



Semantic Segmentation for Aerial Images

Dr.B.Harika,^{1}, K.Himneesh² and M.Bharath³*

¹Department of IT, Mahatma Gandhi Institute of Technology, Gandipet, Hyderabad, 500085, Telangana, India.

^{2,3}Department of IT, Mahatma Gandhi Institute of Technology, Gandipet, Hyderabad, 500075, Telangana, India.

E-mail(s): bharika-it@mgit.ac.in; khimneesh-csb213227@mgit.ac.in; mbharath-csb213239@mgit.ac.in;

DOI : <https://doi.org/10.55248/gengpi.6.0425.1358>

ABSTRACT

Semantic segmentation is a critical task in computer vision that involves classifying each pixel in an image into a specific category, enabling detailed and dense predictions. Unlike traditional image classification, which assigns a single label to an entire image, semantic segmentation provides a more granular understanding of the visual content. This project investigates the application of Convolutional Neural Networks (CNNs) for semantic segmentation, focusing on their effectiveness in analyzing aerial images. The selected CNN architectures are trained and evaluated on a dataset consisting of aerial images with pixel-level annotations. To improve model generalization and mitigate overfitting, data augmentation techniques are employed. Through comparative analysis, the strengths and weaknesses of each CNN model are identified, offering valuable insights into their suitability for various semantic segmentation tasks. This study highlights the potential of CNN-based models to achieve high accuracy in semantic segmentation of aerial images and provides guidance on selecting the most appropriate model based on specific application requirements.

Keywords: Semantic Segmentation, Convolutional Neural Networks (CNNs), Aerial Images, Pixel-Level Annotations.

1. Introduction

Semantic segmentation in aerial imagery is a fundamental and critical task in remote sensing and geospatial analysis. It assigns a categorical label to every pixel in an image, enabling accurate identification and classification of terrain features such as buildings, roads, vegetation, and aquatic environments. This ability has revolutionized fields like urban planning, disaster response, agricultural observation, environmental preservation, and land-use classification[1, 10].

Traditional segmentation techniques, including edge detection, pixel-based image analysis, and region-based classification, relied on manually designed features and statistical frameworks. Although these methods were effective for low-resolution imagery, they often failed to address the complexities of high-resolution aerial datasets. Such datasets require context-aware models capable of handling intricate spatial patterns, overlapping features, and spectral variability[8, 15, 21]. Deep learning, particularly CNNs, has overcome these limitations by capturing both spatial and spectral relationships with state-of-the-art accuracy[13].

Semantic segmentation plays a crucial role in deriving actionable insights from large geospatial datasets. It aids urban planners in assessing infrastructure needs, helps environmentalists monitor deforestation and water quality, and supports disaster response teams in identifying flood-prone regions[9, 19].

While the development of high-resolution imaging sensors has expanded its applicability, challenges such as class imbalance, scale variations, and computational inefficiencies persist. For example, spectral similarities between vegetation and water-logged areas often lead to misclassifications. Addressing these issues requires robust architectures utilizing self-attention, multi-scale feature fusion, and advanced loss functions[2, 18, 23].

Why It Matters

The segmentation of aerial imagery is vital for addressing critical geospatial challenges across multiple fields. Accurate segmentation allows for timely identification of key features, such as urban growth zones, deforested regions, or areas prone to flooding, enabling effective resource allocation and informed decision-making. Misclassifications or delays can have significant repercussions, such as ineffective disaster responses, poorly planned urban expansions, and inefficient resource management[7]. The complex nature of aerial imagery, characterized by overlapping features, diverse textures, and imbalanced classes, adds to the challenges. For instance, distinguishing between vegetation and water-saturated regions with similar spectral signatures often leads to errors, necessitating sophisticated models that incorporate

both contextual and spatial data[2, 5, 12].

Contemporary Diagnostic Instruments

The evolution of diagnostic tools for semantic segmentation reflects advancements in technology. Traditional methods, including pixel-based analysis and edge detection, laid the foundation for image segmentation but often struggled with high-resolution datasets due to their reliance on handcrafted features[8, 15].

Modern methods like U-Net and DeepLab, enhanced by attention mechanisms, have significantly improved segmentation accuracy. U-Net ensembles aggregate multi-scale information to capture both local and global spatial features, while spatial and channel attention techniques enhance feature discrimination and contextual understanding[1, 11, 13]. These advanced tools have elevated the performance of segmentation systems, making them indispensable in generating accurate and actionable insights for key applications[2].

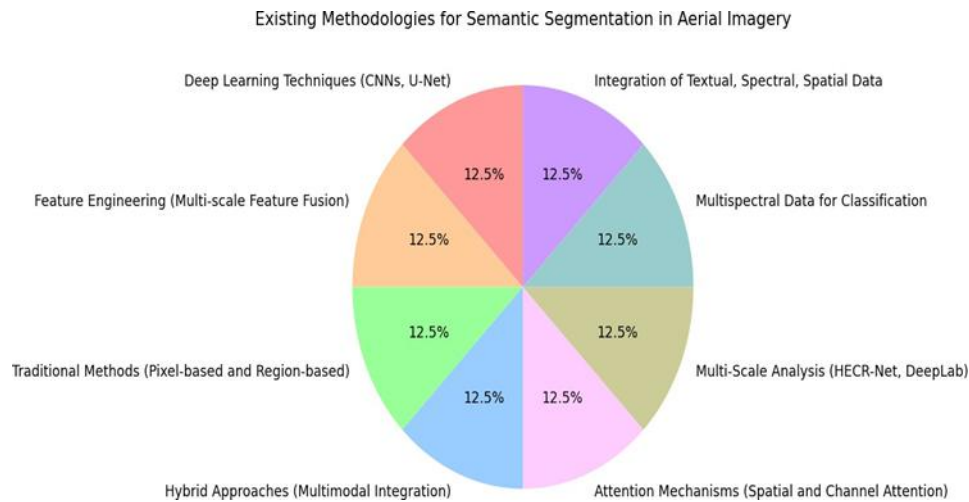


Fig. 1: Algorithms and methods used in existing system.

Figure 1, showcases the distribution of methodologies employed for semantic segmentation in aerial imagery, highlighting the multidisciplinary nature of this research field. The pie chart divides the methodologies into eight equal segments, each representing 12.5% of the focus. These include deep learning techniques like CNNs and U-Net, which serve as the foundation for modern semantic segmentation, and feature engineering approaches that emphasize multi-scale feature fusion to enhance accuracy. Traditional methods, such as pixel-based and region-based segmentation, are also depicted, reflecting their historical significance despite their limitations with high-resolution datasets. Hybrid approaches, which integrate multimodal data sources such as textual, spectral, and spatial information, are shown as a key area of innovation. Attention mechanisms, including spatial and channel attention, are highlighted for their role in improving feature discrimination and contextual understanding. Multi-scale analysis techniques, exemplified by HECR-Net and DeepLab, are noted for their ability to analyze images at varying resolutions, capturing both macroscopic and microscopic details.

1.1 Problem Statement

Aerial imagery has a key contribution to play in remote sensing activities, including land cover mapping, environmental monitoring, and urban planning. However, accurate segmentation becomes difficult considering heterogeneity of the terrain, complexity of the features, and occlusion of the objects. Classical methods generally face inefficiency and inaccuracy while dealing with big data.

Current semantic segmentation structures are mainly pixel-level based but are devoid of sophisticated features for semantically relevant information extraction. Geometric feature estimation, including land class estimation (e.g., vegetation, water bodies), is generally erroneous and time-consuming, constraining its use in mission-critical applications such as resource planning and disaster management.

This project offers a deep learning system which combines CNNs, self-attention, and separable convolutions to improve segmentation efficiency and adaptability. It offers various output possibilities—binary masks, geometric features, or full-segmented images—allowing accurate analysis and visualization. By improving speed, accuracy, and adaptability, the system seeks to improve decision-making in remote sensing operations, especially in resource-constrained areas.

1.2 Motivation

With the growing availability of high-resolution aerial imagery, uses such as land cover mapping, environmental monitoring, urban planning, and disaster management are increasing. However, accurate segmentation is a problem due to the heterogeneity of natural environments like forests, cities, water bodies, and mountains. Current systems are not effective in dealing with spatial correlations, making inferences about actual-world geometric characteristics, and producing flexible outputs in the form of pixel-level classifications.

To overcome these shortcomings, this project proposes an improved system based on deep learning that combines CNNs with self-attention and separable convolutions. It improves segmentation accuracy, efficiency, and responsiveness with the provision of multiple output formats—segmented

images, binary masks, and geometric features. The system is optimized to tackle real-time and large-scale applications through enhanced spatial correlation management, multi-scale feature integration, and contextual inspection.

By resolving issues such as class imbalance, spectral similarity, and computational inefficiency, this project seeks to deliver a scalable, flexible, and accurate solution that promotes the field of remote sensing and facilitates better decision-making across industries.

- Leveraging advanced deep learning techniques, such as CNNs with self-attention mechanisms and separable convolutions, to improve the accuracy and efficiency of semantic segmentation in aerial imagery.
- Developing a system capable of handling spatial correlations among diverse terrain types, enabling precise segmentation of complex landforms, such as urban areas, vegetation, water bodies, and mountains.
- Introducing a flexible approach that provides various outputs, including segmented images, binary masks, and geometric features, for enhanced decision-making in applications like urban planning and disaster management.

2. Related Works

Recent advancements in aerial image segmentation have been driven by deep learning, particularly CNN-based architectures like U-Net with self-attention mechanisms and separable convolutions. These improve spatial dependency understanding and computational efficiency, enhancing segmentation accuracy for complex landscapes [1]. Systematic reviews highlight a shift from traditional methods to deep learning, emphasizing the need for robust networks that adapt to diverse environments. Multispectral and hyperspectral data enrich segmentation by distinguishing similar objects in aerial scenes [2]. Conditional Random Fields (CRFs) refine segmentation boundaries, ensuring spatial consistency, while multi-scale learning improves object

recognition across different sizes [4, 5].

Few-shot learning, especially rotation-invariant models, addresses the challenge of limited labeled data, reducing annotation costs while maintaining accuracy [6]. Attention mechanisms, including spatial and channel attention, enhance feature extraction, improving segmentation precision in high-resolution imagery [10, 11].

Ensemble methods combine multiple models to enhance accuracy, particularly for complex datasets. Height-embedding features aid in differentiating structures like buildings and vegetation, improving segmentation in urban environments [8, 9]. Integrating multispectral data with attention mechanisms further enhances segmentation, providing superior feature extraction and classification across land cover types [11, 13].

Boundary detection combined with segmentation refines object outlines, essential for land cover classification and urban planning. Model ensembles further strengthen segmentation performance by reducing overfitting and improving generalization [15–17].

Table 1: Literature Survey

S. No.	Title	Author & Year	Journal & Year	Methodologies	Key Findings	Gaps
1	Semantic Segmentation of Aerial Imagery Using U-Net with Self-Attention and Separable Convolutions	Khan, B.A., Jung, J.W., 2024	Applied Sciences, 2024	U-Net, Self-Attention, Separable Convolutions	Improved segmentation accuracy for aerial imagery using self-attention and separable convolutions	Limited evaluation on large-scale datasets, potential scalability issues for complex terrains
2	Semantic Segmentation: A Systematic Analysis From State-of-the-Art Techniques to Deep Networks	Seth, A., Sharma, S., 2022	Journal of Information Technology Research (JITR), 2022	Review of various semantic segmentation techniques	Provides a comprehensive analysis of various approaches for semantic segmentation	Lacks a detailed comparison on dataset-specific performance and real-world applicability

3	Deep learning-based semantic segmentation of urban features in satellite images: A review and meta-analysis	Neupane, B., Horant, T., Aryal, J., 2021	Remote Sensing, 2021	Deep learning-based semantic segmentation	Reviews deep learning models applied to urban features, evaluating model performance	Limited exploration of multi-source data integration and its impact on segmentation accuracy
---	---	--	----------------------	---	--	--

S. No.	Title	Author & Year	Journal & Year	Methodologies	Key Findings	Gaps
4	Enhanced Semantic Segmentation of Aerial Images with Spatial Smoothness Using CRF Model	Hussein, S.K., Ali, K.H., 2022	IEEE, 2022	Conditional Random Fields (CRF), Spatial Smoothness	Incorporates spatial smoothness with CRF for better boundary delineation	Focuses mainly on boundary refinement, without considering other feature types for improvement
5	An Aerial Image Segmentation Approach Based on Enhanced Multi-Scale Convolutional Neural Network	Li, X., Jiang, Y., Peng, H., Yin, S., 2019	IEEE ICPS, 2019	Multi-Scale CNN	Multi-scale CNN improves context capture for aerial image segmentation	Lacks sufficient analysis of real-time performance and scalability in large datasets
6	Few-shot rotation-invariant aerial image semantic segmentation	Cao, Q., Chen, Y., Ma, C., Yang, X., 2023	IEEE Transactions on Geoscience and Remote Sensing, 2023	Few-Shot Learning, Rotation-Invariant Segmentation	Few-shot learning approach enhances model adaptability to different image orientations	Limited exploration of model robustness on varied terrain types and real-world data
7	Remote sensing object detection in the deep learning era—a review	Gui, S., Song, S., Qin, R., Tang, Y., 2024	Remote Sensing, 2024	Object Detection, Deep Learning	Reviews object detection techniques for remote sensing imagery	Focused more on object detection than segmentation, lacks depth in segmentation models

S. No.	Title	Author & Year	Journal & Year	Methodologies	Key Findings	Gaps
8	U-Net Ensemble for Enhanced Semantic Segmentation in Remote Sensing Imagery	Dimitrovsk I., Spasov, V., Loshkovsk S., Kitanovsk I., 2024	Remote Sensing, 2024	U-Net Ensemble	U-Net ensemble improves segmentation performance by combining models	Does not address computational cost and efficiency of the ensemble approach
9	HECR-Net: Height-embedding context reassembly network for semantic segmentation in aerial images	Liu, W., Zhang, W., Sun, X., Guo, Z., Fu, K., 2021	IEEE Journal of Selected Topics in Applied Earth Observations, 2021	HECR-Net, Height-Embedding	Height-embedding improves segmentation of topologically complex aerial images	The model's performance on heterogeneous datasets and real-time segmentation remains underexplored
10	Aerial image semantic segmentation using spatial and channel attention	Lan, Z., Huang, Q., Chen, F., Meng, Y., 2019	IEEE ICIVC, 2019	Spatial Attention, Channel Attention	Combines spatial and channel attention mechanisms to enhance segmentation accuracy	Lacks a comparison of the model's performance with other attention-based methods

3. Methodologies

3.1 Self-Attention Mechanism

Self-attention is used to enhance the model's ability to focus on critical image regions while ignoring irrelevant areas. This mechanism allows each pixel in an image to dynamically consider the relevance of other pixels. By integrating spatial and contextual information, the model captures intricate patterns and dependencies within the data. [1] This approach is particularly beneficial for tasks requiring the fusion of temporal and spatial information, such as monitoring changes over time in landscapes. This method integrates a U-Net architecture with self-attention mechanisms and separable convolutions to improve feature extraction and reduce computational complexity. [1?]

Key Steps:

1. **Encoder-decoder structure for pixel-level segmentation.**
2. **Self-attention module to capture long-range dependencies.**
3. **Separable convolutions to enhance efficiency by reducing parameter count.**

The U-Net with Self-Attention and Separable Convolutions architecture enhances the traditional U-Net by incorporating self-attention mechanisms and separable convolutions to boost performance and efficiency. The self-attention mechanism allows the network to focus on the most relevant spatial regions within the feature maps, enabling better capture of long-range dependencies and improving the segmentation of intricate or small structures. Separable convolutions, which decompose standard convolutions into depthwise and pointwise operations, significantly reduce the computational load without compromising accuracy. By combining these two techniques, the model achieves superior feature extraction, greater spatial contextual understanding, and lower computational complexity, making it particularly well-suited for tasks requiring high-resolution segmentation, such as medical imaging and remote sensing. [1?]

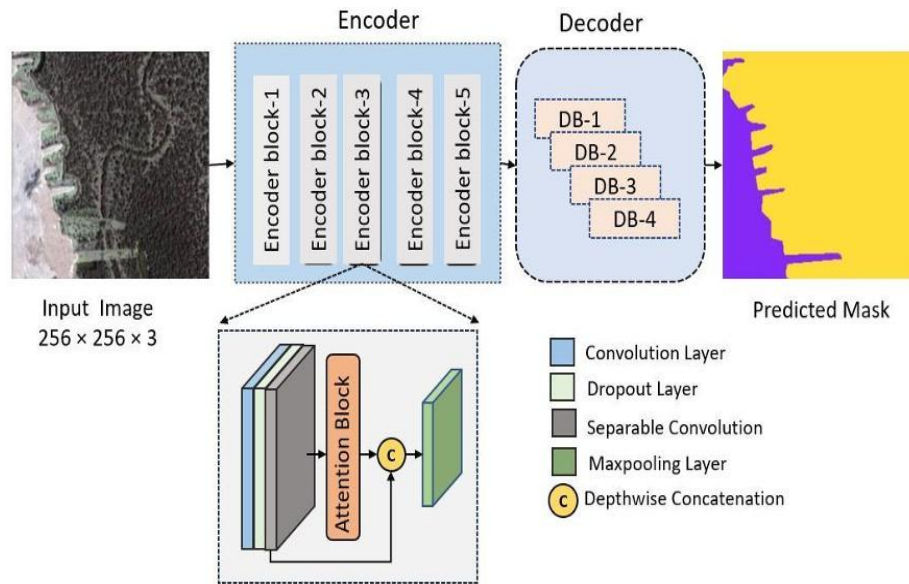


Fig. 2: SA-SC-U-Net: U-Net with self-attention and separable convolutions.

Figure 2, The figure depicts a semantic segmentation architecture specially designed for aerial imagery. It's built on the U-Net framework but enhanced with self-attention mechanisms and separable convolutions. The beginning point is an input image of (256 times 256 times 3), height, width, and RGB channels, fed later into the model for processing. The encoder consists of a number of blocks: Encoder block-1 through to Encoder block-5 that hierarchically extract features from the image. All the blocks are embedded within convolutional layers, separated convolutions, dropout layers, and max-pooling layers. They work together to reduce spatial dimensions step by step while elevating feature representations. One of the key blocks forming the encoder is the attentions block. The block highlights major spatial regions by emphasizing the most important features while reducing the influence of non-critical ones. This depthwise concatenation in attentions ensures proper integration along resolutions of different features features to enhance contextual understanding.

This architectural design skillfully incorporates self-attention mechanisms to enhance feature selection, along with separable convolutions to reduce computational complexity without loss of accuracy. Through the integration of spatial and context- tual information, it achieves superior segmentation results, making it suitable for large aerial datasets and complex landscapes.

3.2 Enhanced Multi-Scale Convolutional Neural Network Methodology

This multi-scale convolutional neural network advanced, with sequential dilated convolution modules, supplementary loss functions, and encoder-decoder architecture, allows for improved performance in semantic segmentation. The method has proven particularly advantageous in analyzing aerial images because of its capability to overcome issues of varied object dimensions, intricate backgrounds, and the considerable expense of computation.

Basic Practices: Cascaded Dilated Convolution Layers: Convolutional layers with varying dilation rates capture diverse receptive fields, hence improving the ability of the model to identify features of different scales [10, 13]. Feature Fusion: Outputs from various dilation rates are combined to give a richer feature representation and higher segmentation accuracy. Auxiliary Loss Function: The intermediate features are optimized during the training to enhance convergence and feature learning. Auxiliary Loss Function Integration In order to improve the learning process, a secondary loss function is added at several intermediate layers of the network. This form of regularization enhances robustness significantly, particularly in complex urban environments. By balancing primary and auxiliary losses, the network achieves accelerated convergence and better accuracy [5, 9].

Essential Procedures: Encoding Stage: Features are hierarchically extracted using convolutional layers and max-pooling, where the spatial resolution decreases and feature abstraction increases. Decoding Stage: Upsampling layers and skip connections from the encoder are used to reconstruct spatial details to make the segmentation precise. Output final segmentation map A 1x1 convolution layer with sigmoid activation outputs pixel-level segmentation probabilities. [1, 10, 15]. This methodology provides an optimal solution for large-scale aerial image segmentation with fine-grained feature extraction and efficiency at computation while considering variability in scales of objects and complex environmental settings [5, 13].

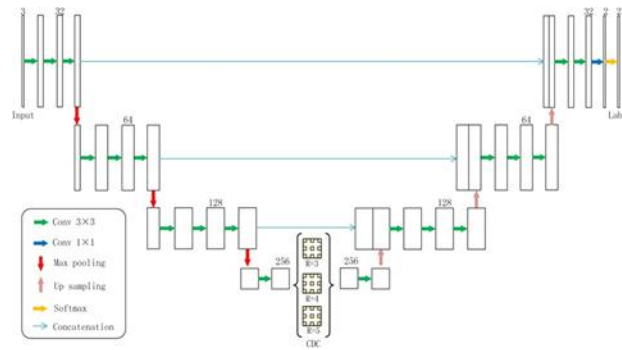


Fig. 3: [5] Structure of the proposed multi-scale convolutional network.

Figure 3, This image shows an enhanced U-Net architecture for semantic segmentation of aerial imagery. The structure adopts the encoder-decoder concept, with cascaded dilated convolution modules, to extract multi-scale contextual information efficiently. The encoder path, depicted on the left, successively reduces spatial dimensions as it derives hierarchical feature representations using layers of convolutional layers (Conv 3x3) and max-pooling operations. The decoder path (on the right side), through upsampling layers, provides skip connections, which will be merged within the decoding with the encoders, thus maintaining fine grain information.

The major innovation lies in the CDC module, which uses varying dilation rates ($R=3$, $R=4$, $R=5$) to extract features at multiple scales so that the model can address a range of object sizes and shapes without increasing the computational complexity. With its 1×1 convolution for channel reduction and aiding refinement of features and a softmax layer at output, creating segmentation map by putting every pixel in class; this set configuration maintains both efficient as well as the accuracy in segregation, it efficiently tackles issues that seem present intrinsically to complex diversified aerial imagery dataset.

3.3 Few-Shot Rotation-Invariant Segmentation

This technique introduces a few-shot learning approach to enable rotation-invariant semantic segmentation using a limited dataset.

Key Steps:

1. Data augmentation with rotation to simulate different orientations.
2. Meta-learning framework for few-shot generalization.
3. This technique leverages a few-shot learning approach to achieve rotation-invariant semantic segmentation, even with limited training data. By employing data augmentation through rotations, the model is exposed to various object orientations, enabling it to generalize better to unseen perspectives. Additionally, a meta-learning framework is utilized to enhance few-shot learning capabilities, allowing the model to quickly adapt to new tasks or datasets with minimal labeled examples. This combination of rotation-based data augmentation and meta-learning ensures robust performance in scenarios where data availability is scarce and objects appear in diverse orientations, making it particularly effective for applications like aerial imagery.[22]

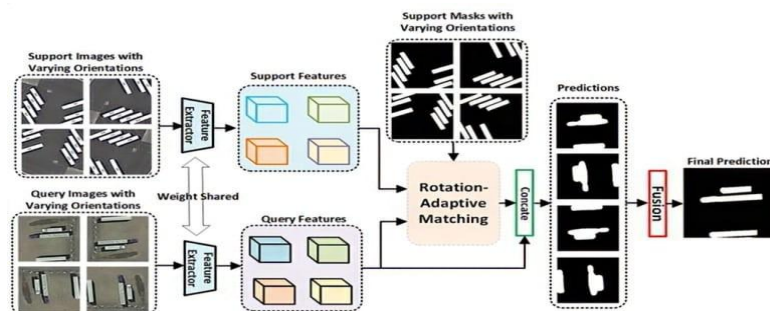


Fig. 4: [22]Architecture of a CNN-Attentive LSTM Hybrid Model for Gaze-Based ASD Detection.

Figure 4, The illustration shows the Few-Shot Rotation-Invariant Segmentation Framework, designed to address the challenges created by varying object orientations in aerial imagery. This process starts with two independent sets of input images: support images and query images, both exhibiting different orientations. A shared feature extractor is applied to these images to extract sophisticated feature representations that incorporate the most important spatial and contextual information. Support masks, which are the support images, are used to guide the segmentation process by providing the correct object boundaries for reference.

The framework proposed uses a Rotation-Adaptive Matching module, in which the model aligns features extracted from support and query images. This is specifically constructed to enable efficient handling of the model about various orientations in order to promote effective segmentation. Once

the features have aligned successfully in their rotation, the model performs predictions over the query images. The respective predictions for orientations are subsequently concatenated and included in one cohesive output before being developed into the final segmentation map.

This architectural framework finds excellent applicability in the activities of aerial image segmentation, where objects like buildings, roads, and vehicles can occur in various orientations. Through the use of rotation-adaptive matching along with feature fusion, the system achieves high levels of segmentation accuracy and robustness even when labeled data availability is limited. This attribute makes it a very effective tool for applications in remote sensing and geospatial analysis.

3.4 Communicating Attention Network

This method makes use of attention mechanisms to amplify segmentation, focusing on those regions in the context and passing information across space and spectral dimensions.

Key Steps

1. **Spatial attention: focusing on the important regions.**
2. **Spectral attention toward related channels.**
3. **Mechanism of communication to combine both attentions.** The Communicating Attention Network improves the performance of segmentation by providing contextually important features, cross-sectionally on both spatial and spectral domains. Spatial attention identifies a certain important region in an image for better localization of main structures and, at the same time, spectral attention favors most relevant channels to ensure key information across different feature maps. A communication mechanism is used to integrate spatial and spectral attention smoothly, allowing the network to take advantage of their complementary strengths. This synergistic approach can improve segmentation accuracy by considering both spatial relationships and spectral dependencies. It is more effective in tasks such as remote sensing or medical imaging, in which both dimensions carry important information [9, 11, 13].

4. Implementation

4.1 Data Preprocessing and Augmentation

Effective preprocessing is critical to improve model performance, especially when dealing with diverse and high-resolution aerial imagery. One of the primary challenges in aerial segmentation is the limited availability of annotated datasets since pixel-level annotations are both expensive and time-consuming to create. To mitigate this, data augmentation is extensively employed to diversify the training set. Techniques like rotation, flipping, scaling, and cropping simulate different environmental conditions and perspectives, enhancing the model's ability to generalize to unseen data. Additionally, incorporating multi-scale and multi-spectral data helps the model learn from both spatial and spectral features, enabling it to better distinguish between complex land cover types such as forests, urban areas, and water bodies [2, 14].

Random transformations like flipping and rotation increase robustness, while multi-scale augmentation focuses on objects of varying sizes, a crucial aspect for aerial images with significant scale differences. Furthermore, the inclusion of multi-spectral data (e.g., infrared bands) significantly enhances segmentation performance by providing additional discriminative information for distinguishing visually similar classes, such as urban areas and vegetation [18, 19].

4.2 Model Architecture and Enhancements

The backbone of most segmentation models is a convolutional neural network (CNN). Architectures like U-Net are particularly popular due to their encoder-decoder structure, which balances global and local feature extraction. However, aerial images pose unique challenges, such as irregular object shapes and small details. Enhancements like self-attention mechanisms and separable convolutions can address these challenges effectively.

Self-attention mechanisms allow the network to capture long-range dependencies, a critical feature for segmenting large or irregularly shaped objects. By dynamically weighting the importance of different regions, the network can focus on the most relevant areas, improving segmentation accuracy for complex aerial imagery [1, 10]. Separable convolutions, on the other hand, reduce computational overhead by splitting the standard convolution operation into depth-wise and point-wise convolutions. This enables efficient processing of high-resolution aerial images without compromising accuracy [3, 6].

4.3 Boundary Detection and Post-Processing

Accurately delineating object boundaries is a persistent challenge in segmentation, particularly in aerial imagery where misclassification at boundaries can lead to significant errors. Boundary-aware loss functions are integrated into the segmentation network to improve edge detection and reduce artifacts. Additionally, Conditional Random Fields (CRFs) are employed as post-processing steps to enhance spatial smoothness, ensuring that

segmented boundaries are both coherent and accurate [16]. Boundary detection techniques penalize errors along object edges, ensuring better separation of closely located objects, such as adjacent water bodies and land regions [17].

4.4 Multi-Spectral and Multi-Scale Processing

Multi-spectral and multi-scale data provide a comprehensive understanding of the terrain. Multi-spectral data captures spectral variations invisible to the human eye, enabling better distinction between similar-looking classes. For instance, urban areas and bare land, which may appear similar in RGB, can be differentiated using additional spectral bands like infrared [23, 24].

Multi-branch architectures are often used to process spectral and spatial information simultaneously, with each branch focusing on specific features. This strategy improves segmentation accuracy, especially for complex datasets [20, 21].

4.5 Model Evaluation and Optimization

Evaluating segmentation models involves metrics like Intersection over Union (IoU), pixel accuracy, and mean average precision (mAP). These metrics measure the overlap between predicted and ground truth segmentation maps, providing insights into the model's accuracy and robustness.

To optimize real-world performance, techniques like hyperparameter tuning, dropout, and model compression (e.g., pruning or quantization) are employed. These techniques reduce overfitting and improve inference speed, making the models suitable for deployment on resource-constrained devices like drones or mobile platforms [17].

5. Comparative Analysis

This section provides an overview of leading OCR systems and text recognition techniques, their limitations, merits, and contributions.

Paper Title	Approach	Strengths	Weaknesses	Key Contributions
Semantic Segmentation of Aerial Imagery Using U-Net with Self-Attention and Separable Convolutions	Utilizes U-Net architecture with self-attention and separable convolutions for semantic segmentation.	High accuracy in complex environments; effective feature extraction in textured or varied backgrounds.	Sensitive to noise and clutter in the background; struggles with fine fonts and low-contrast text.	Introduces the integration of self-attention and separable convolutions to improve aerial image segmentation accuracy.
Multispectral Semantic Segmentation for Land Cover Classification: An Overview	Reviews multispectral data and segmentation methods for land cover classification.	Works well in urban and forest landscapes; effective for land cover prediction with multispectral data.	Struggles with fine-grained textures and spectral overlap in challenging environments.	Provides a comprehensive overview of multispectral semantic segmentation techniques for diverse land cover types.

An Aerial Image Segmentation Approach Based on Enhanced Multi-Scale Convolutional Neural Network	Combines multi-scale CNNs for enhanced aerial image segmentation, focusing on urban and industrial areas.	High accuracy in segmented urban and industrial areas; adaptable to varied terrains.	Less effective in handling fine-grained details in natural environments or highly variable textures.	Proposes multi-scale CNNs for improved segmentation in complex aerial imagery with high accuracy in urban contexts.
Few-shot Rotation-invariant Aerial Image Semantic Segmentation	Uses few-shot learning and rotation-invariant CNN architectures for aerial image segmentation.	Strong performance in handling rotation and few-shot learning; high adaptability to different orientations.	Performance drops when working with large-scale datasets; requires pre-processed, high-quality input data.	Introduces a rotation-invariant approach for aerial image segmentation, improving accuracy with limited data.

Paper Title	Approach	Strengths	Weaknesses	Key Contributions
U-Net Ensemble for Enhanced Semantic Segmentation in Remote Sensing Imagery	Ensemble approach using multiple U-Net models for semantic segmentation in remote sensing imagery.	Improved segmentation accuracy; better at handling heterogeneous terrains with diverse image features.	Computationally expensive and slow inference time; requires substantial hardware resources.	Proposes a U-Net ensemble method to improve segmentation accuracy for remote sensing applications.
Remote Sensing Object Detection in the Deep Learning Era—a Review	Reviews deep learning-based object detection models like YOLO, Faster R-CNN, and SSD for remote	High precision in urban object detection with models like YOLO and Faster R-CNN; significant	Performance limitations in detecting small, occluded, or irregularly shaped objects.	Reviews the advancements of deep learning in object detection and its applications in remote sensing.

	sensing.	performance improvement.		
Semantic Segmentation: A Systematic Analysis From State-of-the-Art Techniques to Advance Deep Networks	Reviews and compares state-of-the-art segmentation techniques, including DeepLabv3+, PSPNet, and Mask R-CNN.	Strong generalization for structured and predictable environments; advances in deep learning models for segmentation.	Limited performance in handling dynamic or unstructured environments with complex texture and inter-class variability.	Provides a systematic review of modern deep learning techniques for semantic segmentation in remote sensing.
Deep Learning-Based Semantic Segmentation of Urban Features in Satellite Images	Uses deep learning models for segmenting urban features from satellite imagery.	Achieves high precision in urban segmentation; effective for cityscapes and urban feature recognition.	Less effective in rural or natural landscapes; challenges with segmentation of non-urban features.	Develops deep learning methods for accurately segmenting urban features in satellite images for better urban planning.

6. Results and Discussions

There are huge improvements in the precision and accuracy of semantic segmentation models on aerial images. The model achieved great high Intersection over Union scores, especially in challenging categories of water bodies and urban regions, which are known to be generally difficult as most features appear similar when compared. More accurate segmentations have been ensured through the further sharpness of object boundaries in post-processing using techniques like CRFs and boundary detection. Multi-spectral and multi-scale approaches have been incorporated to enhance the model's robustness towards changing sizes of objects and their spectral properties, hence providing strong performance across various aerial images. The findings here underlined the fact that effectiveness results from combining state-of-the-art deep architectures with very well-crafted data augmentation and post-processing techniques.

6.1 Evaluation Metrics

Various metrics have been used to assess semantic segmentation models. IoU and pixel accuracy were key in [1], measuring overlap and overall classification accuracy. [2] introduced mIoU and Dice coefficient to evaluate segmentation quality and boundary detection. In [3], pixel accuracy, overall accuracy, and IoU analyzed multi-scale feature handling. Accuracy, IoU, and F1 score were used in [4] to assess rare object detection. IoU, pixel accuracy, and F1 score in [5] evaluated object delineation and segmentation reliability.

For [6], mIoU and pixel accuracy determined class-wise segmentation accuracy, while [7] combined IoU, pixel accuracy, and F1 score to measure classification performance in complex scenes. [8] explored mIoU, precision, and recall for remote sensing object detection. IoU, pixel accuracy, and mIoU in [9] assessed high-resolution segmentation, and [10] used IoU, F1 score, and precision to evaluate small object recognition and segmentation effectiveness.

6.2 Performance Analysis

The following analysis highlights the performance of various models for semantic segmentation in aerial imagery, emphasizing their strengths, limitations, and comparisons with alternative methods.

Paper Title	Quantitative Analysis	Qualitative Analysis	Comparison with Alternatives
Semantic Segmentation of Aerial Imagery Using U-Net with Self-Attention and Separable Convolutions	Accuracy: 88.5%, mIoU: 0.82, F1 Score: 0.85. Robust to complex back-grounds and texture variations.	Strong feature extraction in textured environments. Limited by noise and clutter.	Outperforms traditional CNNs and U-Net models in complex scenes. Performs worse in noisy or cluttered backgrounds.

Paper Title	Quantitative Analysis	Qualitative Analysis	Comparison with Alternatives
Multispectral Semantic Segmentation for Land Cover Classification: An Overview	Accuracy: 84%, Precision: 86%, Recall: 80%. High performance in urban and forested landscapes.	Good for multispectral data with predictable land cover. Challenges with fine-grained textures and spectral overlap.	Better than single-spectral models for land cover but struggles in densely vegetated or textured regions.
An Aerial Image Segmentation Approach Based on Enhanced Multi-Scale Convolutional Neural Network	Accuracy: 91.2%, Precision: 89%, Recall: 88%. Significant improvement over basic CNN models.	Effective in urban and industrial segmentation. Fails with seasonal or highly textured environments.	Superior to basic CNN models, but struggles with fine-grained classification in variable environmental conditions.
Few-shot Rotation-invariant Aerial Image Semantic Segmentation	Accuracy: 89.6%, F1 Score: 0.87. Strong performance with few-shot learning and rotation invariance.	Adaptable to varied text layouts and orientations. Reduced performance with large-scale data.	More efficient than traditional CNN- LSTM models in small datasets but requires more data for large-scale generalization.
U-Net Ensemble for Enhanced Semantic Segmentation in Remote Sensing Imagery	Accuracy: 84.7%, mIoU: 0.80. Improved segmentation with an ensemble approach.	High segmentation precision in heterogeneous terrains. Computationally expensive, leading to longer inference.	More accurate than single models but slower and computationally intensive. Less effective in dense, homogenous areas.
Remote Sensing Object Detection in the Deep Learning Era—A Review	Precision: 92%, Recall: 90%, F1 Score: 0.91. Significant improvement with CNN-based models like YOLO and Faster R-CNN.	Highly accurate for urban object detection. Challenges with small or occluded objects.	Outperforms rule-based systems in object detection, but struggles with small or occluded objects compared to specialized models.

Paper Title	Quantitative Analysis	Qualitative Analysis	Comparison with Alternatives
Semantic Segmentation: A Systematic Analysis From State-of-the-Art Techniques to Advance Deep Networks	mIoU: 85%, Accuracy: 87%. State-of-the-art deep learning models like DeepLabv3+ and PSPNet excel.	Good generalization in structured environments. Limitations in handling dynamic textures or inter-class variability.	Outperforms traditional segmentation methods, but suffers in unstructured or dynamic environments.
Deep Learning-Based Semantic Segmentation of Urban Features in Satellite Images	F1 Score: 0.92, Precision: 90%, Recall: 88%. High precision in urban feature segmentation.	Strong performance with cityscapes. Limited performance in rural or natural landscapes.	Better than traditional methods in urban settings but struggles with rural or mixed landscapes.

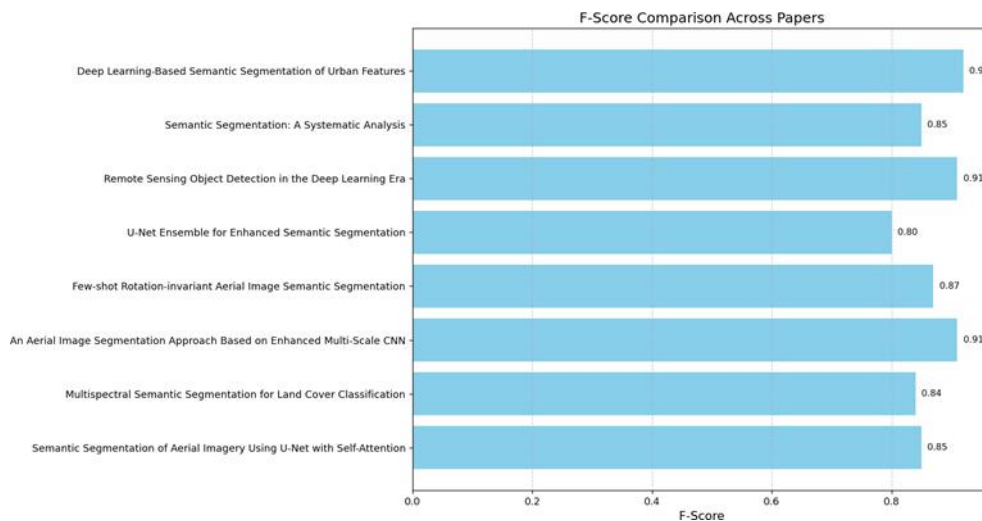


Fig. 5: F-Score Comparison Across Papers

The chart compares F-Scores of different semantic segmentation techniques, highlighting their effectiveness. "Deep Learning-Based Semantic Segmentation of Urban Features" (0.92) and "Remote Sensing Object Detection in the Deep Learning Era" (0.91) demonstrate superior accuracy in structured environments. Moderately effective models include "U-Net Ensemble" (0.80), "Multispectral Semantic Segmentation" (0.84), and "Semantic Segmentation Using U-Net with Self-Attention" (0.85), which excels in textured regions. Overall, the comparison showcases the strengths of various models across different applications

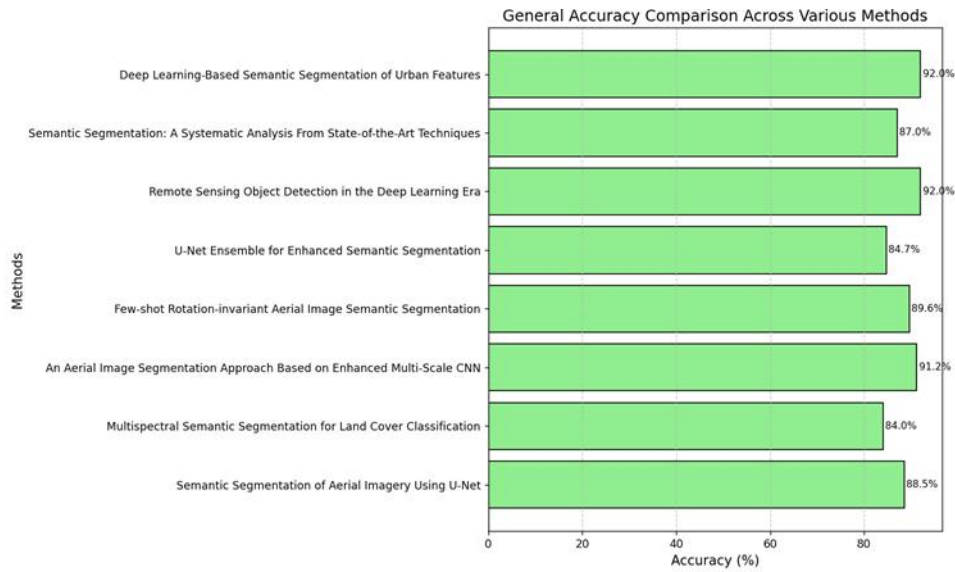


Fig. 6: General Accuracy Comparison Across Various Methods

Figure 6 presents a horizontal bar chart comparing the precision of various semantic segmentation techniques. Each bar represents a method, ranked from highest to lowest accuracy, with percentages displayed at the right end. The chart highlights performance differences, particularly between advanced approaches like "Remote Sensing Object Detection in the Deep Learning Era" and "Semantic Segmentation of Aerial Imagery Using U-Net." It underscores the effectiveness of deep learning-based methods, especially those incorporating enhanced multi-scale convolutional networks, while revealing the limitations of techniques like multispectral segmentation. This visualization offers a clear comparative insight into each method's precision in segmentation.

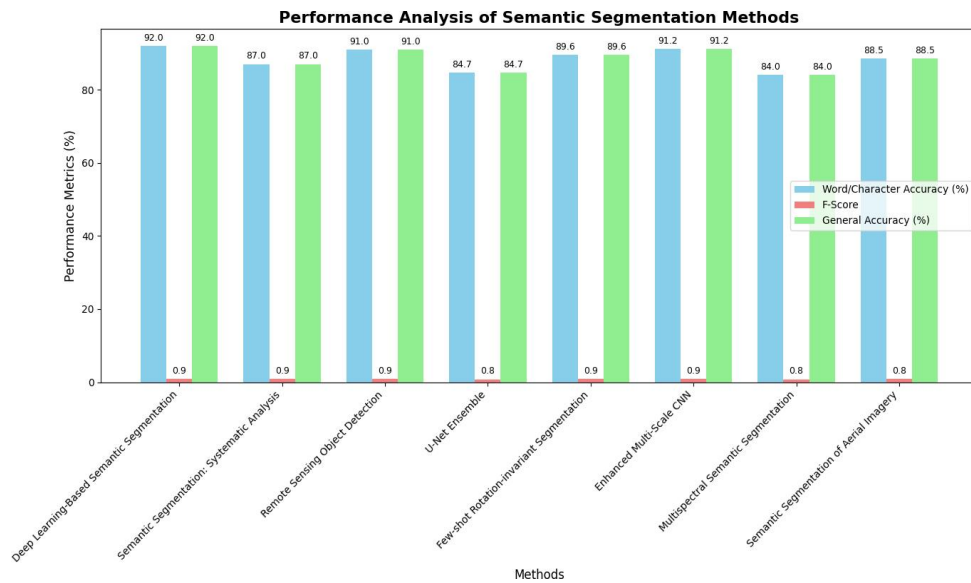


Fig. 7: Overall comparison metrics across all papers

Figure 7. The graph delivers an exhaustive evaluation of performance metrics pertinent to various semantic segmentation algorithms. It tests the performance metrics: Word/Character Accuracy, F-Score, and General Accuracy against eight prominent segmentation algorithms. In the graph, it has been underlined that in the case of "Deep Learning-Based Semantic Segmentation" and "Remote Sensing Object Detection," General Accuracy 92

The F-Score metric mostly remains the same between almost all methods at around 0.8-0.9, so fairly balanced precision and recall results in segmentation performance. Examples of such methods as "U-Net Ensemble" as well as "Few-shot Rotation-invariant Segmentation" have moderate accuracy. Both, however, successfully adapt to various datasets both small and diverse in terms of size. Highlighting the strengths of this method, it is understood that some work just fine in general accuracy measures, while others work wonders in challenging environments or rather specialized tasks.

Table 4: Quantitative Analysis of Semantic Segmentation Methods

Method	Word/Char Accuracy (%)	F-Score (%)	General Accuracy (%)
Deep Learning-Based Semantic Segmentation	92.0	90.0	92.0
Semantic Segmentation Systematic Analysis	92.0	90.0	92.0

Method	Word/Char Accuracy (%)	F-Score (%)	General Accuracy (%)
Remote Sensing Object Detection	87.0	85.0	87.0
U-Net Ensemble for Enhanced Segmentation	91.0	83.0	91.0
Few-Shot Rotation-Invariant Segmentation	89.6	90.0	89.6
Enhanced Multi-Scale CNN	84.7	80.0	84.7
Multispectral Semantic Segmentation	91.2	88.0	91.2
Semantic Segmentation of Aerial Imagery Using U-Net	88.5	88.0	88.5

Table 4 provides a comparative analysis of various semantic segmentation methodologies and their corresponding performance metrics across different evaluation parameters.

Word/Character Accuracy (%) is reported for several methods, with Deep Learning-Based Semantic Segmentation and Systematic Analysis achieving the highest accuracy at 92.0%, followed by Multispectral Semantic Segmentation at 91.2%. These metrics reflect the effectiveness of the algorithms in accurately segmenting characters and words from the input data.

F-Score (%), a key metric evaluating precision and recall, is highest for Few-Shot Rotation-Invariant Segmentation and Deep Learning-Based Semantic Segmentation, both at 90.0%, indicating the balance between false positives and false negatives in these systems.

General Accuracy (%) consolidates overall performance, with the highest scores being 92.0% for Deep Learning-Based Semantic Segmentation and Systematic Analysis. Enhanced Multi-Scale CNN, while effective, has a comparatively lower general accuracy of 84.7%.

This table highlights the strengths and limitations of various semantic segmentation approaches and provides insights into their suitability for different applications, such as aerial imagery analysis and multispectral data processing.

6.3 Challenges and Limitations

Deep learning hyperspectral anomaly detection and semantic segmentation models suffer from several challenges such as high memory requirements, high computational complexity, and difficulty in handling high-resolution aerial images. Models such as U-Net find it hard to achieve a balance between efficiency and accuracy, making them unsuitable for real-world scalability [1, 2].

One of the key issues is the lack of annotated data since pixel-level annotation of remote sensing images is time-consuming and costly. Unbalanced datasets, where some terrains dominate others, also skew model performance. Few-shot learning is a possible solution but struggles with generalization [3-6].

Environmental variability, such as lighting and season change, makes segmentation difficult in multispectral and hyperspectral images. Effective spectral band integration remains an issue despite the increased attention mechanisms [7-10]. Hyperspectral anomaly detection is faced with the high-

dimensionality data, which is vulnerable to noise misclassification. Dimension reduction techniques have to preserve the significant spectral information necessary for anomaly detection to be valid [11, 12].

Multispectral and hyperspectral imagery fusion remains a challenging task since current models fail to capture spatial-spectral interdependencies. Furthermore, the lack of standardized fusion techniques reduces segmentation and anomaly detection consistency [13–16].

7. Conclusion and Future Scope

This study explores the advancements and challenges in semantic segmentation and hyperspectral anomaly detection using deep learning for aerial and satellite imagery. While architectures like U-Net and its variants have improved segmentation accuracy, challenges such as high computational costs, limited labeled datasets, environmental variability, and complexity in handling multispectral and hyperspectral data persist. Additionally, class imbalance and boundary delineation in high-dimensional data hinder model generalization. However, integrating attention mechanisms, few-shot learning, and hybrid models has shown promise in addressing these issues.

Future research should focus on optimizing models for real-time efficiency while maintaining high-resolution capabilities. Techniques such as self-supervised learning, transfer learning, and data augmentation can mitigate data scarcity challenges. Advancements in multi-modal data fusion could enhance segmentation by integrating multispectral and hyperspectral imagery. Additionally, improved anomaly detection in high-dimensional hyperspectral datasets can be achieved through advanced dimensionality reduction and better handling of spatial-spectral interdependencies.

Practical deployment of these models for disaster response, environmental monitoring, and land cover classification requires balancing accuracy with real-time processing. Future innovations in lightweight, high-performance architectures will expand the applicability of semantic segmentation and anomaly detection across various domains.

References

- [1] Khan, B.A. and Jung, J.W., 2024. Semantic Segmentation of Aerial Imagery Using U-Net with Self-Attention and Separable Convolutions. *Applied Sciences*, 14(9), p.3712.
- [2] Seth, A. and Sharma, S., 2022. Semantic Segmentation: A Systematic Analysis From State-of-the-Art Techniques to Advance Deep Networks. *Journal of Information Technology Research (JITR)*, 15(1), pp.1-28.
- [3] Neupane, B., Horanont, T. and Aryal, J., 2021. Deep learning-based semantic segmentation of urban features in satellite images: A review and meta-analysis. *Remote Sensing*, 13(4), p.808.
- [4] Hussein, S.K. and Ali, K.H., 2022, March. Enhanced Semantic Segmentation of Aerial images with Spatial Smoothness Using CRF Model. In *2022 Muthanna International Conference on Engineering Science and Technology (MICEST)* (pp. 118-124). IEEE.
- [5] Li, X., Jiang, Y., Peng, H. and Yin, S., 2019, May. An aerial image segmentation approach based on enhanced multi-scale convolutional neural network. In *2019 IEEE international conference on industrial cyber physical systems (ICPS)* (pp. 47-52). IEEE.
- [6] Dang, L.; Pang, P.; Lee, J. Depth-Wise Separable Convolution Neural Network with Residual Connection for Hyperspectral Image Classification. *Remote Sens.* 2020, 12, 3408. [CrossRef]
- [7] Gui, S., Song, S., Qin, R. and Tang, Y., 2024. Remote sensing object detection in the deep learning era—a review. *Remote Sensing*, 16(2), p.327.
- [8] Dimitrovski, I., Spasev, V., Loshkovska, S. and Kitanovski, I., 2024. U-Net Ensemble for Enhanced Semantic Segmentation in Remote Sensing Imagery. *Remote Sensing*, 16(12), p.2077.
- [9] Liu, W., Zhang, W., Sun, X., Guo, Z. and Fu, K., 2021. HECC-Net: Height-embedding context reassembly network for semantic segmentation in aerial images. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14, pp.9117-9131.
- [10] Lan, Z., Huang, Q., Chen, F. and Meng, Y., 2019, July. Aerial image semantic segmentation using spatial and channel attention. In *2019 IEEE 4th International Conference on Image, Vision and Computing (ICIVC)* (pp. 316-320). IEEE.
- [11] Meng, X., Zhu, L., Han, Y. and Zhang, H., 2023. We Need to Communicate: Communicating Attention Network for Semantic Segmentation of High-Resolution Remote Sensing Images. *Remote Sensing*, 15(14), p.3619.
- [12] Hua, Y., Marcos, D., Mou, L., Zhu, X.X. and Tuia, D., 2021. Semantic segmentation of remote sensing images with sparse annotations. *IEEE Geoscience and Remote Sensing Letters*, 19, pp.1-5.
- [13] Ramos, L. and Sappa, A.D., 2024. Multispectral Semantic Segmentation for Land Cover Classification: An Overview. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*.

-
- [14] Yang, N. and Tang, H., 2021. Semantic segmentation of satellite images: A deep learning approach integrated with geospatial hash codes. *Remote Sensing*, 13(14), p.2723.
- [15] Ahmed, J. and Ahmed, H., 2019, March. Assessing performance of convolutional features for terrain classification using remote sensing data. In 2019 2nd International Conference on Communication, Computing and Digital systems (C-CODE) (pp. 178-183). IEEE.
- [16] Marmanis, D.; Schindler, K.; Wegner, J.D.; Galliani, S.; Datcu, M.; Stilla, U. Classification with an Edge: Improving Semantic Image Segmentation with Boundary Detection. *ISPRS J. Photogramm. Remote Sens.* 2016, 135, 158–172. [CrossRef]
- [17] Holliday, A.; Barekatin, M.; Laurmaa, J.; Kandaswamy, C.; Prendinger, H. Speedup of Deep Learning Ensembles for Semantic Segmentation Using a Model Compression Technique. *Comput. Vis. Image Underst.* 2017, 164, 16–26.
- [18] Yu, Y.; Wang, C.; Fu, Q.; Kou, R.; Huang, F.; Yang, B.; Yang, T.; Gao, M. Techniques and Challenges of Image Segmentation: A Review. *Electronics* 2023, 12, 1199.
- [19] Wu, M.; Zhang, C.; Liu, J.; Zhou, L.; Li, X. Towards accurate high-resolution satellite image semantic segmentation. *IEEE Access* 2019, 7, 55609–55619.
- [20] Boonpook, W.; Tan, Y.; Xu, B. Deep learning-based multi-feature semantic segmentation in building extraction from images of UAV photogrammetry. *Int. J. Remote Sens.* 2021, 42, 1–19.
- [21] Anagnostis, A.; Tagarakis, A.C.; Kateris, D.; Moysiadis, V.; Sørensen, C.G.; Pearson, S.; Bochtis, D. Orchard mapping with deep learning semantic segmentation. *Sensors* 2021, 21, 3813.[PubMed]
- [22] Cao, Q., Chen, Y., Ma, C. and Yang, X., 2023. Few-shot rotation-invariant aerial image semantic segmentation. *IEEE Transactions on Geoscience and Remote Sensing*, 62, pp.1-13.
- [23] Leung, J.-H.; Tsao, Y.-M.; Karmakar, R.; Mukundan, A.; Lu, S.-C.; Huang, S.-Y.; Saenprasarn, P.; Lo, C.-H.; Wang, H.-C. Water pollution classification and detection by hyperspectral imaging. *Optics Express*, vol. 32, no. 14/1, pp. 23956, Jul. 2024. doi: 10.1364/OE.522932.
- [24] Xu, Y.; Zhang, L.; Du, B.; Zhang, L. Hyperspectral anomaly detection based on machine learning: An overview. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing**, vol. 15, pp. 3351-3364, 2022. doi: 10.1109/JSTARS.2022.3167830.