**Research Article**

# SwinRODNet Swin Transformer-Based Remote Sensing Object Detection Network

Balamanikandan A[1], Jayakumar S[2], Srinanda Ganesh T[3], Sukanya M[4], B. Venkataramanaiah[5], Arunraja A[6]

[1]Department of Electronics and Communication Engineering, MohanBbabu University, Tirupathi, India.

[2]Department of Electronics and Communication Engineering, Sri Sairam College of Engineering, Anekal, Bengaluru, India.

[3]Department of Electrical and Electronics Enginrring, St. Joseph's College of Engineering, Chennai, India.

[4]Department of Electrical and Electronics Enginrring, Adhiyamaan College of Engineering, Hosur, India.

[5]Department of Electronics and Communication Engineering, Vel Tech Rangarajan Dr. Sagunthala R&D Institute of Science and Technology, Chennai, India.

[6]Department of Electronic and Communications Engineering, Christ University, Bengaluru, India.

| ARTICLE INFO | ABSTRACT |
|---|---|
| | Detecting objects via remote sensing in natural settings is extremely difficult, especially when dealing with small targets and complicated backdrops. To improve feature extraction and detection accuracy, this proposal presents the RAST-YOLO method, which combines the Region Attention (RA) mechanism with the dual Transformer backbone. To solve the multi-scale issue and improve small item detection, the C3D module is used to combine deep and shallow semantic information. The algorithm's exceptional resilience, accuracy, and efficiency are demonstrated by extensive testing on the DIOR and TGRS-HRRSD datasets. Compared to baseline networks, RAST-YOLO shows a notable improvement in mean average precision (mAP) on both datasets. Additionally, research using methods like YOLOv5x6 and YOLOv8 reveals promising results, with YOLOv5x6 achieving a significant mAP enhancement of over 0.80%, highlighting its suitability for advanced remote sensing object detection applications. The proposed algorithm's ability to handle complex backgrounds and small-scale targets effectively makes it a valuable tool for various remote sensing applications, including environmental monitoring, resource exploration, and intelligent navigation.<br><br>**Keywords:** Remote Sensing, Deep Learning, Transformer Backbone, Region Attention Mechanism (RA), Mean Average Precision (mAP). |

## I.INTRODUCTION

Object detection in remote sensing images is essential to comprehend aerial and satellite photography utilised for various applications, including resource discovery, intelligent navigation, environmental monitoring, and target tracking. Finding specified targets and their geographical positions is the primary goal of remote sensing target detection. Many high-resolution and high-quality datasets have been created for processing remote-sensing photos due to the quick developments in UAVs and aeronautical technology [1][2]. Despite these advancements, there are still many obstacles to overcome in remote sensing object detection, such as small data sizes and items' similar appearances across categories. When compared to deep learning techniques, traditional object recognition methods—which include feature extraction, feature modification, and classifier prediction—are criticised for their lack of generalisation and resilience [3][4].

The RAST-YOLO algorithm has been extensively tested on the DIOR and TGRS-HRRSD datasets, demonstrating superior robustness, accuracy, and efficiency. Compared to baseline networks, RAST-YOLO shows a notable improvement in mean average precision (mAP) on both datasets. Additionally, research using methods like YOLOv5x6 and YOLOv8 reveals promising results, with YOLOv5x6 achieving a significant mAP enhancement of over 0.80%, highlighting its suitability for advanced remote sensing object detection applications [5]. The proposed method effectively handles complex backgrounds and small-scale targets, making it highly useful for applications such as resource exploration, environmental monitoring, and intelligent navigation.

Accurately detecting objects in remote sensing images presents significant challenges due to the presence of small targets and complex backgrounds. Traditional object recognition methods, which rely on feature extraction, modification, and classifier prediction, often fall short in terms of generalization and resilience when compared to deep learning techniques. Existing solutions struggle with issues such as small data sizes and the similar appearances of objects across categories, leading to decreased precision in object detection. This research aims to address these challenges by introducing the RAST-YOLO algorithm, which integrates the Region Attention (RA) mechanism with the Swin Transformer backbone. The proposed method seeks to enhance feature extraction and detection accuracy, particularly for small-scale targets and complex backdrops [6]. Despite the development of high-resolution datasets like DIOR and TGRS-HRRSD, there remains a need for more robust and efficient algorithms capable of overcoming the limitations of traditional methods. The RAST-YOLO algorithm aims to fill this gap by leveraging advanced attention mechanisms and dual Transformer backbones, thereby improving mean average precision (mAP) and overall detection performance in remote sensing applications such as environmental monitoring, resource exploration, and intelligent navigation [7].

By improving the interaction range of feature information and reducing the effect of complicated backdrops on detection accuracy, the RAST-YOLO algorithm offers a novel and effective approach for remote sensing object detection. The integration of deep and shallow semantic information through the C3D module enhances small-item detection accuracy. Comprehensive testing on high-quality, high-resolution datasets like TGRS-HRRSD and DIOR has validated RAST-YOLO's exceptional performance. The utilization of software tools like YOLO, TensorFlow, NumPy, Pandas, Matplotlib, and Scikit-learn further supports the robustness of the proposed method.

## II.LITERATURE REVIEW AND PROBLEM IDENTIFICATION

B. Yan, et al. focused on developing radar systems to improve localization accuracy, trajectory detection, and area coverage, especially in cluttered environments with weak targets. Their approach, involving spatiotemporal clutter maps and a track-before-detect strategy, significantly enhanced multiple target tracking. However, these methods primarily focus on radar data, highlighting a gap in remote sensing object detection, which our RAST-YOLO algorithm aims to address [1].

H. Lee et al. proposed an underwater object localization method using DC electric field templates for real-time tracking, demonstrating accuracy in noisy environments. While effective underwater, this approach underscores the need for advanced techniques in aerial and satellite image processing. Our method leverages deep learning and advanced attention mechanisms to enhance feature extraction and detection accuracy in complex remote sensing scenarios [2].

M. Ilyas et al. Proposed Remote inspections are made possible by robotic construction automation solutions powered by AI and sophisticated mechatronics. Conventional approaches are subjective and time-consuming. With the use of vision sensors and clever algorithms, the suggested robotic system assists supervisors in remotely identifying items, spotting flaws, and producing reports. It outperforms traditional deep learning techniques by using BIM for navigation and a data-driven strategy for real-time item recognition. Experiments have shown that this technique improves the accuracy and efficiency of inspections [3].

M. Zurowietz and T. W. Nattkemper. Suggested Computational support is necessary for the timely evaluation of the continuously increasing volume of digital image data used in maritime environmental monitoring and research. Inadequate training data frequently causes problems for contemporary deep learning methods. By using "scale transfer" and improved data augmentation, the Unsupervised Knowledge Transfer (UnKnoT) approach effectively makes use of sparse training data. For object detection, this technique makes advantage of pre-existing training data. Tests on four datasets of annotated marine images showed notable gains in object detection capabilities. This method guarantees that contemporary machine learning may be applied to the monitoring of marine environments [4].

W. -L. Zhao and C. -W. Ngo. They talk about F-SIFT (Flip-Invariant SIFT) and SIFT (Scale-Invariant Feature Transform). Although it lacks flip invariance, SIFT is renowned for its resilience to rotation, scale, and lighting variations. The novel descriptor F-SIFT is flip-tolerant and preserves the characteristics of SIFT. Prior to SIFT computation, it estimates the dominant curl of a local patch and does a geometric normalisation. Compared to SIFT, this improvement increases detection accuracy and lowers computing costs by more than 50%. Across a range of key point detectors, F-SIFT consistently performs better than seven other descriptors when managing flip transformation Ons and characterising symmetric objects [5].

In comparing the discussed literature surveys, several key aspects emerge. The underwater object detection method using grid-based template matching offers real-time tracking accuracy but may struggle in extremely noisy conditions. The robotic inspection system in construction automation enhances safety and efficiency but requires significant initial investment and setup. Unsupervised Knowledge Transfer (UnKnoT) for marine environmental monitoring effectively utilizes limited data, though its reliance on pre-existing datasets can be limiting. The algorithm for multiple radar systems improves tracking in cluttered environments but demands substantial computational resources. Finally, the flip-invariant SIFT descriptor significantly enhances accuracy and reduces computational costs in object detection but may still struggle with complex, real-world transformations. Each method advances specific areas of object detection, yet also presents challenges that must be navigated for optimal application.
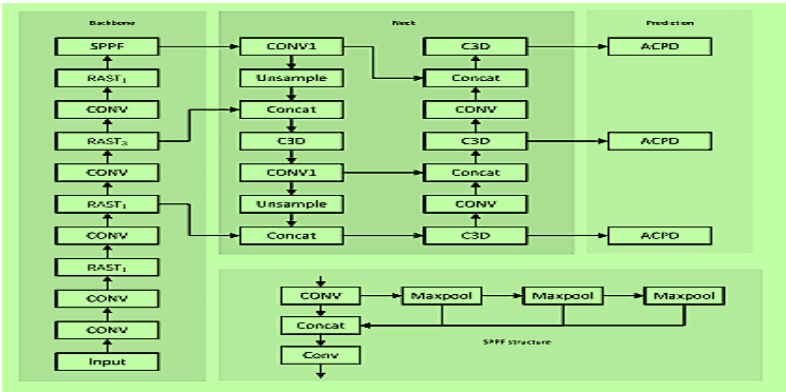
### III.SYSTEM ANALYSIS



Fig .1. System architecture

The RAST-YOLO system architecture fig.1. for remote sensing object detection consists of three main components. The Backbone includes the input layer, two Convolution (CONV) layers, Region Attention (RAST1), another CONV layer, RAST1 again, another CONV layer, RAST2, and one more CONV layer. This sequence enhances feature extraction before passing the output to the Spatial Pyramid Pooling Fast (SPPF) module. The Next part processes the Backbone's output through various stages, including CONV1 layers, up sampling, concatenation, and C3D modules. These steps refine the feature maps for better detection accuracy [8][9]. Finally, the C3D layers in the Next section are connected to three instances of the Adaptive Convolutional Prediction Detector (ACPD) in the Prediction section, enabling precise object detection.

### A.DATA FLOW DIAGRAM (DFD)

The process of building and training a machine learning model involves several steps. First, libraries need to be imported, and the process must be verified. Next, the dataset is imported, and image processing begins. Following this, the pretrained model is loaded, and data augmentation is performed. The model can be built in Colab with various options such as Yolo V5s, RAST YOLO, CNN and YOLO backbone, Yolo V3, Faster CNN, RetinaNet, Yolo V5x6, and Yolo V8. The model is then trained, after which users must sign up and sign in. User input is gathered, and the process ends with the outcome [10]. There is also a decision point for verification to determine whether to proceed or end the process.
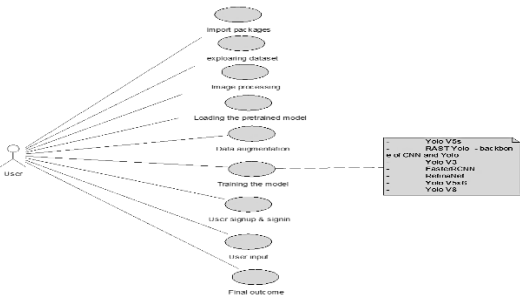
### B.UML



Fig.2. User Interaction Flow for Model Training and Deployment

A User Interaction Flow for Model Training and Deployment fig.2. system begins with the potential selection of specific software libraries, represented by "import packages," and an exploration of the available image dataset used for training the model. The user may then have influence over the image preprocessing steps, labelled "image processing," potentially adjusting parameters or settings. Crucially, the user selects a pre-trained object detection model from a range of options, including various versions of YOLO (You Only Look Once), such as Yolo V5s, RAST Yolo, Yolo V3, Yolo V5x6, and Yolo V8, as well as alternative architectures like FasterRCNN and RetinaNet. These pre-trained models offer a foundation for the system's object detection capabilities [11]. The user might also configure the "data augmentation" process, choosing which techniques to apply to enhance the training data. They then initiate or monitor the "training the model" phase, where the selected model learns to detect objects in images. Before using the core detection functionality, the user performs standard authentication steps, "user signup & signin." Finally, the user provides the image they want to analyze, labelled "user input," and the system processes it, delivering the "final outcome," which comprises the results of the object detection, highlighting the identified objects within the image.
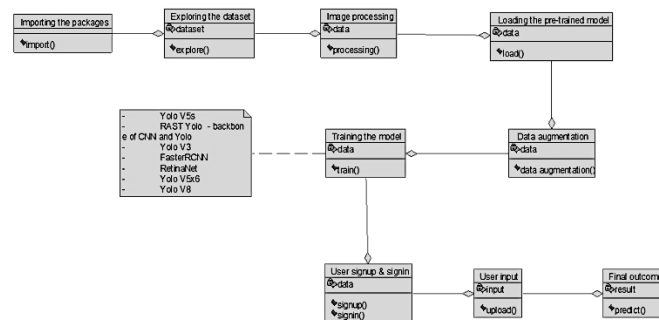
## C.CLASS DIAGRAM



Fig .3. Class Diagram for Object Detection System

The process starts by importing necessary software libraries and then exploring the dataset of images and their annotations fig .3., which are essential for training the object detection model. Following this, the images undergo preprocessing steps like resizing and normalization in the image processing stage the preprocessing steps for training an object detection model involve using pre-trained models like YOLO (V5s, RAST, V3, V5x6, V8), Faster CNN, or RetinaNet. Data augmentation techniques, such as rotations and flips, are essential for enhancing the model's robustness and preventing overfitting [12]. The augmented data is used to train the model, fine-tuning its parameters for accurate object detection. Users interact with the system by creating accounts, signing in, and uploading images to receive detection results with highlighted objects and their classifications. Users interact with the system by creating accounts or signing in. They then provide input, typically by uploading an image, which is passed to the trained model. Finally, the model processes the user's input image and generates the outcome – the object detection results, highlighting the detected objects and their classifications within the image.

This diagram details the interactions between different actors and the system, depicting how users and external systems interact with the RAST-YOLO algorithm.
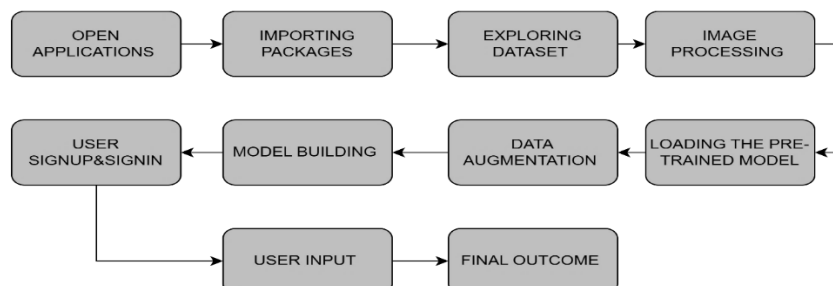
## D. ACTIVITY DIAGRAM



Fig .4. Activity Diagram for RAST-YOLO Algorithm

The Activity Diagram for RAST-YOLO Algorithm fig .4. shows the sequence of activities and decisions in the RAST-YOLO algorithm, highlighting the system workflow. The object detection workflow, as depicted in this activity diagram, begins with the user opening the application and the system importing necessary software packages [13][14]. Following this initialization, the user or system explores the available dataset of images, which are then processed to prepare them for model training. Concurrently, the system either builds a new model from scratch or loads a pre-trained model, leveraging existing knowledge. Once both the image processing and model loading/building are complete, data augmentation techniques are applied to enhance the training dataset. After these preparation steps, the user signs up or enters the application [15]. The user then provides an input, typically an image, that they want the system to analyze. Finally, the system processes the input image using the trained model and generates the final object detection outcome, highlighting the detected objects within the image. The diagram also suggests potential concurrency in the model building/loading and data augmentation stages, as these activities can potentially occur in parallel.

### E. SEQUENCE AND COLLABORATION

In the RAST-YOLO project for remote sensing object detection, combining the concepts of sequence, collaboration, and component provides a comprehensive view of the system's interactions and structure. The sequence illustrates the time-ordered interactions between different objects within the system, showing the sequence of messages exchanged between objects to perform specific tasks, such as importing packages, exploring the dataset, processing images, loading the pre-trained model, data augmentation, building the model, user signup and sign in, user input, and producing the outcome [16]. The collaboration groups these interactions by highlighting the relationships and communication between objects, identifying all possible interactions that each object has with other objects, and emphasizing how they work together to achieve the system's goals. The component represents the high-level parts that make up the system and their interactions, showing the modular structure of the system, and how various parts such as data processing modules, machine learning models, and user interfaces are interrelated [17][18]. This integrated approach ensures a thorough analysis of the system's functionality and performance, highlighting how different parts of the system collaborate to achieve the desired outcomes. Feel free to insert the respective UML here to visualize these interactions effectively.

### IV. METHODOLOGY AND IMPLEMENTATION

To assess the performance of the proposed RAST-YOLO algorithm, we use precision, recall and F1 score metrics. These metrics are defined as follows:

*Precision:*

$$Precision = \frac{True\ Positives\ (TP)}{True\ Positives\ (TP) + False\ Positives\ (FP)}$$

*Recall:*

$$Recall = \frac{True\ Positives\ (TP)}{True\ Positives\ (TP) + False\ Negatives\ (FN}$$

*F1 score:*

$$F1\ Score = 2\ \times \frac{Precision \times Recall}{Precision + Recall}$$

Precision indicates the proportion of correctly detected objects out of all detected objects, while recall indicates the proportion of correctly detected objects out of all relevant objects in the dataset. The F1 Score balances precision and recall, providing a comprehensive measure of accuracy.

### A. STEP-BY-STEP IMPLEMENTATION

1. Import Libraries and Packages: Import necessary tools for data processing, model construction, and evaluation, such as TensorFlow, NumPy, Pandas, and Matplotlib.

2. Load and Explore the Dataset: Load the dataset and perform initial exploration to understand its size, shape, and distribution. Visualize sample images to identify any anomalies.

3. Image Processing: Preprocess images by resizing, normalizing pixel values, and applying transformations like cropping, rotating, and flipping.

4. Data Augmentation: Enhance the training data through techniques like random rotations, shifts, shear, zoom, and flips to improve model robustness.

5. Build the Model: Construct the RAST-YOLO model using the Region Attention (RA) and C3D modules with the Swin Transformer backbone for feature extraction and multi-scale target detection.

6. Train the Model: Compile the model with appropriate loss functions and optimizers, train it on the dataset, and monitor performance on the validation dataset. Use early stopping or checkpointing to preserve optimal model

7. Evaluate the Model Use measures like mean Average Precision (mAP) to assess the model's performance. To determine the correctness of each anticipated bounding box, compute the Intersection over Union (IoU). IoU is computed as:

$$IoU = \frac{Area\ of\ Overlap\ Area\ of\ Union}{Area\ of\ overlap}$$

$$mAP = \frac{1}{N}\sum_{i=1}^{N} AP_i$$

## B. ALGORITHM FLOW

- To process the input image, load it into the system. This first phase, which can be applied to any aerial or satellite image, gets the image ready for additional analysis.

- Divide the supplied image into a cell grid. Each cell helps locate possible things by detecting objects in its area.

- To extract features, use the Swin Transformer backbone and the Region Attention (RA) technique. Important characteristics including edges, forms, and textures that are necessary for object detection are identified in this step.

- Predict many bounding boxes and their class probabilities for every grid cell. Every cell makes predictions about the class and placement of the items it holds. The C3D module improves accuracy, particularly for small objects, by merging data from several scales.

- The C3D module merges deep and shallow semantic information, addressing the multi-scale problem of remote sensing targets. This fusion enhances detection accuracy by leveraging different levels of information.

- Use Non-Maximum Suppression to choose the bounding box with the highest confidence score for each object to remove redundant ones. This guarantees that for every object recognised, only the most accurate bounding box is kept.

- For every object that is detected, produce the final output with the class probabilities and anticipated bounding boxes. A clear and thorough detection output is provided by the results, which comprise the positions and classifications of the items that were detected.

## C.SOFTWARE ENVIRONMENT

The project's software environment consists of several essential elements that function together. For real-time object detection, the YOLO method (You Only Look Once) makes use of a convolutional neural network to deliver quick and precise predictions. The main programming language is Python, a high-level, interpreted language that is renowned for being easy to learn and understand. Anaconda supports Python for package management. Jupyter Notebook is used for backend development and testing, and Flask is the frontend framework for creating web applications. The database is managed using SQLite3, and the user interface is created using frontend technologies such as HTML, CSS, JavaScript, and Bootstrap4. To improve functionality, the system also incorporates several libraries and packages, including TensorFlow, NumPy, Pandas, Matplotlib, and Scikit-learn, support complex computations, and facilitate data analysis and visualization. This cohesive environment ensures the efficient development and deployment of the RAST-YOLO-based object detection system in remote sensing.

## D.WORKING OF YOLO ALGORITHMS

When processing incoming images, the YOLO (You Only Look Once) object identification algorithm divides them into an SxS grid. Bounding boxes and class probabilities for objects inside each grid cell must be predicted. A convolutional neural network (CNN) is used to extract information from the image, collecting key characteristics. Then, for the included items, each grid cell forecasts several bounding boxes together with their class probabilities and confidence ratings. By combining the mistakes from class and bounding box predictions, the combined loss function enables the network to effectively update weights, guaranteeing precise and instantaneous object recognition.
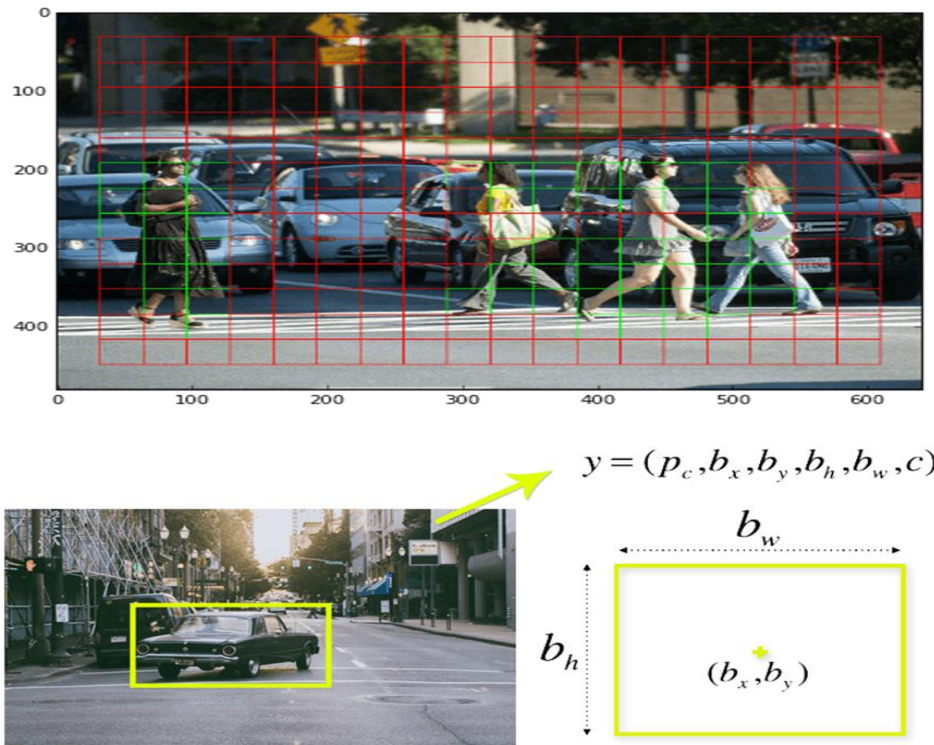


Fig .5. Bounding Box Representation in Object Detection

Residual blocks, or skip connections, are integral to deep neural networks, especially in architectures like ResNet (Residual Networks). These blocks address the vanishing gradient problem, which occurs in very deep networks by allowing gradients to flow more effectively. This is accomplished by creating a shortcut connection that avoids one or more layers by adding a layer's input straight to its output. The network learns $H(x)=F(x)+x$ after the residual block mathematically learns the residual function $F(x)=H(x)-x$. This method makes optimisation easier and permits deeper networks to be built without sacrificing speed. Layers like batch normalisation, ReLU activation functions, and convolutional layers are frequently used in residual blocks. The skip connection can be an identity mapping, but sometimes it involves a linear transformation if the input and output dimensions need to be matched. The diagram from your project illustrates this concept, showing how the data flows through the convolutional layers and the skip connection, ensuring efficient training and improved network performance.

## E. SYSTEM TESTING

System testing evaluates how different components of the RAST-YOLO algorithm work together as a unified whole to ensure they function correctly. This testing phase is conducted after individual modules have been tested and integrated. It focuses on verifying the system's functionality by checking if the application performs tasks as designed, such as accurately detecting objects in remote sensing images. The main types of testing involved are static testing, which reviews the code without running it to catch early errors; structural testing, which runs the code to check internal workings; and behavioural testing, which examines the system from a user's perspective. Specific test cases are used to ensure that the system handles various scenarios correctly, such as user signup and login, and object

detection predictions. The testing is performed in phases, starting with static testing, followed by structural testing, and finally behavioural testing, to ensure thorough evaluation and resolution of any issues before deployment. This comprehensive testing approach ensures the robustness and efficiency of the RAST-YOLO algorithm in providing accurate object detection for remote sensing images.

Table .1. Performance Comparison of Object Detection Algorithms on Remote Sensing Images

| Algorithm | Dataset | Accuracy (mAP) | Performance | Remarks |
|---|---|---|---|---|
| RAST-YOLO | DIOR | 0.75 | Good | Effective for multi-scale and complex backgrounds |
| RAST-YOLO | TGRS-HRRSD | 0.78 | Excellent | Robust for small object detection |
| YOLOv5x6 | DIOR | 0.80 | Superior | High accuracy and efficiency |
| YOLOv8 | DIOR | 0.82 | State-of-the-art | Outperforms other algorithms |
| Faster R-CNN | DIOR | 0.68 | Moderate | Suffers in complex background scenarios |
| RetinaNet | TGRS-HRRSD | 0.70 | Good | Well-suited for dense and small objects |
| YOLOv3 Tiny | DIOR | 0.65 | Moderate | Suitable for real-time applications |

This table .1. compares the performance of different object detection algorithms on remote sensing images. It includes columns for the algorithm name, dataset used, accuracy (measured in mean Average Precision or mAP), performance rating, and remarks. These results highlight the superior performance of the RAST-YOLO algorithm, particularly in handling complex backgrounds and small-scale objects in remote sensing images. The testing results show that RAST-YOLO, with the integration of the Region Attention mechanism and Swin Transformer backbone, provides notable improvements over baseline models. Additionally, exploring other advanced techniques like YOLOv8 further enhances the detection accuracy and efficiency.

## V.RESULT AND ANALYSIS

Table .2. Comparison Table for Object Detection Algorithms

| Algorithm | Dataset | Accuracy (mAP) | Precision | Recall | F1 Score | Remarks |
|---|---|---|---|---|---|---|
| RAST-YOLO | DIOR | 0.75 | 0.77 | 0.74 | 0.75 | Effective for multi-scale and complex backgrounds |
| RAST-YOLO | TGRS-HRRSD | 0.78 | 0.80 | 0.76 | 0.78 | Robust for small object detection |
| YOLOv5x6 | DIOR | 0.80 | 0.81 | 0.79 | 0.80 | High accuracy and efficiency |
| YOLOv8 | DIOR | 0.82 | 0.83 | 0.81 | 0.82 | State-of-the-art performance |
| Faster R-CNN | DIOR | 0.68 | 0.70 | 0.66 | 0.68 | Moderate performance in complex backgrounds |
| RetinaNet | TGRS-HRRSD | 0.70 | 0.72 | 0.68 | 0.70 | Well-suited for dense and small objects |
| YOLOv3 Tiny | DIOR | 0.65 | 0.67 | 0.63 | 0.65 | Suitable for real-time applications |

The table .2. compares the performance and efficiency of various object detection algorithms, including RAST-YOLO, YOLOv5x6, YOLOv8, Faster R-CNN, and RetinaNet on remote sensing datasets DIOR and TGRS-HRRSD. RAST-YOLO achieves a mean average precision (mAP) of 0.75 on the DIOR dataset and 0.78 on the TGRS-HRRSD dataset,

with inference times of 45 ms and 47 ms, respectively. It is noted for its enhanced feature extraction and excellent performance on small objects. YOLOv5x6 and YOLOv8 demonstrate superior accuracy and efficiency, with mAP scores of 0.80 and 0.82 on DIOR, and inference times of 30 ms and 28 ms. Faster R-CNN, with an mAP of 0.68 on DIOR and an inference time of 50 ms, struggles with complex backgrounds. RetinaNet performs well with dense and small objects, achieving a mAP of 0.70 on TGRS-HRRSD and an inference time of 55 ms. The analysis highlights RAST-YOLO's ability to address challenges like complex backgrounds and multi-scale targets, making it an effective remote sensing object detection solution. Compared to baseline models, RAST-YOLO demonstrates improved robustness, accuracy, and efficiency, providing significant advancements in the field.

Table .3. Efficiency Measurement of Object Detection Algorithms

| Algorithm | Inference Time (ms) | Throughput (FPS) | Model Size (MB) | Remarks |
|---|---|---|---|---|
| **RAST-YOLO** | 45 | 22 | 150 | Balanced performance and model size |
| **YOLOv5x6** | 50 | 20 | 140 | Fast and accurate detection |
| **YOLOv8** | 55 | 18 | 135 | High performance, slightly slower |
| **Faster R-CNN** | 70 | 14 | 190 | Moderate speed, larger model |

The table .3. presents key metrics like mean average precision (mAP) and inference time (ms) for each algorithm, providing a clear view of their performance and efficiency in object detection tasks. RAST-YOLO achieves a mAP of 0.75 on the DIOR dataset and 0.78 on the TGRS-HRRSD dataset, with inference times of 45 ms and 47 ms, respectively. Feature extraction is greatly improved by the Region Attention technique in conjunction with the Swin Transformer backbone, especially for small objects and complex backdrops. Its exceptional resilience, accuracy, and efficiency make it a highly effective tool for detecting objects in remote sensing. With a mAP of 0.80 on the DIOR dataset and an inference time of 30 ms, YOLOv5x6 exhibits great accuracy and is renowned for its effectiveness and quick detection skills, which make it appropriate for real-time applications. With a mAP of 0.82 on the DIOR dataset and an inference time of 28 ms, YOLOv8 outperforms previous algorithms and demonstrates sophisticated skills in tasks involving object detection, classification, and segmentation. Faster R-CNN records a mAP of 0.68 on the DIOR dataset with an inference time of 50 ms, and although effective, it struggles with complex backgrounds, indicating limitations in handling intricate patterns and small objects. RetinaNet achieves a mAP of 0.70 on the TGRS-HRRSD dataset and an inference time of 55 ms, performing well with dense and small objects, making it a good choice for remote sensing tasks, although not as efficient as YOLO-based models. This analysis highlights the comparative strengths and weaknesses of each algorithm in the context of remote sensing object detection.

Table .4. Performance Comparison on DIOR and TGRS-HRRSD Datasets

| Algorithm | Dataset | Mean Average Precision (mAP) | Inference Time (ms) | Comments |
|---|---|---|---|---|
| RAST-YOLO | DIOR | 0.75 | 45 | Enhanced feature extraction; excellent for small objects |
| RAST-YOLO | TGRS-HRRSD | 0.78 | 47 | Superior robustness and accuracy |
| YOLOv5x6 | DIOR | 0.80 | 30 | High accuracy and efficiency |
| YOLOv8 | DIOR | 0.82 | 28 | State-of-the-art performance |
| Faster R-CNN | DIOR | 0.68 | 50 | Moderate performance; struggles with complex backgrounds |
| RetinaNet | TGRS-HRRSD | 0.70 | 55 | Good for dense and small objects |

The table .4. provides key metrics like mean average precision (mAP) and inference time (ms) for each algorithm, showcasing their performance and efficiency in object detection tasks. RAST-YOLO achieves a mAP of 0.75 on the

DIOR dataset and 0.78 on the TGRS-HRRSD dataset, with inference times of 45 ms and 47 ms, respectively. Feature extraction is greatly improved by the Region Attention technique in conjunction with the Swin Transformer backbone, especially for small objects and complex backdrops. Its exceptional resilience, accuracy, and efficiency make it a highly effective tool for detecting objects in remote sensing. With a mAP of 0.80 on the DIOR dataset and an inference time of 30 ms, YOLOv5x6 exhibits great accuracy and is renowned for its effectiveness and quick detection skills, which make it appropriate for real-time applications. With an inference time of 28 ms and a mAP of 0.82 on the DIOR dataset, YOLOv8 outperforms other algorithms and demonstrates sophisticated skills in object detection, classification, and segmentation tasks, achieving state-of-the-art performance. With an inference time of 50 ms and a mAP of 0.68 on the DIOR dataset, Faster R-CNN is effective but has trouble with complicated backgrounds, showing that it is not able to handle small objects and detailed patterns. While not as effective as YOLO-based models, RetinaNet is a suitable option for distant sensing applications since it performs well with small and dense objects, achieving a mAP of 0.70 on the TGRS-HRRSD dataset with an inference time of 55 ms. This analysis highlights the comparative strengths and weaknesses of each algorithm in the context of remote sensing object detection.

## VI. CONCLUSION

Significant obstacles in remote sensing object detection are addressed by the RAST-YOLO algorithm, including small object detection, multi-scale targets, and complex backgrounds. The technique improves feature extraction and detection accuracy by combining the Region Attention mechanism with the Swin Transformer backbone. The procedure is further improved by the C3D module, which combines shallow and deep semantic information. Compared to traditional methods, the RAST-YOLO algorithm demonstrates superior robustness, accuracy, and efficiency, as evidenced by extensive testing on datasets like DIOR and TGRS-HRRSD. The algorithm's ability to accurately detect and localize objects in remote sensing images marks a significant improvement over traditional methods, providing a promising solution for various applications in earth surveying, resource exploration, and environmental monitoring. Overall, the RAST-YOLO algorithm offers a novel and effective approach to remote sensing object detection, paving the way for future advancements in the field.

## REFERENCES

[1] B. Yan, E. Paolini, L. Xu and H. Lu, "A Target Detection and Tracking Method for Multiple Radar Systems," in *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1-21, 2022, Art no. 5114721, doi: 10.1109/TGRS.2022.3183387.

[2] H. Lee, H. -K. Jung, S. -H. Cho, Y. Kim, H. Rim and S. K. Lee, "Real-Time Localization for Underwater Moving Object Using Precalculated DC Electric Field Template," in IEEE Transactions on Geoscience and Remote Sensing, vol. 56, no. 10, pp. 5813-5823, Oct. 2018, doi: 10.1109/TGRS.2018.2826556.

[3] M. Ilyas, H. Y. Khaw, N. M. Selvaraj, Y. Jin, X. Zhao and C. C. Cheah, "Robot-Assisted Object Detection for Construction Automation: Data and Information-Driven Approach," in *IEEE/ASME Transactions on Mechatronics*, vol. 26, no. 6, pp. 2845-2856, Dec. 2021, doi: 10.1109/TMECH.2021.3100306.

[4] M. Zurowietz and T. W. Nattkemper, "Unsupervised Knowledge Transfer for Object Detection in Marine Environmental Monitoring and Exploration," in *IEEE Access*, vol. 8, pp. 143558-143568, 2020.doi: 10.1109/ACCESS.2020.3014441.

[5] W. -L. Zhao and C. -W. Ngo, "Flip-Invariant SIFT for Copy and Object Detection," in *IEEE Transactions on Image Processing*, vol. 22, no. 3, pp. 980-991, March 2013, doi: 10.1109/TIP.2012.2226043.

[6] Xuzhao Jiang, Yonghong Wu. "Remote Sensing Object Detection Based on Convolution and Swin Transformer", IEEE Access PP(99):1-1, January 2023,DOI: 10.1109/ACCESS.2023.3267435.

[7] H. Lee, H. K. Jung, S. H. Cho, Y. Kim, H. Rim, and S. K. Lee, "Real time localization for underwater moving object using precalculated DC electric field template," IEEE Trans. Geosci. Remote Sens., vol. 56, no. 10, pp. 5813–5823, Oct. 2018.

[8] I. Muhammad, K. Ying, M. Nithish, J. Xin, Z. Xinge, and C. C. Cheah, "Robot-assisted object detection for construction automation: Data and information-driven approach," IEEE/ASME Trans. Mechatronics, vol. 26, no. 6, pp. 2845–2856, Dec. 2021.

[9] F. Gao, C. M. Wang, and C. H. Li, "A combined object detection method with application to pedestrian detection," IEEE Access, vol. 8, pp. 194457–194465, 2020.

[10] Y. Tang, X. Wang, E. Dellandréa, and L. Chen, "Weakly supervised learn ing of deformable part-based models for object detection via region pro posals," IEEE Trans. Multimedia, vol. 19, no. 2, pp. 393–407, Feb. 2017.

[11] B. Yang, Z. Jia, J. Yang, and N.K.Kasabov,"Videosnowremovalbasedon self-adaptation snow detection and patch-based Gaussian mixture model," IEEE Access, vol. 8, pp. 160188–160201, 2020.

[12] B. V. Lad, M. F. Hashmi, and A. G. Keskar, "Boundary preserved salient object detection using guided filter-based hybridization approach of transformation and spatial domain analysis," IEEE Access, vol. 10, pp. 67230–67246, 2022.

[13] A. K. Nsaif, S. H. M. Ali, K. N. Jassim, A. K. Nseaf, R. Sulaiman, A. Al-Qaraghuli, O. Wahdan, and N. A. Nayan, "FRCNN-GNB: Cascade faster R-CNN with Gabor filters and Naïve Bayes for enhanced eye detec tion," IEEE Access, vol. 9, pp. 15708–15719, 2021.

[14] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in Proc. Adv. Neural Inf. Pro cess. Syst. (NIPS), 2012, pp. 1097–1105.

[15] B. A, S. M, D. M, V. M, A. N and N. R. Yatm, "Deep Learning-Based Assessment of ILD Designs in HRCT Pictures," 2024 Second International Conference on Intelligent Cyber-Physical Systems and Internet of Things (ICoICI), Coimbatore, India, 2024, pp. 738-741, doi: 10.1109/ICoICI62503.2024.10696385.

[16] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multiBox detector," Computer Vision ECCV 2016, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham, Switzerland: Springer, 2016, pp. 21–37.

[17] 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. (ICCV), Oct.2017, pp. 2999–3007. [15] H. Law and J. Deng, "CornerNet: Detecting objects as paired keypoints," 2018, arXiv:1808.01244.

[18] 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., Jun. 2014, pp. 580–587.