionally, the dataset includes 91 different categories of "stuff" or background categories, such as sky, grass, and water. The dataset also includes annotations for object segmentation masks, object bounding boxes, and keypoint annotations.
- VisDrone: The VisDrone dataset is a large-scale benchmark dataset for visually understanding aerial scenes. The dataset consists of more than 6,000 video clips

and more than 2 million images captured by various types of drones in different locations and scenarios, covering a wide range of aspects of aerial vision. The dataset contains rich annotations for object detection, tracking, counting, and image quality assessment, which are important for both academic research and industrial applications.

Experimental Setting: Experimental Configuration: All the experiments were executed on the Linux operating system, an NVIDIA GTX3060 GPU with 12 GB of memory, and data training and object detection were performed in PyTorch 1.12.0, torchvision0.13.0, and CUDA10.2 environments. On two public datasets, COCOtrain2017 and VisDrone2019-DETtrain, we conducted training and testing on the COCOval2017 and VisDrone2019-DET-val validation sets, respectively. In the experiments with the COCO dataset, YOLO-Ant was trained without the use of pretrained weights, and the model was tested using the weights that achieved the highest mAP.5:.95 after training to full convergence. For the Antenna dataset and the Visdrone dataset, training was more challenging due to the smaller number of samples and the difficulty posed by small targets. To improve the model fit and enhance the accuracy, all the models were trained using pretrained weights (models pretrained on the COCO dataset). To ensure a fair comparison, all the experiments involving model training, testing, and other evaluations are conducted using images of size 640x640 as the input. All pretrained weights for comparative models in the experiments, as well as settings for hyperparameters such as learning rates and optimizers, were sourced from the official project or MMDetection's official documentation.

Evaluation Criteria: In the object detection field, the following metrics are used to judge the detection accuracy of a model [57]: accuracy (P), recall (R), average precision (AP), mean AP (mAP), and intersection over union (IoU). The accuracy is the proportion of positive cases predicted by the model that are actually positive. Recall is the proportion of all true positive cases that the model correctly predicts as positive cases. AP is the weighted average of accuracy and recall. The mAP is the mean AP value for different categories. The IoU is the ratio of the intersection area of the predicted frame to the true frame to the concatenated area and is usually used to evaluate the localization accuracy of an object detection model. mAP.5 indicates the mAP when the IoU threshold is 0.5; mAP.5:.95 evaluates the performance of the algorithm over a range of IoU thresholds from 0.5 to 0.95, as the detection model performance under different IoU thresholds is considered, allowing a more comprehensive assessment of the detection and localization capabilities of the model.

$$P = \frac{T_P}{T_P + F_P}$$
(12)

$$R = \frac{T_P}{T_P + F_N}$$
(13)

$$AP = \int_0^1 P(R)dR \tag{14}$$

$$mAP = \frac{\sum_{i=1}^{N} AP(i)}{N}$$
(15)

In the formula, TP (true positive) is the correct detection result, FP (false positive) is the wrong detection result, and FN (false negative) is the wrong undetected result, which means missed detection. N indicates the number of detection task categories.

In addition to model accuracy, the number of parameters and giga floating-point operations per second (GFLOPs) are crucial considerations for determining the lightweight nature of a model. FLOPs indicate the number of floating-point operations performed by a model in a single forward propagation, serving as a measure of its computational complexity. This metric enables the comparison of computational overhead across different models. GFLOPS represents the billion floatingpoint operations per second, with 1 GFLOPs equaling 1,000 MFLOPS. In general, the more complex the model structure is, the higher the number of FLOPs. When designing an object detection model, a trade-off between accuracy and computational complexity is needed to achieve better detection performance and higher computational efficiency by making a corresponding choice according to the specific task requirements.

B. Experimental Results and Analysis

Baseline Model Selection: Starting from the practical task of antenna interference source detection, we select classic or cutting-edge algorithms (with a focus on lightweight methods) for comparison on our antenna dataset. The selection process encompasses considerations from three main aspects: model parameter count, computational complexity, and detection accuracy. It is evident from Table I that with the increased popularity of transformers in the computer vision domain in recent years, most detection methods have become closely associated with them. However, this trend has led to increasingly larger models, rendering them unsuitable for lightweight real-world applications. In terms of detection accuracy, several cuttingedge transformer models such as PVT-tiny [61] and DAB-DETR, not only fail to meet the lightweight criteria but also exhibit detection performances inferior to those of traditional single-stage detectors, such as FCOS [58] and YOLO. The reason behind this phenomenon lies in the inherent limitations of transformer mechanisms for small object detection.

Taking a holistic view, among this group of the lightest models (enclosed by the red box in Fig. 7(c)), there is a significant disparity in model accuracy, with most mAP.5 values falling below 0.4. In contrast, among some larger models (enclosed by the green box in Fig. 7(c)), the accuracy on the antenna dataset is higher and more consistent, with mAP.5 values generally hovering at approximately 0.5. We observed that lightweight models such as YOLOv4-tiny and EfficientDet-D1 are constrained by their simplistic network structures and computational limitations, making it challenging for them to achieve outstanding performance on antenna datasets characterized by small targets and complex backgrounds. On the other hand, larger models have better feature processing capabilities, resulting in overall superior accuracy. However, to meet the demand for deploying models on lowlevel computing platforms, lightweight models remain a crucial focus. From this perspective, even though RT-DETR [66]

model	Year	Param.	GFLOPs	mAP.5	mAP.5:.95	Yagi Antenna mAP.5	Plate Log Antenna mAP.5	Patch Antenna mAP.5
RetinaNet(ResNet50) [26]	ICCV'2017	37.97M	191.43	0.477	0.237	-	-	-
YOLOv3-tiny [25]	'2018	8.85M	13.17	0.32	0.156	0.288	0.645	0.0274
FCOS(ResNet50) [58]	ICCV'2019	32.12M	60.59	0.533	0.264	-	-	-
YOLOv4-tiny [30]	'2020	5.88M	16.18	0.28	0.134	-	-	-
YOLOv5-s [13]	'2020	7.23M	16.49	0.619	0.339	-	-	-
YOLOx-s [59]	CVPR'2021	8.94M	26.76	0.579	0.350	0.650	0.781	0.306
YOLOv7-tiny [31]	CVPR'2023	6.23M	13.86	0.42	0.217	0.476	0.642	0.141
EfficientDet-D1 [29]	CVPR'2020	6.56M	11.51	0.267	0.156	-	-	-
TOOD(ResNet50) [60]	ICCV'2021	32.02M	60.79	0.499	0.265			
PVT-tiny [61]	ICCV'2021	23.0M	50.94	0.472	0.246	-	-	-
DAB-DETR(ResNet50) [62]	ICLR'2022	43.70M	33.90	0.531	0.227	-	-	-
YOLOv8-s [63]	'2023	11.2M	28.4	0.548	0.316	-	-	-
DINO(ResNet50) [64]	ICLR'2023	47.54M	93.77	0.595	0.331	-	-	-
Mask R-CNN(ConvNeXt-V2-B [65])	'2023	110.59M	198.56	0.642	0.348	-	-	-
RT-DETR(ResNet18) [66]	'2023	20M	60	0.653	0.362	-	-	-
ours	'2023	6.13M	16.18	0.692	0.374	0.739	0.825	0.512

TABLE I: The performance of different models on the Antenna dataset

TABLE II: Baseline model pruning

		YOLOv5-s		YOLOv5s -lite		
Layer Name	Output size	modules[in channels,out channels]×n	params	modules[in channels,out channels]×n	params	
stage 17	80×80	C3[256,128]×1	90880	C3[256,128]×3	173312	
stage 18	80×80	Conv[128,128]×1	147712	Conv[128,128]×1	147712	
stage 20	40×40	C3[256,256]×1	296448	C3[256,128]×3	173312	
stage 21	40×40	Conv[256,256]×1	590336	Conv[256,128]×1	147712	
stage 23	20×20	C3[512,512]×1	1182720	C3[256,128]×3	173312	
GFL	OPs	16.5		15.7		
Parms	s(all)	7235389		5398845		
mA	P.5	0.619		0.640		

achieves the highest accuracy among numerous comparative models, its computational complexity remains significantly higher than that of YOLO models when using a ResNet18 backbone. Additionally, due to its complex transformer architecture, the model deviates from the lightweight direction. While YOLOv7-tiny boasts lower computational complexity and model parameter count while achieving an mAP.5 above 0.4, its detection accuracy still lags behind that of YOLOv5-s. Therefore, YOLOv5-s is appropriately selected as the baseline model for subsequent experiments.

Baseline Model pruning: In the neck structure of the YOLOv5-s model, the C3 module and convolutional layers exhibit a wide and shallow shape (deeper channel counts and fewer convolution module stackings) as the downsampling factor increases. Specifically, when the 640×640 input image is downsampled to 80 \times 80 (P3), 40 \times 40 (P4), and 20 \times 20 (P5) in the neck structure, the number of stacked C3 modules is 1, and the channel numbers of the C3 modules and the convolutional layers in the three feature maps are 128, 256, and 512, respectively. When training YOLOv5-s on the antenna dataset, the visualization of the outputs of the P3, P4, and P5 feature layers revealed that increasing the number of channels in the neck section did not necessarily imply the extraction of richer features. As shown in Fig. 6(c), among the 256 channels in the P4 feature map, there were varying numbers of highly redundant patterns, and this redundancy was even more prevalent in the 512 channels of the P5 feature map. These

redundant feature maps impeded the model's detection speed and, to some extent, reduced the model's detection accuracy, particularly when these highly repetitive features were ineffective. Therefore, we conducted experiments by reducing the output channels of the C3 modules and convolutional layers corresponding to the P4 and P5 feature layers by half. The results demonstrated that this pruning operation not only made the model lighter but also improved the detection accuracy, as shown in Fig. 6(b). We reasonably conclude that widening the model by increasing only the number of channels is not effective for feature extraction and fusion on the antenna dataset; it increases the number of model parameters, reduces the computational speed, and even brings negative gains.

Based on this discovery, we reduced the channel dimensions of the P3, P4, and P5 to 128. While making the model "narrower" in terms of its width, we enhanced its learning capability and robustness by increasing the number of stacked C3 layers. After pruning the model's neck structure into a "narrow and deep" configuration, the overall model not only became more lightweight but also achieved even higher detection accuracy, as shown in Table II.

Ablation Experiment: As shown in Table III and Fig. 7(a), the baseline model was tested via ablation experiments on the COCO dataset. After pruning YOLOv5-s, both the number of parameters and the computational complexity were effectively reduced. Simultaneously, all the detection indices are comparable to or even surpass the original version, among



Fig. 6: (a) shows an example of an input antenna image and detection result, (b) compares the performance results of the baseline model after channel pruning on the neck, and (c) shows the visualization results of the P4 feature layer of the baseline model were presented. These visualizations unveiled the presence of a considerable degree of information redundancy among the 256 channels. This redundancy was notably characterized by a pronounced repetition of feature maps. To illustrate this phenomenon, we have highlighted two specific groups as exemplary instances.

TABLE III: Ablation Experiment on the COCO Dataset

model	Param.	GFLOPs	mAP.5	mAP.5:.95	mAP.5(small)	mAP.5(medium)	mAP.5(large)
YOLOv5-s [13]	7.23M	16.49	0.572	0.374	0.212	0.423	0.490
YOLOv5s-pruning	5.39M	15.67	0.570	0.385	0.222	0.430	0.498
ours(DSLK-Block)	5.35M	14.97	0.584	0.395	0.236	0.441	0.508
ours(DSLK-Block+DSLKVit-Block)	6.13M	16.18	0.599	0.410	0.245	0.455	0.535

	TABLE IV	: Ablation	1 Experi	ment on the	e Antenna Dat	aset
nodel	Param.	GFLOPs	mAP.5	mAP.5:.95	Yagi Antenna	Plate Log Antenna

model	Param.	GFLOPs	mAP.5	mAP.5:.95	Yagi Antenna	Plate Log Antenna	Patch Antenna
					mAP.5	mAP.5	mAP.5
YOLOv5-s [13]	7.23M	16.49	0.619	0.339	0.682	0.794	0.383
YOLOv5s-pruning	5.39M	15.67	0.640	0.356	0.666	0.804	0.451
ours(DSLK-Block)	5.35M	14.97	0.663	0.363	0.727	0.793	0.468
ours(DSLK-Block+DSLKVit-Block)	6.13M	16.18	0.692	0.374	0.739	0.825	0.512

which the mAP.5:.95 improved significantly, which shows that the structural pruning operations discussed in Chapter III are reasonable. On this basis, the model adds the feature processing module DSLK-Block to replace the C3 structure of the original model, which further reduces the number of parameters by 26% compared to the original YOLOv5-s. The computational complexity is also reduced by 1.52 GFLOPs. In terms of accuracy, this version of the model shows a more significant improvement in both mAP.5 and mAP.5:.95, with 11.3% and 4.26% improvements in small and medium-sized targets, respectively; these findings show that DSLK-Block plays a significant role in small object detection. Furthermore, when we constructed the PAN structure in the neck using the transformer module based on DSLK-Block, the final model was obtained. Due to the balanced relationship between computing resources and accuracy improvement, DSLKVitBlock is applied to only two feature layers of smaller sizes, P4 and P5, corresponding to medium- and large-sized objects, respectively, in the detection task. The final model has an increase in the number of parameters and the complexity of operations compared to the version using only the DSLK-Block, but both are lower than YOLOv5-s. The experimental results also show that, compared to YOLOv5-s, the detection accuracy of large targets is improved by 9.2%, which leads to a 9.6% improvement in mAP.5:.95; this is a more significant effect, while other indicators are also improved by approximately 0.03.

Table IV shows the ablation experiments of the proposed model on the antenna interference source dataset. The overall performance improvement is similar to that of the ablation experiments on the COCO dataset. The pruned version achieves comparable levels of all the metrics compared to the original



Fig. 7: Figure (a) shows the corresponding changes in accuracy and number of parameters during the ablation experiment of YOLO-Ant on the COCO dataset. Figure (b) shows the performance of two modules, DSLK-Block and DSLKVit-Block, retrofitting different models on the antenna dataset. Figure (c) shows the performance of different models on the antenna dataset in three dimensions: accuracy, computational complexity and number of parameters. The models enclosed by the green box generally exhibit higher accuracy compared to those enclosed by the red box but deviate from the lightweight design principle.

model	Year	Param.	GFLOPs	mAP.5	mAP.5:.95	Yagi Antenna mAP.5	Plate Log Antenna mAP.5	Patch Antenna mAP.5
RetinaNet(ResNet50) [26]	ICCV'2017	37.97M	191.43	0.477	0.237	-	-	-
+DSLK-Block		36.80M	184.92	0.487	0.241	-	-	-
+DSLK-Block+DSLKVit-Block		38.26M	187.85	0.492	0.243	-	-	-
YOLOv3-tiny [25]	'2018	8.85M	13.17	0.32	0.156	0.288	0.645	0.0274
+DSLK-Block		6.72M	8.94	0.36	0.145	0.414	0.596	0.0682
+DSLK-Block+DSLKVit-Block		7.58M	11.45	0.38	0.17	0.37	0.609	0.16
FCOS(ResNet50) [58]	ICCV'2019	32.12M	60.59	0.533	0.264	-	-	-
+DSLK-Block		31.15M	58.49	0.549	0.268	-	-	-
+DSLK-Block+DSLKVit-Block		31.94M	58.49	0.561	0.268	-	-	-
YOLOv4-tiny [30]	'2020	5.88M	16.18	0.28	0.134	-	-	-
+DSLK-Block		3.34M	14.53	0.292	0.145	-	-	-
+DSLK-Block+DSLKVit-Block		5.06M	15.76	0.322	0.156	-	-	-
YOLOx-s [59]	CVPR'2021	8.94M	26.76	0.579	0.350	0.650	0.781	0.306
+DSLK-Block		8.32M	26.62	0.603	0.334	0.679	0.789	0.341
+DSLK-Block+DSLKVit-Block		9.61M	27.81	0.623	0.350	0.673	0.763	0.434
YOLOv7-tiny [31]	CVPR'2023	6.23M	13.86	0.42	0.217	0.476	0.642	0.141
+DSLK-Block		5.77M	12.93	0.446	0.227	0.512	0.696	0.13
+DSLK-Block+DSLKVit-Block		5.79M	12.89	0.502	0.255	0.594	0.706	0.206

TABLE V: Validity of the proposed module on other models (on Antenna Dataset)

version; both the version with only DSLK-Block and the final model show significant improvements compared to YOLOv5s, and the final model achieves the best performance in all the metrics. Notably, the version using only the DSLK-Block has the most significant improvement in Yagi antenna detection compared to the other two antennas, while the final version using the DSLKVit-Block has a more significant improvement in the plate log antenna and patch antenna. We believe that the efficient feature extraction capability of DSLK-Block plays a crucial role due to the more complex shape of the Yagi antenna. The plate log antenna and patch antenna are more fixed in shape and color, and simple local convolution cannot extract more effective information; therefore, a transformer, a selfattention mechanism that focuses more on global information, can effectively extract the feature information of the target and background and combine the target with the surrounding environment information, making the improvement of detection accuracy more obvious.

In general, combining the ablation experiments performed

by the model on the COCO dataset and the antenna dataset, we can conclude the following. a. Pruning on the neck structure based on the original version can effectively reduce the model redundancy and still ensure accuracy while effectively reducing the number of parameters and computational complexity. b. Our proposed DSLK-Block can effectively improve the feature extraction capability of models for small-sized targets, as well as detect complex environments, and is lightweight. c. The DSLKVit-Block employed in the final model offers a significant advantage in effectively utilizing global information, leading to improvements in accuracy across all aspects. Since this structure acts on the output layer for larger size object detection, it makes the model improvement for such objects more obvious and directly leads to a significant improvement in the mAP.5:.95 metric.

Validity of the proposed module for other models: In our ablation experiment, we demonstrated that modifying the model structure and adding two modules, DSLK-Block and DSLKVit-Block, can effectively improve the baseline model

TABLE VI: Comparative Experiments on the COCO Dataset

model	Param.	GFLOPs	mAP.5	mAP.5:.95	mAP.5(small)	mAP.5(medium)	mAP.5(large)
EfficientDet-D1 [29]	6.6M	11.51	0.586	0.396	0.179	0.443	0.560
YOLOv5-s [13]	7.2M	16.5	0.568	0.374	-	-	-
YOLOv6-tiny [67]	15.0M	36.7	0.566	0.403	-	-	-
DAMO-YOLO-tiny [68]	8.5M	18.1	0.580	0.418	0.230	0.461	0.585
PPYOLOE-s [69]	7.9M	17.4	0.605	0.430	0.232	0.464	0.569
YOLOv7-tiny [31]	6.2M	13.7	0.567	0.387	0.188	0.424	0.519
Mobile-Former-508M(RetinaNet) [46]	8.4M	-	0.583	0.380	0.229	0.412	0.497
EdgeViT-XXS(RetinaNet) [41]	13.1M	-	0.590	0.387	0.224	0.420	0.516
PVT-tiny(RetinaNet) [61]	23.0M	50.94	0.569	0.367	0.226	0.388	0.500
ConT-m-tiny(RetinaNet) [70]	27.0M	217.2	0.581	0.379	0.230	0.406	0.504
our	6.1M	16.2	0.599	0.410	0.245	0.455	0.535

TABLE VII: Comparative Experiments on the VisDrone Dataset

model	Param.	GFLOPs	mAP.5	mAP.5:.95
FCOS [58]	-	-	0.258	0.142
VFNet [71]	-	-	0.288	0.168
TOOD [60]	-	-	0.294	0.181
RetinaNet(ResNet50) [26]	37.97M	191.43	0.214	0.118
YOLOx-s [59]	8.97M	26.93	0.254	0.133
YOLOv5-s [13]	7.23M	16.5	0.284	0.154
YOLOv7-tiny [31]	6.23M	13.86	0.298	0.153
ours	6.13M	16.18	0.295	0.163

in terms of parameter size, computational complexity, and accuracy. To further demonstrate the effectiveness and applicability of these two modules, we conducted various model comparisons on the antenna dataset. After each model was tested on the antenna dataset, we added the two modules in turn to the model before further comparison; the experimental results are shown in Table V and Fig. 7(b). The two modules we proposed perform well on various commonly used detectors, and the improvement effect they bring is consistent with the experimental trend on the baseline. The introduction of the DSLK-Block to replace the original convolution modules resulted in a significant improvement in the lightweight nature of all the models. However, due to differences in the complexity of their respective neck structure designs, there were substantial variations in the change in model parameter count with the introduction of the DSLKVit-Block. Nevertheless, without exception, the introduction of both modules led to a noticeable enhancement in the model's detection accuracy. Fig. 8 shows a comparison of the detection performances of our model and the baseline model on the antenna dataset. The DSLK-Block proved instrumental in detecting numerous small targets, particularly patch antennas, which the baseline model failed to identify. When faced with the common background of balcony fences, the baseline model is prone to mistaken local regions for the Yagi antenna, resulting in many false alarms. When encountering the problem of large interclass differences presented by plate log antennas in multiangle shooting, the baseline model's detection resulted in many false positives. To address these issues, our proposed method combines a CNN and a transformer, introducing the DSLKVit-Block module to effectively resolve the problem. A comparison of the detection figures clearly reveals that our method results in fewer false alarms and misses after the correct targets are detected. Finally, we compare the performance of YOLO-Ant with that of many classical and cutting-edge detectors on the antenna dataset

TABLE VIII: Model Detection Speed Test

model	Param.(M)	GFLOPs	FPS
YOLOv5-s [13]	7.23	16.5	32.43
YOLOv7-tiny [31]	6.23	13.86	43.82
DETR(ResNet50) [11]	41.58	79.52	3.89
RT-DETR(ResNet18) [66]	20.00	60.00	10.89
ours	6.1	16.2	35.87

in three dimensions, namely, the model parameter number, accuracy, and computational complexity, as shown in Fig. 7(c), and YOLO-Ant achieves the best performance.

Experiments on Public Datasets: To demonstrate the generalizability, robustness, and capability of YOLO-Ant, which can excel not only on antenna datasets but also on public datasets, we conducted comparisons with numerous models on two common datasets, COCO and VisDrone. As shown in Tables VI and VII, we compare our model with several other recent outstanding one-stage detectors and several lightweight transformer-based detectors. The results indicate that YOLO-Ant has the fewest parameters and relatively lower computational complexity, outperforming other transformerbased detection models by a significant margin. In terms of detection accuracy, our approach slightly lags behind the top-performing PPYOLO-s in terms of mAP.5 and slightly lags behind PPYOLO-s and DAMO-YOLO-tiny in terms of mAP.5:.95. However, in small object detection, our approach achieves the highest precision among all the models. Fig. 9 shows the detection performance comparison between our model and YOLOv5-s on the COCO dataset. On the VisDrone dataset, YOLO-Ant also achieves competitive results. Overall, considering the parameter count, computational complexity, and detection accuracy, YOLO-Ant is one of the top contenders among lightweight object detectors.

Model Detection Speed Test: In this experiment, we conducted a speed comparison between YOLO-Ant and several transformer models, as well as YOLO series models known for their lightweight design. We converted all the model weights trained on the antenna dataset into ONNX format and tested them using the OPENVINO tool on a platform with low-computational power(Intel $Core^{TM}$ i9-11900 CPU). Testing was performed on the validation set of the antenna dataset. The comparative results are presented in Table VIII.

The table shows that YOLOv7-tiny achieves an FPS (frames per second) rate exceeding 40, making it the fastest model. However, this lightweight design comes at the cost of reduced detection accuracy. On the other hand, while YOLO-Ant



Fig. 8: Comparison of detection performance between YOLO-Ant and baseline models on the antenna interference source dataset.



Fig. 9: Comparison of detection performance between YOLO Ant and baseline models in the COCO dataset. The three images from left to right in each group represent the ground truth, the detection effect of baseline and YOLO-Ant, respectively.

exhibits a slightly slower detection speed than YOLOv7-tiny does, it outperforms the baseline model YOLOv5-s to a small extent. Additionally, the lightweight design of the proposed model is significantly superior to that of the two transformerbased models. This demonstrates the success of our structural design, both in terms of accuracy and lightweight efficiency.

Due to the extensive use of depthwise separable convolutions in our designed modules, which require more memory bandwidth, the I/O (input/output) read speed of the device became the speed bottleneck of the model. We believe that in future work, further optimizations can be made in this regard. This approach enables YOLO-Ant to achieve not only the highest detection accuracy in antenna detection but also further enhance its lightweight design.

V. CONCLUSION

Taking UAV detection of antenna interference sources as the starting point, we propose a lightweight object detector with improved YOLOv5. Initially, to ensure the model is lightweight and adaptable to subsequent modifications, the new network is first pruned based on YOLOv5, which effectively reduces the number of model parameters and computational complexity while ensuring accuracy. To address the challenges posed by small target sizes and complex backgrounds in antenna detection tasks, we propose an efficient and lightweight convolutional module called DSLK-Block. Furthermore, we introduce a lightweight transformer structure integrated with DSLK-Block, which is applied to the network's neck. This combination significantly enhances the network's ability to extract and process features. The new model not only is effective on the antenna dataset but also achieves competitive results on public datasets such as COCO. In future work, we will further explore the integration of our current efforts with drone inspection technology. Simultaneously, we will incorporate traditional equipment such as spectrum analyzers required for conventional antenna interference source identification. The subsequent objective will be to refine the current approach to create a comprehensive intelligent unmanned detection system.

REFERENCES

- X. Lu, J. Ji, Z. Xing, and Q. Miao, "Attention and feature fusion ssd for remote sensing object detection," *IEEE Transactions on Instrumentation* and Measurement, vol. 70, pp. 1–9, 2021.
- [2] Z. Yang, Z. Xu, and Y. Wang, "Bidirection-fusion-yolov3: An improved method for insulator defect detection using uav image," *IEEE Transactions on Instrumentation and Measurement*, vol. 71, pp. 1–8, 2022.
- [3] K. Hao, G. Chen, L. Zhao, Z. Li, Y. Liu, and C. Wang, "An insulator defect detection model in aerial images based on multiscale feature pyramid network," *IEEE Transactions on Instrumentation and Measurement*, vol. 71, pp. 1–12, 2022.
- [4] K. Han, Y. Wang, Q. Tian, J. Guo, C. Xu, and C. Xu, "Ghostnet: More features from cheap operations," in *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, 2020, pp. 1580– 1589.
- [5] M. Tan and Q. Le, "Efficientnetv2: Smaller models and faster training," in *International conference on machine learning*. PMLR, 2021, pp. 10096–10106.
- [6] J.-S. Lim, M. Astrid, H.-J. Yoon, and S.-I. Lee, "Small object detection using context and attention," in 2021 International Conference on Artificial Intelligence in Information and Communication (ICAIIC). IEEE, 2021, pp. 181–186.

- [7] X. Yang, J. Yang, J. Yan, Y. Zhang, T. Zhang, Z. Guo, X. Sun, and K. Fu, "Scrdet: Towards more robust detection for small, cluttered and rotated objects," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8232–8241.
- [8] C. Yang, Z. Huang, and N. Wang, "Querydet: Cascaded sparse query for accelerating high-resolution small object detection," in *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, 2022, pp. 13 668–13 677.
- [9] S.-J. Ji, Q.-H. Ling, and F. Han, "An improved algorithm for small object detection based on yolo v4 and multi-scale contextual information," *Computers and Electrical Engineering*, vol. 105, p. 108490, 2023. [Online]. Available: https://www.sciencedirect.com/ science/article/pii/S0045790622007054
- [10] Y. Cao, J. Xu, S. Lin, F. Wei, and H. Hu, "Gcnet: Non-local networks meet squeeze-excitation networks and beyond," *arXiv*, 2019.
- [11] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16.* Springer, 2020, pp. 213–229.
- [12] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr *et al.*, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 6881–6890.
- [13] G. Jocher, A. Stoken, A. Chaurasia, J. Borovec, NanoCode012, TaoXie, Y. Kwon, K. Michael, L. Changyu, J. Fang, A. V, Laughing, tkianai, yxNONG, P. Skalski, A. Hogan, J. Nadar, imyhxy, L. Mammana, AlexWang1900, C. Fati, D. Montes, J. Hajek, L. Diaconu, M. T. Minh, Marc, albinxavi, fatih, oleg, and wanghaoyang0106, "ultralytics/yolov5: v6.0 - YOLOv5n 'Nano' models, Roboflow integration, TensorFlow export, OpenCV DNN support," Oct. 2021. [Online]. Available: https://doi.org/10.5281/zenodo.5563715
- [14] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [16] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," arXiv preprint arXiv:1409.1556, 2014.
- [17] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [18] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision* and pattern recognition, 2016, pp. 770–778.
- [19] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.
- [20] R. Girshick, "Fast r-cnn," in Proceedings of the IEEE international conference on computer vision, 2015, pp. 1440–1448.
- [21] R. Faster, "Towards real-time object detection with region proposal networks," Advances in neural information processing systems, vol. 9199, no. 10.5555, pp. 2969239–2969250, 2015.
- [22] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779– 788.
- [23] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Computer Vision– ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14.* Springer, 2016, pp. 21–37.
- [24] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 7263–7271.
- [25] —, "Yolov3: An incremental improvement," 2018.
- [26] 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.
- [27] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, pp. 211–252, 2015.

- [28] G. Ghaisi, T.-Y. Lin, R. Pang, and Q. V. L. NAS-FPN, "Learning scalable feature pyramid architecture for object detection."
- [29] M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 10781–10790.
- [30] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," arXiv preprint arXiv:2004.10934, 2020.
- [31] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," arXiv preprint arXiv:2207.02696, 2022.
- [32] J. Li, Y. Zhang, P. Yun, G. Zhou, Q. Chen, and R. Fan, "RoadFormer: Duplex transformer for rgb-normal semantic road scene parsing," *arXiv* preprint arXiv:2309.10356, 2023.
- [33] R. Fan, S. Guo, and M. J. Bocus, Autonomous driving perception. Springer Nature, 2023.
- [34] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 7132–7141.
- [35] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.
- [36] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, 2019, pp. 3146– 3154.
- [37] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [38] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [39] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 10012–10022.
- [40] S. Mehta and M. Rastegari, "Mobilevit: light-weight, generalpurpose, and mobile-friendly vision transformer," arXiv preprint arXiv:2110.02178, 2021.
- [41] J. Pan, A. Bulat, F. Tan, X. Zhu, L. Dudziak, H. Li, G. Tzimiropoulos, and B. Martinez, "Edgevits: Competing light-weight cnns on mobile devices with vision transformers," in *Computer Vision–ECCV 2022:* 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XI. Springer, 2022, pp. 294–311.
- [42] D. Meng, X. Chen, Z. Fan, G. Zeng, H. Li, Y. Yuan, L. Sun, and J. Wang, "Conditional detr for fast training convergence," in *Proceedings* of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 3651–3660.
- [43] F. Li, H. Zhang, S. Liu, J. Guo, L. M. Ni, and L. Zhang, "Dn-detr: Accelerate detr training by introducing query denoising," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13619–13627.
- [44] J. Guo, K. Han, H. Wu, Y. Tang, X. Chen, Y. Wang, and C. Xu, "Cmt: Convolutional neural networks meet vision transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12175–12185.
- [45] Z. Peng, W. Huang, S. Gu, L. Xie, Y. Wang, J. Jiao, and Q. Y. Conformer, "Local features coupling global representations for visual recognition. in 2021 ieee," in *CVF International Conference on Computer Vision*, *ICCV*, 2021, pp. 357–366.
- [46] Y. Chen, X. Dai, D. Chen, M. Liu, X. Dong, L. Yuan, and Z. Liu, "Mobile-former: Bridging mobilenet and transformer," in *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 5270–5279.
- [47] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint* arXiv:1704.04861, 2017.
- [48] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11976–11986.
- [49] X. Ding, X. Zhang, J. Han, and G. Ding, "Scaling up your kernels to 31x31: Revisiting large kernel design in cnns," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11963–11975.

- [50] X. Pan, C. Ge, R. Lu, S. Song, G. Chen, Z. Huang, and G. Huang, "On the integration of self-attention and convolution," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 815–825.
- [51] H. Zhang, W. Hu, and X. Wang, "Parc-net: Position aware circular convolution with merits from convnets and transformer," in *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October* 23–27, 2022, Proceedings, Part XXVI. Springer, 2022, pp. 613–630.
- [52] L. Guan, L. Jia, Z. Xie, and C. Yin, "A lightweight framework for obstacle detection in the railway image based on fast region proposal and improved yolo-tiny network," *IEEE Transactions on Instrumentation* and Measurement, vol. 71, pp. 1–16, 2022.
- [53] Y. Dai, W. Liu, H. Wang, W. Xie, and K. Long, "Yolo-former: Marrying yolo and transformer for foreign object detection," *IEEE Transactions* on *Instrumentation and Measurement*, vol. 71, pp. 1–14, 2022.
- [54] L. Wang, H. Qin, X. Zhou, X. Lu, and F. Zhang, ⁴R-yolo: A robust object detector in adverse weather," *IEEE Transactions on Instrumentation and Measurement*, vol. 72, pp. 1–11, 2022.
- [55] M. A. A. Al-qaness, A. A. Abbasi, H. Fan, R. A. Ibrahim, S. H. Alsamhi, and A. Hawbani, "An improved yolo-based road traffic monitoring system," vol. 103, pp. 211–230, 2021.
- [56] F. Dang, D. Chen, Y. Lu, and Z. Li, "Yoloweeds: A novel benchmark of YOLO object detectors for multi-class weed detection in cotton production systems," *Comput. Electron. Agric.*, vol. 205, p. 107655, 2023. [Online]. Available: https://doi.org/10.1016/j.compag. 2023.107655
- [57] S. Guo et al., "UDTIRI: An online open-source intelligent road inspection benchmark suite," *IEEE Transactions on Intelligent Transportation* Systems, DOI: 10.1109/TITS.2024.3351209.
- [58] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional onestage object detection," arXiv preprint arXiv:1904.01355, 2019.
- [59] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, "Yolox: Exceeding yolo series in 2021," arXiv preprint arXiv:2107.08430, 2021.
- [60] C. Feng, Y. Zhong, Y. Gao, M. R. Scott, and W. Huang, "Tood: Taskaligned one-stage object detection," 2021.
- [61] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 568–578.
- [62] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, and L. Zhang, "Dab-detr: Dynamic anchor boxes are better queries for detr," 2022.
- [63] G. Jocher, A. Chaurasia, and J. Qiu, "YOLO by Ultralytics," Jan. 2023. [Online]. Available: https://github.com/ultralytics/ultralytics
- [64] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. M. Ni, and H.-Y. Shum, "Dino: Detr with improved denoising anchor boxes for end-toend object detection," 2022.
- [65] S. Woo, S. Debnath, R. Hu, X. Chen, Z. Liu, I. S. Kweon, and S. Xie, "Convnext v2: Co-designing and scaling convnets with masked autoencoders," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 16133–16142.
- [66] W. Lv, Y. Zhao, S. Xu, J. Wei, G. Wang, C. Cui, Y. Du, Q. Dang, and Y. Liu, "Detrs beat yolos on real-time object detection," 2023.
- [67] C. Li, L. Li, H. Jiang, K. Weng, Y. Geng, L. Li, Z. Ke, Q. Li, M. Cheng, W. Nie, Y. Li, B. Zhang, Y. Liang, L. Zhou, X. Xu, X. Chu, X. Wei, and X. Wei, "Yolov6: A single-stage object detection framework for industrial applications," 2022.
- [68] X. Xu, Y. Jiang, W. Chen, Y. Huang, Y. Zhang, and X. Sun, "Damoyolo: A report on real-time object detection design," *arXiv preprint* arXiv:2211.15444, 2022.
- [69] S. Xu, X. Wang, W. Lv, Q. Chang, C. Cui, K. Deng, G. Wang, Q. Dang, S. Wei, Y. Du, and B. Lai, "Pp-yoloe: An evolved version of yolo," 2022.
- [70] H. Yan, Z. Li, W. Li, C. Wang, M. Wu, and C. Zhang, "Contnet: Why not use convolution and transformer at the same time?" arXiv preprint arXiv:2104.13497, 2021.
- [71] H. Zhang, Y. Wang, F. Dayoub, and N. Sünderhauf, "Varifocalnet: An iou-aware dense object detector," arXiv preprint arXiv:2008.13367, 2020.



Xiaoyu Tang (Member, IEEE) received the B.S. degree from South China Normal University in 2003 and the M.S. degree from Sun Yat-sen University in 2011. He is currently pursuing the Ph.D. degree with South China Normal University. He is working with the School of Physics, South China Normal University, where he engaged in information system development. His research interests include machine vision, intelligent control, and the Internet of Things. He is a member of the IEEE ICICSP Technical Committee.



Rui Fan (Senior Member, IEEE) received the B.Eng. degree in Automation from the Harbin Institute of Technology in 2015 and the Ph.D. degree (supervisors: Prof. John G. Rarity and Prof. Naim Dahnoun) in Electrical and Electronic Engineering from the University of Bristol in 2018. He worked as a Research Associate (supervisor: Prof. Ming Liu) at the Hong Kong University of Science and Technology from 2018 to 2020 and a Postdoctoral Scholar-Employee (supervisors: Prof. Linda M. Zangwill and Prof. David J. Kriegman) at the University of

California San Diego between 2020 and 2021. He began his faculty career as a Full Research Professor with the College of Electronics & Information Engineering at Tongji University in 2021, and was then promoted to a Full Professor in the same college, as well as at the Shanghai Research Institute for Intelligent Autonomous Systems in 2022.

Prof. Fan served as an Associate Editor for ICRA'23 and IROS'23/24, an Area Chair for ICIP'24, and a Senior Program Committee Member for AAAI'23/24. He is the general chair of the AVVision community and organized several impactful workshops and special sessions in conjunction with WACV'21, ICIP'21/22/23, ICCV'21, and ECCV'22. He was honored by being included in the Stanford University List of Top 2% Scientists Worldwide in both 2022 and 2023, recognized on the Forbes China List of 100 Outstanding Overseas Returnees in 2023, and acknowledged as one of Xiaomi Young Talents in 2023. His research interests include computer vision, deep learning, and robotics.



Xingming Chen received his B.Eng. degree from the School of Physics and Telecommunication Engineering at South China Normal University in 2021. Currently, he is pursuing the M.E. degree with the Department of Electronics and Information Engineering at the same institution. His research focuses on computer vision and deep learning.



Jintao Cheng received his B.Eng. degree from the School of Physics and Telecommunications Engineering at South China Normal University in 2021. His research focuses on computer vision, SLAM, and deep learning.



Chengxi Zhang (Member, IEEE) was born in Shandong, China, in February 1990. He received the B.S. and M.S. degrees in microelectronics and solid-state electronics from Harbin Institute of Technology, Harbin, China, in 2012 and 2015, respectively, and the Ph.D. degree in control science and engineering from Shanghai Jiao Tong University, Shanghai, China, in 2019.

He was a Postdoctoral Fellow with the School of Electronic and Information Engineering, Harbin Institute of Technology (Shenzhen Campus), Shen-

zhen, China, from 2020 to 2022. His research interests are embedded system software and hardware design, information fusion, and control theory.



Jin Wu (Member, IEEE) was born in Zhenjiang, China, in 1994. He received a B.S. degree from the University of Electronic Science and Technology of China, Chengdu, China. From 2013 to 2014, He was a visiting student with Groep T, Katholieke Universiteit Leuven (KU Leuven). He is currently pursuing a Ph.D. degree in the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology (HKUST), Hong Kong. He has co-authored over 120 technical papers in representative journals and conference proceedings.

He was awarded the outstanding reviewer of IEEE TRANSACTIONS ON INSTRUMENTATION AND MEASUREMENT in 2021. He is now a Review Editor of Frontiers in Aerospace Engineering and an invited guest editor for 5 special issues of MDPI. He is also in the IEEE Consumer Technology Society (CTSoc), as a committee member and publication liaison. He was a committee member for the IEEE CoDIT conference in 2019, a special section chair for the IEEE ICGNC conference in 2021, a special session chair for the 2023 IEEE International Conference on Intelligent Transportation Systems (ITSC), a Track Chair for the 2024 IEEE International Conference on Consumer Electronics (ICCE) and a Chair for the 2024 IEEE CTSoc Gaming, Entertainment and Media (GEM) conference. He was selected as the World's Top 2% Scientist by Stanford University and Elsevier, in the 2020, 2021 and 2022 year round.



Zebo Zhou was born in November, 1982 in Yongchuan, Chongqing, China. He received the B.Sc. and M.Sc. degrees from the School of Geodesy and Geomatics, Wuhan University, Wuhan, China, in 2004 and 2006, respectively, and the Ph.D. degree from the College of Surveying and Geoinformatics, Tongji University, Shanghai, China, in 2009. In 2009 and 2015, he was a Visiting Fellow with Surveying and Geospatial Engineering Group, within the School of Civil and Environmental Engineering, University of New South Wales, Kensington, NSW,

Australia. He is currently an Associate Professor with the School of Aeronautics and Astronautics, University of Electronic Science and Technology of China, Chengdu, China. His research interests include GNSS navigation and positioning, GNSS/INS integrated navigation, and multisensor fusion.

He has been an In Charge of projects of the National Natural Science Foundation of China and has taken part in the National 863 High-tech Founding of China. He was the Guest Editor of several special issues published on *International Journal of Distributed Sensor Networks and Asian Journal of Control*. He has been presenting related works on the annual conference of the institute of navigation, the annual Chinese satellite navigation conferences.