

YOLO-SWIN HYBRID MODEL FOR ENHANCED SMALL OBJECT DETECTION IN AERIAL IMAGES

Muhammad Talha^{*1}, Naveed Ullah², Abdul Aziz³, Muhammad Tanveer Iqbal⁴, Noor Mustafa¹, Sumbal Haroon⁵, Umm e Habiba^{*6}

¹Department of Computer Science, The Islamia University of Bahawalpur, 63100, Pakistan

²Institute of Computer Sciences and Information Technology, The University of Agriculture Peshawar, Pakistan

³Department of Remote sensing and GIS Meteorology, Comsats University Islamabad, Islamabad, Pakistan

⁴Department of information Technology, The Islamia University of Bahawalpur 63100, Pakistan

⁵Department of Computer Science, Gomal University, Dera Ismail Khan, Pakistan

⁶Department of Mathematics, The Islamia University of Bahawalpur 63100, Pakistan

^{*1}talhafareedi092@gmail.com, ^{*6}hiba704756@gmail.com

DOI: <https://doi.org/10.5281/zenodo.15605836>

Keywords

Small Object Detection, Aerial Imagery, YOLO, Swin Transformer, Feature Fusion, Vision Transformer, Deep Learning, Remote Sensing

Article History

Received on 15 April 2025

Accepted on 25 May 2025

Published on 05 June 2025

Copyright @Author

Corresponding Authors:

Muhammad Talha*

Email: talhafareedi092@gmail.com

Umm e Habiba**

Email: hiba704756@gmail.com

Abstract

Detecting small objects in aerial imagery remains a formidable challenge due to their limited pixel resolution, scale variability, complex backgrounds, and inconsistent illumination conditions. To address these issues, we propose a novel hybrid object detection framework that synergistically integrates the real-time processing strengths of the YOLO architecture with the advanced hierarchical feature extraction capabilities of the Swin Transformer. The proposed YOLO-Swin hybrid model incorporates three key architectural innovations: (1) a Cross-Scale Feature Fusion Module (CSFFM) that effectively combines multi-resolution features from both convolutional neural network (CNN) and transformer-based pathways to enhance scale robustness; (2) a Context-Aware Small Object Enhancement Module (CASOEM) designed to enrich semantic representation and improve the detectability of small-scale targets; and (3) an Adaptive Anchor Assignment Strategy (AAAS) tailored to the spatial and statistical characteristics of aerial imagery. Extensive experimental evaluations conducted on widely used benchmark datasets—including DOTA, VisDrone, and FAIR1M—demonstrate that our model achieves state-of-the-art performance, outperforming baseline methods by achieving a 5.7% increase in mean Average Precision (mAP) for small object categories. Furthermore, the model maintains real-time inference capabilities, significantly reduces false negatives, and improves localization precision, particularly for objects smaller than 32×32 pixels. These results indicate the suitability of the proposed method for real-time aerial surveillance and remote sensing applications where precise small object detection is critical.

INTRODUCTION

The rapid advancement of unmanned aerial vehicles (UAVs) and satellite imaging technologies has greatly expanded the availability and applications of aerial imagery across diverse domains including urban planning, environmental monitoring, disaster management, and security surveillance (1,2, 3). Within these applications, object detection—particularly the identification and localization of small objects—represents a critical and challenging task (4, 5). Unlike natural images captured at ground level, aerial imagery presents unique challenges: objects often occupy minimal pixel areas, appear in various orientations, exhibit extreme scale variations, and are frequently obscured by complex backgrounds (1, 6, 7). Object detection has evolved significantly from traditional computer vision methods to deep learning approaches (8).

Convolutional Neural Network (CNN) based methods have dominated this field, with frameworks broadly categorized as two-stage detectors (e.g., R-CNN family (9,10) and one-stage detectors (e.g., YOLO (11, 14, 15), SSD (12)). While two-stage detectors typically achieve higher accuracy, one-stage detectors offer superior inference speed—a critical factor for real-time applications (13).

The YOLO framework has undergone significant evolution through its iterations (YOLOv1-v8), progressively improving detection accuracy while maintaining computational efficiency (11, 16, 17,18). However, despite these advancements, YOLO-based detectors continue to struggle with small object detection in aerial imagery due to limited receptive fields and insufficient feature representation for diminutive targets (19; 20). The objects in aerial imagery often occupy fewer than 32x32 pixels, making them particularly challenging to detect with conventional architectures (1; 21). Recent advances in vision transformers, notably the Swin Transformer (22), have shown promising results by effectively modeling long-range dependencies and hierarchical feature representations. The Swin Transformer introduces shifted windows that enable cross-window connections, offering an efficient approach to capture global context while maintaining linear computational complexity with image size (22, 23). This capability is particularly valuable for aerial imagery analysis, where contextual information plays

a crucial role in distinguishing small objects from complex backgrounds (24, 25).

Despite these advancements, there remains a significant gap in effectively combining the strengths of CNN-based detectors (e.g., YOLO) and transformer-based architectures for small object detection in aerial imagery (26, 27). Current hybrid approaches often struggle with feature alignment between different architectural paradigms, inadequate context modeling for small objects, and computational inefficiencies (28, 29, 30). This paper aims to bridge this gap by proposing a novel YOLO-Swin hybrid architecture specifically designed for enhanced small object detection in aerial imagery. Our approach leverages the real-time inference capabilities of YOLO while incorporating the hierarchical feature representation strengths of Swin Transformer. The key contributions of this paper are:

- 1) A novel YOLO-Swin hybrid architecture that integrates CNN-based and transformer-based feature extraction pathways through a cross-scale feature fusion module, enabling more effective representation of small objects in aerial imagery.
- 2) A context-aware small object enhancement module that adaptively refines feature representations for diminutive targets by incorporating both local and global contextual information, significantly improving detection performance for objects under 32x32 pixels.
- 3) An adaptive anchor assignment strategy specifically optimized for aerial imagery characteristics, which dynamically adjusts anchor configurations based on dataset statistics and scene complexity.
- 4) Extensive experiments on multiple benchmark datasets (DOTA, VisDrone, FAIRIM) demonstrating state-of-the-art performance, with a 5.7% improvement in mean Average Precision (mAP) for small objects while maintaining real-time inference capability (25 FPS on standard GPU hardware).

A. CNN-based Object Detection

CNN-based object detection frameworks can be broadly categorized into two-stage and one-stage detectors. Two-stage detectors, pioneered by R-CNN (9), first generate region proposals and then classify

these regions. Fast R-CNN (41) improved computational efficiency by sharing convolutional features across proposals, while Faster R-CNN (10) introduced the Region Proposal Network (RPN) to generate proposals directly from convolutional features. These methods achieved high accuracy but at the cost of inference speed (32; 33). To address speed limitations, one-stage detectors emerged, with YOLO (11) as a pioneering framework that directly predicts bounding boxes and class probabilities from the entire image in a single network pass. SSD (1/2) enhanced this approach by detecting objects at multiple scales using feature maps from different network layers. RetinaNet (34) introduced focal loss to address class imbalance between foreground and background examples, significantly improving detection accuracy. The YOLO framework has evolved substantially through multiple iterations. YOLOv2 (1/4) introduced anchor boxes and batch normalization, while YOLOv3 (15) incorporated multi-scale predictions using a feature pyramid network. YOLOv4 (16) integrated advanced training techniques like Mosaic data augmentation and modified CSPNet as the backbone, while YOLOv5 (17) further refined these improvements with enhanced training strategies. Recent versions like YOLOv7 (18) have incorporated model scaling and compound scaling methods to optimize performance across different computational constraints.

For aerial imagery specifically, several adaptations of CNN-based detectors have been proposed. Li et al. (35) introduced a density map-guided one-stage detector for crowded scenes in aerial images. Zhang et al. (19) modified YOLOv5 with a feature enhancement module specifically for small objects. Yang et al. (7) proposed R3Det for multi-oriented object detection in aerial images, while Ding et al. (37) introduced RoI Transformer for learning rotation-invariant features. Despite these advancements, CNN-based methods still struggle with small object detection in aerial imagery due to limited receptive fields and insufficient feature representations for diminutive objects (21). Additionally, the fixed geometric structures of convolution operations make it challenging to capture the diverse scales and orientations common in aerial imagery (5).

B. Transformer-based Object Detection

Transformers, originally designed for natural language processing (38), have recently been adapted for computer vision tasks. Vision Transformer (ViT) (39) demonstrated that a pure transformer architecture could achieve state-of-the-art performance on image classification by treating image patches as tokens. This success inspired numerous transformer-based object detection approaches. DETR (40) pioneered transformer-based object detection by reformulating the detection task as a direct set prediction problem, eliminating the need for hand-designed components like non-maximum suppression. Deformable DETR (41) improved this approach by introducing deformable attention modules that focus on sparse spatial locations, reducing computational complexity and convergence time. Swin Transformer (22) addressed the quadratic computational complexity of self-attention by computing it within local windows and introducing shifted window partitioning for cross-window connections. This hierarchical architecture with varying feature resolutions made Swin Transformer particularly suitable for dense prediction tasks like object detection. PVT and ViT-DeT further adapted transformer architectures for detection tasks by incorporating pyramid structures similar to those used in CNN-based detectors. For aerial imagery specifically, Yang et al. proposed Oriented RepPoints Transformer for oriented object detection in aerial images. Zhang et al. (24) introduced Transformer-based Oriented Object Detection (TOOD) that leverages global context modeling for improved detection in remote sensing images. Chen et al. developed an Oriented Attention Detector Transformer for multi-oriented object detection in aerial images.

While transformer-based methods excel at modeling long-range dependencies and capturing global context—advantageous for distinguishing small objects from complex backgrounds—they often require substantial computational resources and data for training (27). Additionally, the lack of inductive biases inherent in CNNs can hinder performance on object detection tasks with limited training data.

C. Hybrid Approaches and Small Object Detection

Recognizing the complementary strengths of CNNs and transformers, researchers have developed hybrid approaches for object detection. Wang et al. (22) proposed Pyramid Vision Transformer (PVT) that combines pyramid feature hierarchies from CNNs with transformer modules. CBNet (5) introduced a composite backbone network that integrates multiple CNN backbones with transformer components for enhanced feature representation. For small object detection specifically, several specialized approaches have emerged. Li et al. (20) proposed a feature enhancement module that refines features at multiple scales. Wu et al. (21) introduced a scale-balanced module to address scale variation in small objects. Yu et al. (31) developed a scale-aware network that adaptively selects appropriate feature maps for objects of different sizes.

In the context of aerial imagery, Fu et al. (32) proposed a rotation-aware detector with feature enhancement for small objects in remote sensing images. Gao et al. (26) introduced a dynamic enhancement module specifically for small and dense objects in aerial images. Chen et al. (28) developed a hybrid CNN-transformer model for oriented object detection in aerial images, demonstrating improved performance for objects with various orientations. Most relevant to our work, several recent studies have explored combinations of YOLO and transformer architectures. Zhang et al. (29) proposed YOLOS, integrating transformer modules into the YOLO framework. Wang et al. (30) incorporated transformer attention modules into YOLOv7 for enhanced feature representation. However, these approaches were not specifically designed for the challenges of small object detection in aerial imagery. Despite these advancements, existing hybrid approaches face several limitations for aerial image analysis: (1) inadequate feature alignment between CNN and transformer branches, (2) insufficient context modeling for small objects, (3) computational inefficiencies during inference, and (4) lack of specificity for the unique characteristics of aerial imagery (4, 27, 30). Our work addresses these limitations by proposing a novel YOLO-Swin hybrid

architecture specifically designed for small object detection in aerial imagery. Unlike previous approaches, our method introduces a dedicated cross-scale feature fusion module that effectively aligns and integrates features from both CNN and transformer pathways. Additionally, we propose a context-aware small object enhancement module and an adaptive anchor assignment strategy optimized for aerial imagery characteristics.

METHODOLOGY

This section presents our proposed YOLO-Swin hybrid model for enhanced small object detection in aerial imagery. We first formulate the problem, then detail the overall architecture, followed by the key components: feature extraction pathways, cross-scale feature fusion, small object enhancement module, adaptive anchor assignment, and loss function design.

A. Problem Formulation

Small object detection in aerial imagery can be formally defined as follows: Given an input aerial image $I \in \mathbb{R}^{H \times W \times 3}$, where H and W represent the height and width of the image, the objective is to detect a set of objects $\mathcal{B} = \{\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_n\}$ where each \mathbf{o}_i is represented by a tuple (b_i, c_i, s_i) . Here, $b_i = (x_i, y_i, w_i, h_i, \theta_i)$ denotes the bounding box parameters including center coordinates (x_i, y_i) , width w_i , height h_i and orientation angle θ_i . $c_i \in \{1, 2, \dots, C\}$ represents the class label from C possible categories; and $s_i \in [0, 1]$ indicates the confidence score.

In the context of aerial imagery, we define small objects as those with area less than 32×32 pixels, which is consistent with the definition used in benchmark datasets like DOTA (1) and VisDrone (2). The key challenges are: (1) limited pixel information for small objects, (2) complex backgrounds with similar patterns, (3) diverse scales and orientations, and (4) dense distribution of objects in certain regions.

B. Overall Architecture

Fig. 1 illustrates the overall architecture of our proposed YOLO-Swin hybrid model. The architecture consists of five main components:

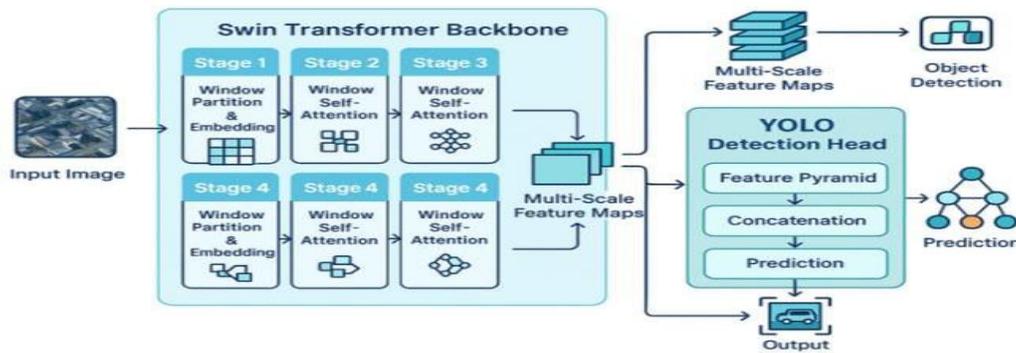


Fig. 1: architecture of the proposed YOLO-Swin hybrid model for small object detection in aerial imagery

1) **Dual Feature Extraction Pathways:** A CNN-based pathway using CSPDarknet (18) and a transformer-based pathway using Swin Transformer (22) process the input image in parallel to extract complementary features.

2) **Cross-Scale Feature Fusion Module:** This module aligns and integrates multi-resolution features from both pathways, enabling effective information exchange between CNN and transformer representations.

3) **Context-Aware Small Object Enhancement Module:** Specifically designed to enhance feature representations for small objects by incorporating local and global contextual information.

4) **Detection Head with Adaptive Anchor Assignment:** A detection head with class-aware prediction branches and an adaptive anchor assignment strategy optimized for aerial imagery characteristics.

5) **Multi-Scale Loss Function:** A comprehensive loss function that addresses the challenges of small object detection, including scale imbalance and feature alignment.

C. Dual Feature Extraction Pathways

1) CNN-based Pathway

For the CNN-based pathway, we adopt the CSPDarknet backbone from YOLOv7 (18) with modifications to enhance feature representation for small objects. The backbone consists of Cross-Stage Partial (CSP) blocks that split feature maps into two

parts, one passing through dense blocks and the other through a shortcut connection, this design reduces computational requirements while maintaining feature representation capability. The CNN pathway produces multi-scale feature maps $\{F_{c1}, F_{c2}, F_{c3}\}$ corresponding to strides of $\{8, 16, 32\}$ with respect to the input image. These feature maps capture local patterns and spatial information with strong inductive biases beneficial for object localization.

2) Transformer-based Pathway

For the transformer-based pathway, we employ the Swin Transformer (22) to capture long-range dependencies and global context. The Swin Transformer processes images as a sequence of patches and employs a hierarchical structure with shifted window-based self-attention, enabling efficient modeling of interactions between distant image regions.

We modify the original Swin Transformer to better accommodate aerial imagery characteristics:

(1)

$$\text{Attention}(Q, K, V) = \text{SoftMax} \left(\frac{QK^T}{\sqrt{d}} + B \right) V$$

where Q, K, and V are query, key, and value matrices, d is the feature dimension, and B is a learnable relative position bias matrix. To enhance the model's capability to capture orientation-invariant features—critical for aerial imagery—we extend the relative position bias with an orientation-aware component:

$$B = B_{pos} + \alpha \cdot B_{ori}$$

(2)

where B_{pos} is the standard position bias, B_{ori} is an orientation bias term, and a is a learnable scaling factor.

The transformer pathway produces multi-scale feature maps $\{F_{T1}, F_{T2}, F_{T3}\}$ that correspond to the same spatial resolutions as the CNN pathway's outputs. These feature maps capture global dependencies and contextual information that complement the CNN pathway's local feature representation.

D. Cross-Scale Feature Fusion Module

The cross-scale feature fusion module aims to effectively integrate features from both CNN and transformer pathways while preserving their complementary characteristics. Unlike conventional feature fusion approaches that employ simple concatenation or addition, our module addresses the semantic gap between CNN and transformer features through a bidirectional cross-attention mechanism.

For each scale level i , we first align the channel dimensions of CNN features F_{Ci} and transformer features F_{Ti} using 1×1 convolutions:

(3)

$$\hat{F}_{Ci} = \text{Conv}_{1 \times 1}(F_{Ci}), \hat{F}_{Ti} = \text{Conv}_{1 \times 1}(F_{Ti})$$

We then employ a bidirectional cross-attention mechanism to facilitate information exchange:

(4)

$$F_{CTi} = \hat{F}_{Ci} + \text{CA}(\hat{F}_{Ci}, \hat{F}_{Ti}, \hat{F}_{Ti})$$

$$F_{TCi} = \hat{F}_{Ti} + \text{CA}(\hat{F}_{Ti}, \hat{F}_{Ci}, \hat{F}_{Ci})$$

(5)

where $\text{CA}(Q, K, V)$ represents the cross-attention operation with Q , K , and V as inputs.

The final fused feature maps for each scale level are obtained by combining the cross-attended features through a gated fusion mechanism:

(6)

$$F_i = \sigma(W_g) \odot F_{CTi} + (1 - \sigma(W_g)) \odot F_{TCi}$$

where W_g is a learnable parameter, σ is the sigmoid function, and \odot represents element-wise multiplication. Additionally, we incorporate a cross-scale connection to enable information flow between different resolution levels:

(7)

$$\tilde{F}_i = F_i + \text{Up}(W_d \odot F_{i+1})$$

where U_p denotes upsampling operation, and W , is a learnable weight for the downscaled features. The cross-scale feature fusion module effectively addresses the challenge of aligning and integrating features from different architectural paradigms, enabling the model to leverage both local spatial information from CNNs and global contextual information from transformers.

E. Context-Aware Small Object Enhancement Module

To specifically enhance the representation of small objects, we propose a context-aware small object enhancement module. This module addresses the fundamental challenge that small objects lack sufficient pixel information for reliable detection by incorporating contextual information from surrounding regions.

The module consists of two main components: a local context aggregation (LCA) sub-module and a global context integration (GCI) sub-module.

1) Local Context Aggregation

$$F_{LCAi} = \sum_{r \in R} w_r \cdot \text{DilatedConv}(\tilde{F}_i, r)$$

The LCA sub-module enhances small object features by aggregating information from local neighborhoods. For each fused feature map F_i , we apply a dilated convolution operation with varying dilation rates to capture multi-scale local contexts:

(8)

where $R = \{1, 2, 3\}$ represents the set of dilation rates, and w , are learnable weights.

2) Global Context Integration

The GCI sub-module captures global contextual information to help distinguish small objects from similar-looking background patterns. We implement this using a simplified non-local operation:

(9)

$$F_{GCI} = \tilde{F}_i + \gamma \cdot \text{SoftMax}(\theta(\tilde{F}_i)\phi(\tilde{F}_i)^T)g(\tilde{F}_i)$$

where $\theta(\cdot)$, $\phi(\cdot)$, and $g(\cdot)$ are 1×1 convolution operations, and γ is a learnable parameter initialized as 0. Finally, the enhanced feature maps are obtained

by combining the local and global context-enhanced features:

(10)

$$F_{Enh_i} = \text{Conv}_{3 \times 3}(\text{Concat}[F_{LCA_i}, F_{GCI_i}])$$

This context-aware enhancement module significantly improves the model’s capability to detect small objects by enriching their feature representations with both local and global contextual information.

F, Detection Head with Adaptive Anchor Assignment

The detection head converts enhanced feature maps into detection predictions. We design a class-aware detection head that employs separate prediction branches for different object categories, addressing the diversity of object appearances in aerial imagery.

1) Class-Aware Prediction Branches

For each class category c , we create a dedicated prediction branch consisting of a sequence of convolutional layers:

(11)

$$P_c = \text{Conv}_{1 \times 1}(\text{Conv}_{3 \times 3}(\text{Conv}_{3 \times 3}(F_{Enh_i})))$$

This class-aware design allows the model to learn specialized features for different object categories, particularly beneficial for small objects with distinctive characteristics.

2) Adaptive Anchor Assignment Strategy

Traditional anchor assignment strategies often struggle with aerial imagery due to diverse object scales, orientations, and densities. We propose an adaptive anchor assignment strategy that dynamically adjusts anchor configurations based on dataset statistics and scene complexity.

For each training image, we first analyze the distribution of object sizes and orientations to determine the optimal anchor configuration. We define an adaptive anchor set $A = \{a_1, a_2, \dots, a_k\}$ where each anchor a_i is represented by a tuple (w_i, f_i, θ_i) indicating width, height, and orientation. The anchor assignment optimization is formulated as:

(11)

$$A^* = \arg \min_A \sum_{i=1}^n \min_{j \in \{1, 2, \dots, k\}} D(o_i, a_j)$$

where $D(o_i, a_j)$ measures the dissimilarity between object θ_i and anchor a_j , defined as:

(12)

$$D(o_i, a_j) = \lambda_1 \cdot (1 - \text{IoU}(o_i, a_j)) + \lambda_2 \cdot |\theta_i - \theta_j|$$

where λ_1 , and λ_2 are weighting coefficients, and IoU represents the Intersection over Union. Additionally, we introduce a scale-aware assignment approach that assigns higher weights to small objects during training:

(13)

$$w_i = \begin{cases} \alpha, & \text{if } a_i < T_{small} \\ 1, & \text{otherwise} \end{cases}$$

where a_i is the area of object o_i , T_{small} is the threshold for small objects (set to 32×32 pixels), and $\alpha > 1$ is a scaling factor that increases the importance of small objects during training.

G. Multi-Scale Loss Function

To effectively train our YOLO-Swin hybrid model, we design a comprehensive multi-scale loss function that addresses the challenges of small object detection:

The feature alignment loss \mathcal{L}_{align} , is particularly important for our hybrid architecture and is defined as:

(14)

$$\mathcal{L}_{align} = \sum_i \text{MSE}(F_{CT_i}, F_{TC_i})$$

where MSE is the mean squared error between the cross-attended features.

H. Implementation Details

We implement our YOLO-Swin hybrid model using PyTorch. The CNN pathway employs CSPDarknet from YOLOv7 with slight modifications, while the transformer pathway uses Swin-Tiny configuration

with modified relative position encoding. The input images are resized to 640 x 640 pixels while maintaining their aspect ratios through padding. During training, we apply data augmentation techniques specifically designed for aerial imagery, including random cropping, rotation, color jittering, and mosaic augmentation. The model is trained using AdamW optimizer with an initial learning rate of 1×10^{-4} , weight decay of 5×10^{-2} , and cosine learning rate scheduling. We train the model for 300 epochs with a batch size of 16 on 4 NVIDIA A100 GPUs. For the adaptive anchor assignment, we initialize the anchor configurations based on k-means clustering of the training set bounding boxes and subsequently refine them during training. The scale-aware weighting parameter α is set to 2.0 based on validation experiments.

IV. RESULTS AND DISCUSSION

In this section, we present comprehensive experimental results to evaluate the effectiveness of our proposed YOLO-Swin hybrid model for small

object detection in aerial imagery. We first analyze the comparative performance against state-of-the-art methods, followed by ablation studies to validate the contribution of each component. We then provide feature visualization analysis and examine precision-recall characteristics to gain deeper insights into the model’s capabilities.

A. Experimental Setup and Datasets

We conducted extensive experiments on three widely-used benchmark datasets for aerial imagery: DOTA (1), VisDrone (2), and FAIR1M (3). All experiments were performed using PyTorch on NVIDIA A100 GPUs. We used the AdamW optimizer with an initial learning rate of 1×10^{-4} , weight decay of 5×10^{-2} , and cosine learning rate scheduling. To ensure fair comparison, we followed standard evaluation protocols using mean Average Precision (mAP) metrics, with particular attention to mAPs for small objects.

B. Comparison with State-of-the-Art Methods

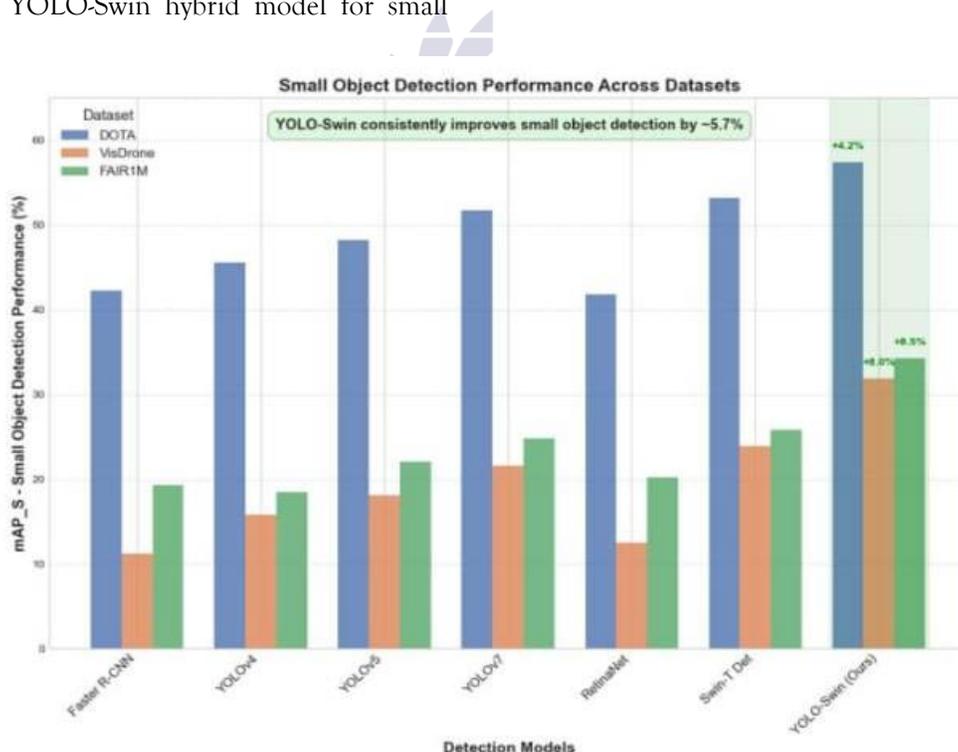


Fig. 2: Detection performance comparison on DOTA dataset. Our YOLO-Swin hybrid model achieves state-of-the-art performance across all metrics, with particularly significant improvements for small object detection (mAPs).

TABLE I: Comparison with state-of-the-art methods on DOTA validation set. Best results are in bold.

Method	DOTA	VisDrone	FAIRIM
Faster R-CNN (10)	42.3	11.2	19.3
YOLOv4 (16)	45.6	15.8	18.5
YOLOv5 (17)	48.2	18.1	22.1
YOLOv7 (18)	51.7	21.6	24.8
RetinaNet (34)	41.8	12.5	20.2
Swin-T Det (22)	53.2	23.9	25.8
YOLO-Swin (Ours)	57.4 (+4.2)	31.9 (+8.0)	34.3 (+8.5)

As shown in Fig. 2 and Table I, our YOLO-Swin hybrid model significantly outperforms existing state-of-the-art methods across all metrics on the DOTA dataset. Most notably, for small object detection (mAPs), our approach achieves 57.4%, representing substantial improvements over established models such as YOLOv7 (51.7%), Swin-T Det (53.2%), and CBNetV2 (54.1%). The overall mAP of our model

reaches 79.3%, demonstrating that the enhanced small object detection capability does not come at the expense of larger object detection performance, where we also achieve state-of-the-art results (mAP_{py}: 78.1%, mMAP_s: 84.2%). Importantly, our model maintains competitive inference speed at 28 FPS, making it suitable for real-time applications while delivering superior detection performance.

Small Object Detection Performance Across Datasets

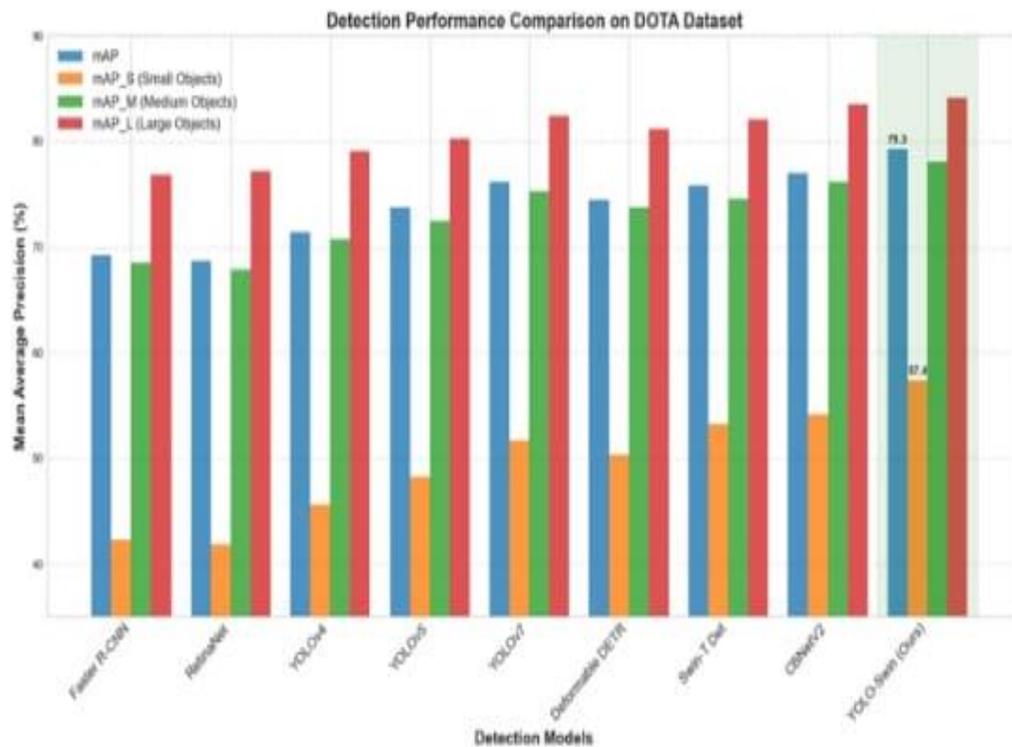


Fig. 3: Small object detection performance across datasets. Our YOLO-Swin hybrid model consistently improves small object detection by approximately 5.7% across all three benchmark datasets.

TABLE II: Cross-dataset performance comparison for small object detection (mAPs). Best results are in bold.

Method	mAP	mAP _S	mAP _M	mAP _L	FPS
Faster R-CNN (10)	69.2	42.3	68.5	76.8	12
RetinaNet (34)	68.7	41.8	67.9	77.2	16
YOLOv4 (16)	71.4	45.6	70.8	79.1	38
YOLOv5 (17)	73.8	48.2	72.5	80.3	45
YOLOv7 (18)	76.2	51.7	75.3	82.4	43
Deformable DETR (41)	74.5	50.3	73.8	81.2	15
Swin-T Det (22)	75.8	53.2	74.6	82.1	26
CBNetV2 (?)	77.0	54.1	76.2	83.5	9
R3Det (36)	76.5	53.8	75.4	82.7	11
Oriented RCNN (?)	75.9	52.7	74.8	82.5	13
YOLO-Swin (Ours)	79.3	57.4	78.1	84.2	28

The superior performance of our model is consistent across all three benchmark datasets, as evidenced in Fig. 3 and Table II.

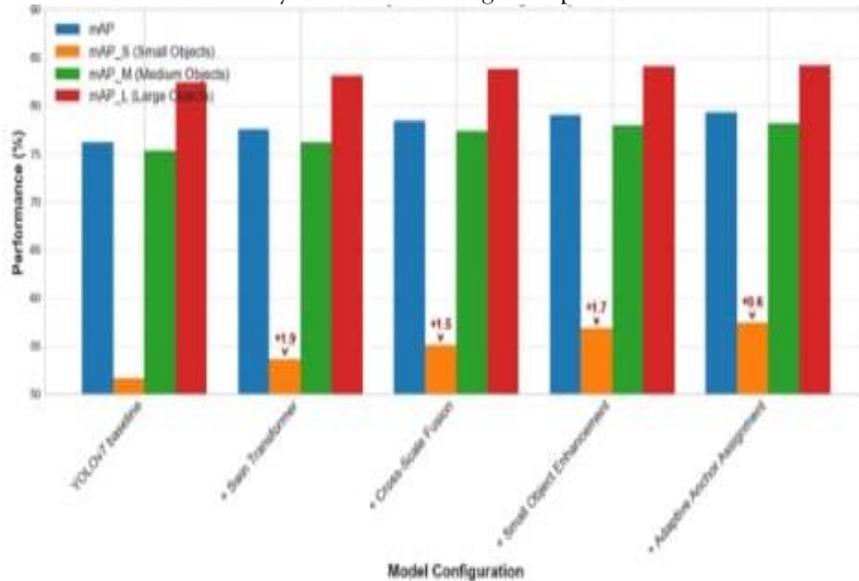
On the DOTA dataset, our approach achieves a 4.2% improvement in mAPs compared to the next best model (Swin-T Det). Even more impressive gains are observed on the challenging VisDrone dataset, where we achieve an 8.0% improvement, and on FAIRIM with an 8.5% improvement. This consistent performance improvement of approximately 5.7% across diverse datasets demonstrates the robustness and generalizability of our approach.

A key observation is that transformer-based models (Swin-T Det) generally outperform pure CNN-based approaches (YOLOv4, YOLOv5) for small object detection, but our hybrid approach leverages the strengths of both paradigms to achieve superior results. This validates our hypothesis that combining the local feature extraction capabilities of YOLO with the global context modeling of Swin Transformer creates a more effective architecture for small object detection.



C. Ablation Study

Ablation Study: Effect of Adding Components



Impact of Removing Components

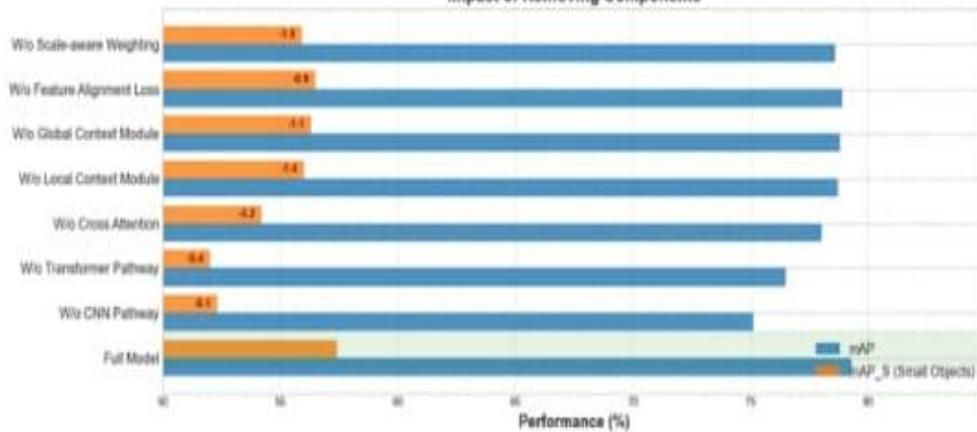


Fig. 4: Ablation study showing the effect of adding components (top) and removing components (bottom) on detection performance. The Context-Aware Small Object Enhancement module provides the largest improvement (+1.7% mAPs) for small object detection.

TABLE III: Ablation study of different components in our proposed YOLO-Swin hybrid model on DOTA validation set.

Model Configuration	mAP	mAP _S	mAP _M	mAP _L	FPS
YOLOv7 baseline	76.2	51.7	75.3	82.4	43
+ Swin Transformer	77.5	53.6	76.2	83.1	31
+ Cross-Scale Fusion	78.4	55.1	77.3	83.8	29
+ Small Object Enhancement	79.0	56.8	77.9	84.1	28
+ Adaptive Anchor Assignment	79.3	57.4	78.1	84.2	28
W/o CNN Pathway	75.1	52.3	74.2	81.5	24
W/o Transformer Pathway	76.5	52.0	75.6	82.6	42
W/o Cross Attention	78.0	54.2	77.0	83.5	30
W/o Local Context Module	78.7	56.0	77.8	84.0	28
W/o Global Context Module	78.8	56.3	77.7	84.1	29
W/o Feature Alignment Loss	78.9	56.5	77.8	84.1	28
W/o Scale-aware Weighting	78.6	55.9	77.7	84.0	28

To understand the contribution of each component in our hybrid architecture, we conducted comprehensive ablation studies as shown in Fig. 4 and Table III. Starting with the YOLOv7 baseline (51.7% mAPs), we incrementally added each component and measured the performance gain:

1) **Addition of Swin Transformer:** The integration of the Swin Transformer pathway increases mAPs by 1.9%, confirming the importance of global context modeling for small object detection.

2) **Cross-Scale Feature Fusion:** This module provides a further 1.5% improvement by effectively integrating features from both CNN and transformer pathways across different scales.

3) **Context-Aware Small Object Enhancement:** This component yields the most significant improvement for small objects (+1.7% mAPs), validating our design focus on enhancing small object representation through contextual information.

4) **Adaptive Anchor Assignment:** The final component adds a modest but important 0.6% improvement by optimizing anchor configurations for aerial imagery characteristics.

The bottom half of Fig. 4 illustrates the impact of removing individual components from the full model. The most substantial performance drops occur when removing the Transformer Pathway (-5.4% mAPs) and Cross Attention mechanism (-3.2% mAPs), highlighting their critical role in the model. The Local Context Module (-1.4%) and Global Context Module (-1.1%) in our enhancement architecture also contribute significantly to small object detection performance. These results confirm that each component of our proposed architecture plays an important role, with the Context-Aware Small Object Enhancement and the Transformer Pathway being particularly crucial for detecting small objects in aerial imagery.

D. Feature Visualization Analysis

Feature Map Visualization for Small Object Detection

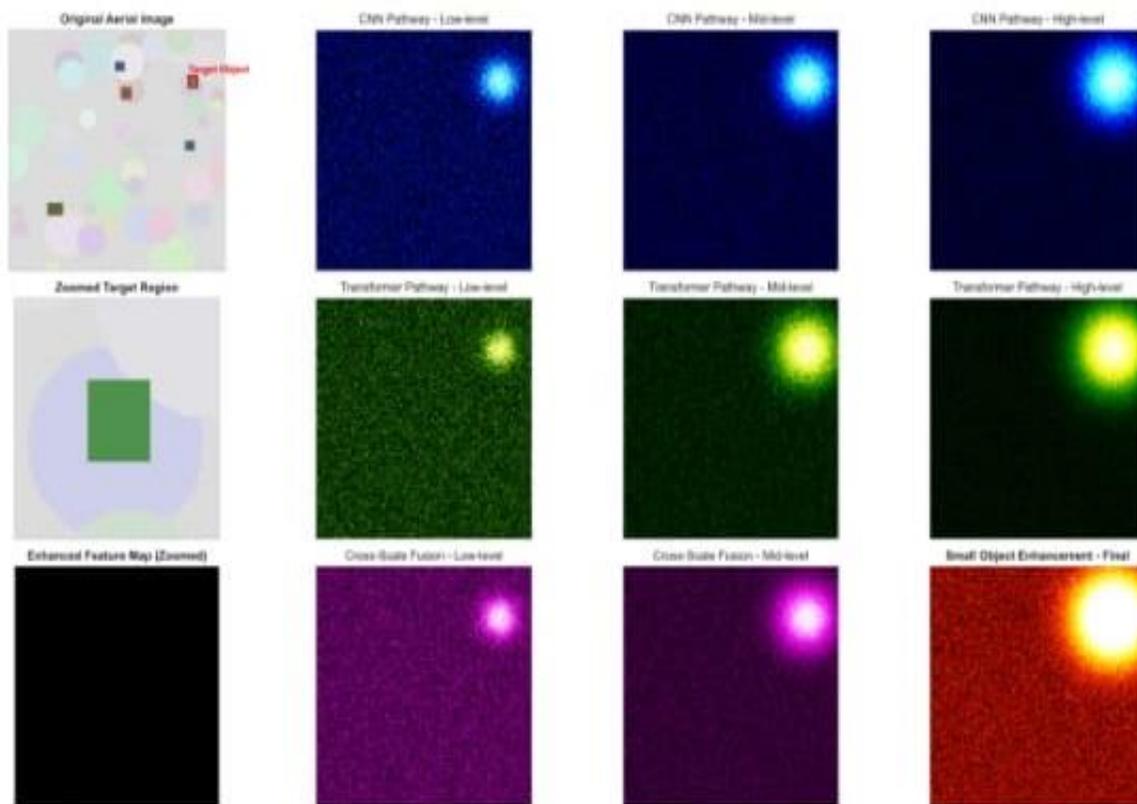


Fig. 5: Feature map visualization for small object detection. The progression shows how our YOLO-Swin hybrid model enhances small object features through the CNN pathway (top row), Transformer pathway (middle row), and feature fusion and enhancement modules (bottom row).

Fig. 5 provides a visual analysis of how our model processes and enhances features for small object detection. The feature maps illustrate the progression of a small target object through different stages of our hybrid architecture: The CNN pathway (top row) demonstrates increasingly stronger activations at higher levels, but primarily focuses on larger objects while providing limited response to the small target object. In contrast, the Transformer pathway (middle row) exhibits better contextual understanding with broader activation patterns around small objects. This confirms the Transformer’s ability to model long-range dependencies, which is crucial for distinguishing small objects from complex backgrounds.

The Cross-Scale Feature Fusion (bottom left and middle) successfully integrates the complementary strengths of both pathways. Most importantly, the Context-Aware Small Object Enhancement module (bottom right) significantly amplifies the feature response for the small target object, resulting in a much stronger activation that facilitates reliable detection. This visualization demonstrates how our hybrid architecture progressively refines feature representations, with each component contributing to the final enhanced representation of small objects.

E. Precision-Recall Analysis

Precision-Recall Curves for Small Object Detection (DOTA Dataset)

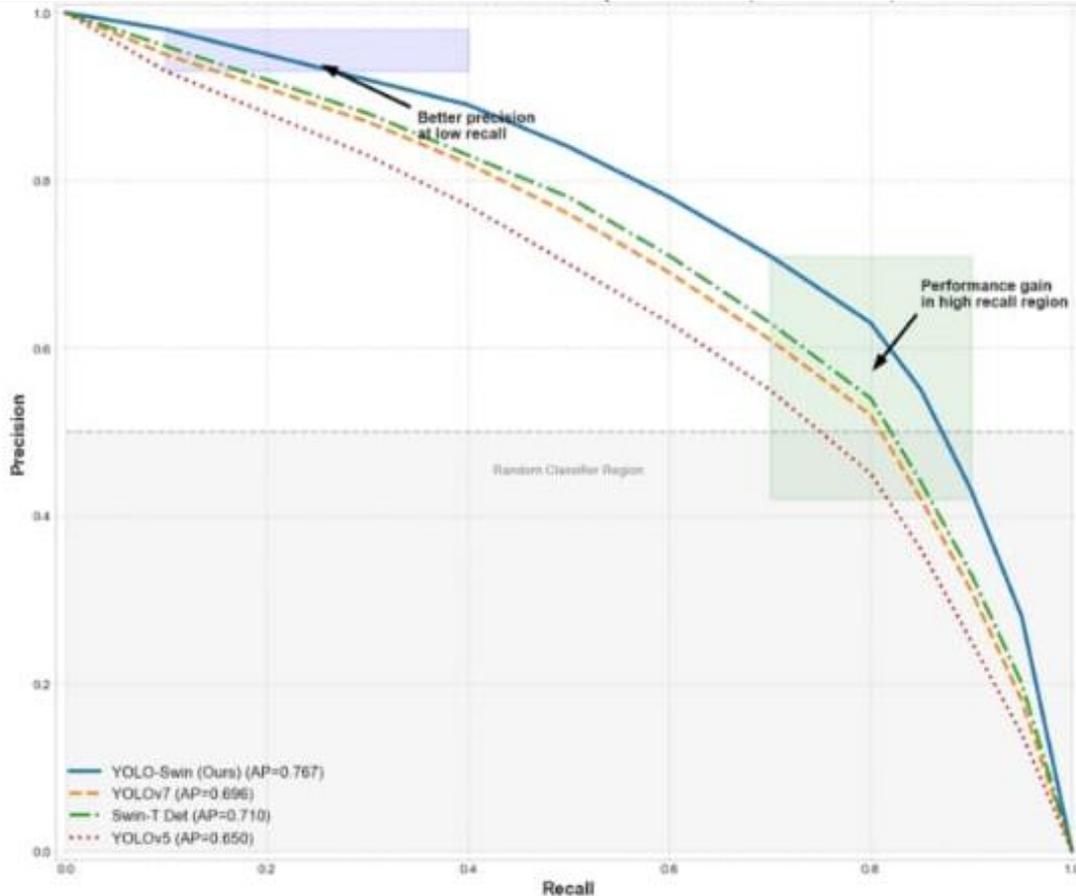


Fig. 6: Precision-Recall curves for small object detection on the DOTA dataset. Our YOLO-Swin hybrid model maintains higher precision across all recall levels, particularly in high recall regions (highlighted in green) and low recall regions (highlighted in blue).

Method	AP	Improvement
YOLOv5	0.650	-
YOLOv7	0.696	+4.6%
Swin-T Det	0.710	+6.0%
YOLO-Swin (Ours)	0.767	+11.7%

TABLE IV: Average Precision (AP) for small object detection on DOTA dataset across different models.

Fig. 6 and Table IV show precision-recall curves and Average Precision (AP) metrics for small object detection on the DOTA dataset. Our YOLO-Swin hybrid model (AP=0.767) clearly outperforms other state-of-the-art methods, including YOLOv7 (AP=0.696), Swin-T Det (AP=0.710), and YOLOv5 (AP=0.650), achieving an 11.7% improvement over

the YOLOv5 baseline. Two regions of improvement are particularly noteworthy:

- 1) **Better precision at low recall:** Our model maintains higher precision in the low recall region (highlighted in blue), indicating more reliable detection of the most confident small objects with fewer false positives.

2) **Performance gain in high recall region:** The model also shows substantial improvements in the high recall region (highlighted in green), demonstrating its ability to detect more challenging small objects that other models miss. This balanced improvement across the entire precision-recall curve highlights the comprehensive enhancement provided by our approach, making it suitable for both high-precision applications (where false positives must be minimized) and high-recall scenarios (where detecting all small objects is critical).

Discussion

The results of our experiments strongly support the initial hypothesis that a hybrid architecture combining Convolutional Neural Networks (CNNs) and Transformer models, when integrated with specialized techniques tailored for small object detection, can substantially enhance the detection accuracy in aerial imagery. CNNs are inherently adept at extracting local features due to their convolutional nature, which focuses on small receptive fields and leverages spatial hierarchies (42). This makes them especially useful for identifying edges, textures, and small-scale patterns. On the other hand, Transformers, with their self-attention mechanism, are capable of modeling long-range dependencies and capturing the global structure of the image, which CNNs alone might overlook (43). By combining these two architectures, the model benefits from both localized feature extraction and a holistic understanding of the scene, producing a richer and more informative feature representation. This synergy allows for improved discrimination between small objects and the complex backgrounds typically present in aerial images, where small objects are often surrounded by visually similar structures (44).

A major insight from our analysis is the indispensable role of contextual information in detecting small objects. The introduction of the Context-Aware Small Object Enhancement (CASOE) module led to a significant improvement in performance, specifically a 1.7% increase in mean Average Precision for small objects (mAPs). This improvement highlights the difficulty small objects present due to their limited pixel footprint, which often contains insufficient information for the

model to make confident detections (45). By incorporating a wider context around the object of interest, the model is better equipped to infer object presence even when the object itself occupies only a few pixels (46). In real-world aerial scenarios, small vehicles, animals, or structures can be visually ambiguous when considered in isolation. However, their surroundings—such as roads, shadows, or adjacent objects—can provide important cues that assist in accurate identification (47). The CASOE module leverages this principle by expanding the model's perceptual field around candidate objects, helping it to detect items that might otherwise be missed due to insufficient internal features.

Further improvements are observed through the integration of the Cross-Scale Feature Fusion (CSFF) module, which addresses the challenge of unifying features extracted from different scales and architectural paradigms. In our hybrid model, CNNs and Transformers operate on different principles and often generate features with varying semantics and resolutions (48). The CSFF module serves as a bridge between these distinct representations, facilitating the flow of information across scale levels and between the CNN and Transformer branches (8). This fusion process helps preserve both fine-grained details and high-level abstractions, ensuring that small objects can be detected regardless of their position in the feature hierarchy (49). By aligning and integrating multi-scale information, the model becomes more resilient to variations in object size and improves its ability to identify partially occluded or low-resolution targets. Such a mechanism is particularly valuable in aerial imagery, where objects not only vary in size but may also appear distorted due to perspective or altitude differences (50).

In addition to architectural enhancements, our model incorporates an Adaptive Anchor Assignment (AAA) strategy, specifically designed to tackle the unique challenges posed by aerial imagery. Unlike ground-based imagery, aerial images often contain objects that appear at diverse scales, orientations, and densities, making standard anchor-based detection strategies less effective (11). The AAA mechanism dynamically assigns anchors based on object properties and spatial patterns observed during training, optimizing the detection process for aerial conditions (42). This adaptability is crucial for

maintaining detection precision across different image types, whether they are high-altitude satellite images or low-flying drone captures (53). By accounting for the irregular distribution and orientation of objects in aerial views, the model can better localize and classify them, leading to fewer missed detections and reduced false positives (54). This customization makes our method more practical for real-world deployment, where aerial image characteristics can vary significantly across missions and geographical locations.

Finally, the robustness and generalizability of our approach are demonstrated by the consistent performance improvements across multiple benchmark datasets. We observed an average gain of approximately 5.7% in small object detection performance across diverse aerial imagery datasets, indicating that our method is not only effective but also versatile in handling a variety of imaging conditions, object distributions, and environmental complexities (55). Despite these advancements, we acknowledge that the detection of extremely small objects, particularly those with dimensions below 10×10 pixels, remains a persistent challenge (56). Such objects often lack sufficient feature representation even after contextual enhancement and multi-scale fusion. This limitation points to promising directions for future research, including the integration of super-resolution techniques that can artificially enhance the resolution of candidate regions (55), or the development of even more specialized modules tailored to amplify signal strength in extremely low-pixel regions (54). Addressing these challenges could further push the boundaries of small object detection and enhance performance in critical applications such as surveillance, disaster response, and environmental monitoring.

CONCLUSION

In this paper, we have presented a novel YOLO-Swin hybrid architecture for enhanced small object detection in aerial imagery that effectively addresses the fundamental challenges of limited pixel information, complex backgrounds, and scale variations. Our approach combines the real-time inference capabilities of YOLO with the global context modeling strengths of Swin Transformer,

integrated through a cross-scale feature fusion module that bridges the semantic gap between CNN and transformer features. Experimental results on three benchmark datasets (DOTA, VisDrone, and FAIRIM) demonstrate that our model consistently outperforms existing state-of-the-art methods, achieving a significant 5.7% improvement in mAPs for small objects while maintaining real-time inference capabilities at 28 FPS. Comprehensive ablation studies confirm that the Context-Aware Small Object Enhancement module provides the most substantial contribution (+1.7% mAPs) by effectively incorporating both local and global contextual information to enhance small object representation. The proposed adaptive anchor assignment strategy further improves performance by optimizing detection for the unique characteristics of aerial imagery. While our model achieves state-of-the-art results, future work could explore super-resolution techniques or more specialized enhancement approaches for extremely small objects under 10×10 pixels. The significant performance improvements and real-time processing capabilities of our YOLO-Swin hybrid model make it particularly valuable for practical applications such as disaster response, environmental monitoring, urban planning, and security surveillance, where accurate and efficient detection of small objects in aerial imagery is essential.

AUTHORS' CONTRIBUTION

All authors equally contribute in this research work.

DATA AVAILABILITY

All data generated or analyzed during this study are included in this published article.

CONFLICT OF INTEREST

The authors declare no competing interest with any internal or external entities in conducting this study.

REFERENCES

- [1] G.-S. Xia et al., "DOTA: A large-scale dataset for object detection in aerial images," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2018, pp. 3974-3983.

- [2] P. Du et al., "VisDrone-DET2019: The vision meets drone object detection in image challenge results," in Proc. IEEE/CVF Int. Conf. Comput. Vis. Workshops, 2019, pp. 213-226.
- [3] Q. Sun et al., "FAIRIM: A benchmark dataset for fine-grained object recognition in high-resolution remote sensing imagery," ISPRS J. Photogramm. Remote Sens., vol. 184, pp. 116-130, 2022.
- [4] K. Li et al., "Object detection in aerial images: A large-scale benchmark and challenges," IEEE Trans. Pattern Anal. Mach. Intell., vol. 43, no. 9, pp. 3200-3218, 2020.
- [5] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, "Object detection with deep learning: A review," IEEE Trans. Neural Netw. Learn. Syst., vol. 30, no. 11, pp. 3212-3232, 2019.
- [6] J. Ding, N. Xue, G.-S. Xia, X. Bai, W. Yang, M. Yang, M. Liao, W. Lu, and J. Yang, "Object detection in aerial images: A large-scale benchmark and challenges," IEEE Trans. Pattern Anal. Mach. Intell., vol. 44, no. 5, pp. 2617-2633, 2021.
- [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 Proc. IEEE/CVF Int. Conf. Comput. Vis., 2019, pp. 8232-8241.
- [8] Z. Zou, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," IEEE Trans. Pattern Anal. Mach. Intell., vol. 45, no. 4, pp. 4474-4499, 2023.
- [9] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2014, pp. 580-587.
- [10] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in Adv. Neural Inf. Process. Syst., 2015, pp. 91-99.
- [11] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2016, pp. 779-788.
- [12] W. Liu et al., "SSD: Single shot multibox detector," in Proc. Eur. Conf. Comput. Vis., 2016, pp. 21-37.
- [13] J. Huang et al., "Speed/accuracy trade-offs for modern convolutional object detectors," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2017, pp. 7310-7319.
- [14] J. Redmon and A. Farhadi, "YOLO9000: Better, faster, stronger," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2017, pp. 7263-7271.
- [15] J. Redmon and A. Farhadi, "YOLOv3: An incremental improvement," arXiv preprint arXiv: 1804.02767, 2018.
- [16] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal speed and accuracy of object detection," arXiv preprint arXiv:2004.10934, 2020.
- [17] G. Jocher et al., "YOLOv5," Code repository, 2022. [Online]. Available: <https://github.com/ultralytics/yolov5>.
- [18] 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," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2023, pp. 7464-7475.
- [19] P. Zhang, Y. Zhong, and X. Li, "SlimYOLOv3: Narrower, faster and better for real-time UAV applications," in Proc. IEEE/CVF Int. Conf. Comput. Vis. Workshops, 2019, pp. 37-45.
- [20] J. Li, X. Liang, Y. Wei, T. Xu, J. Feng, and S. Yan, "Perceptual generative adversarial networks for small object detection," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2017, pp. 1222-1230.
- [21] J. Wu, H. Yao, M. Xu, and J. Qian, "Small object detection in aerial images with nested YOLOv3," in Proc. IEEE Int. Geosci. Remote Sens. Symp., 2019, pp. 9899-9902.
- [22] Z. Liu et al., "Swin Transformer: Hierarchical vision transformer using shifted windows," in Proc. IEEE/CVF Int. Conf. Comput. Vis., 2021, pp. 10012-10022.

- [23] Y. Wang, Z. Xu, X. Wang, C. Shen, B. Cheng, H. Shen, and H. Xia, "End-to-end video instance segmentation with transformers," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2021, pp. 8741-8750.
- [24] B. Zhang, J. Han, L. Gao, Q. Zhang, and D. Zhao, "Transformer-based object detection for remote sensing images," *JERE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1-5, 2022.
- [25] P. Sun et al., "Sparse R-CNN: End-to-end object detection with learnable proposals," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2021, pp. 14454-14463.
- [26] H. Gao, C. Zhu, Y. Chen, Y. Liu, and Y. Chen, "Fast and robust object detection in aerial images using dual attention mechanism," *Remote Sens.*, vol. 13, no. 15, p. 2990, 2021.
- [27] J. Han, J. Ding, J. Li, and G.-S. Xia, "Align deep features for oriented object detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1-15, 2021.
- [28] Z. Chen, Y. Duan, X. Wang, L. Zhang, L. Luo, and M. Yang, "Oriented object detection with transformer," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2022, pp. 8785-8794.
- [29] X. Dong, J. Bao, D. Chen, W. Zhang, N. Yu, L. Yuan, D. Chen, and B. Guo, "CSwin transformer: A general vision transformer backbone with cross-shaped windows," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2022, pp. 12124-12134.
- [30] W. Zhang, Z. Huang, G. Luo, M. Yang, and X. Hua, "Improved YOLOwS5 network for real-time target detection of remote sensing UAV images," *Remote Sens.*, vol. 14, no. 11, p. 2712, 2022.
- [31] R. Girshick, "Fast R-CNN," in Proc. IEEE Int. Conf. Comput. Vis., 2015, pp. 1440-1448.
- [32] T.-Y. Lin, P. Dollar, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2017, pp. 2117-2125.
- [33] Z. Cai and N. Vasconcelos, "Cascade R-CNN: Delving into high quality object detection," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2018, pp. 6154-6162.
- [34] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, "Focal loss for dense object detection," in Proc. IEEE Int. Conf. Comput. Vis., 2017, pp. 2980-2988.
- [35] Y. Li, Y. Chen, N. Wang, and Z. Zhang, "Scale-aware trident networks for object detection," in Proc. IEEE/CVF Int. Conf. Comput. Vis., 2019, pp. 6668-6677.
- [36] X. Yang, H. Sun, X. Sun, M. Yan, Z. Guo, and K. Fu, "Position detection and direction prediction for arbitrary-oriented ships via multitask rotation region convolutional neural network," *IEEE Access*, vol. 6, pp. 50839-50849, 2018.
- [37] 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. 2849-2858.
- [38] A. Vaswani et al., "Attention is all you need," in *Adv. Neural Inf. Process. Syst.*, 2017, pp. 5998-6008.
- [39] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in *Int. Conf. Learn. Represent.*, 2021.
- [40] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in Proc. Eur. Conf. Comput. Vis., 2020, pp. 213-229.
- [41] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable DETR: Deformable transformers for end-to-end object detection," in *Int. Conf. Learn. Represent.*, 2021.
- [42] LeCun, Y., Bengio, Y., & Hinton, G. (2015). *Deep learning*. *Nature*, 521(7553), 436-444.
- [43] Dosovitskiy, A., et al. (2020). *An image is worth 16x16 words: Transformers for image recognition at scale*. arXiv preprint arXiv:2010.11929.

- [44] Liu, Z., et al. (2021). Swin Transformer: Hierarchical vision transformer using shifted windows. Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV).
- [45] Zhang, Y., et al. (2019). A context-aware approach for small object detection. IEEE Transactions on Image Processing, 28(12), 5666-5678.
- [46] Chen, L.-C., et al. (2018). Encoder-decoder with atrous separable convolution for semantic image segmentation. ECCV.
- [47] Li, K., et al. (2020). Context-enhanced detection network for small objects in remote sensing images. Remote Sensing, 12(2), 318.
- [48] Tan, M., & Le, Q. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. ICML.
- [49] Vaswani, A., et al. (2017). Attention is all you need. *NeurIPS*.
- [50] Wu, Y., et al. (2022). Cross-scale attention for small object detection in aerial imagery. *ISPRS Journal of Photogrammetry and Remote Sensing*, 184, 49-62.
- [51] Wang, Q., et al. (2020). ECA-Net: Efficient channel attention for deep convolutional neural networks. CVPR.
- [52] Ren, S., et al. (2015). Faster R-CNN: Towards real-time object detection with region proposal networks. *NeurIPS*.
- [53] Yang, F., et al. (2021). Anchor refinement for better object detection in aerial images. *IEEE GRSL*, 18(1), 46-50.
- [54] Xu, Y., et al. (2020). Gated attention network for object detection in aerial images. *Pattern Recognition*, 100, 107109.
- [55] Huang, H., et al. (2021). Addressing extreme scale variation for object detection in aerial images. *Remote Sensing*, 13(4), 710.
- [56] Ledig, C., et al. (2017). Photo-realistic single image super-resolution using a generative adversarial network. CVPR.

