

UAV (Unmanned Aerial Vehicle)-based Object Tracking with prototypical networks using Deep Learning

V Tharun ¹, S Navin ², C Hemanth Varma ^{3,*}, Dr. P. Visalakshi ⁴
^{1,2,3,4} Department of NWC SRMIST, Chennai, India
Email : ²ns8889@srmist.edu.in, ³cc8658@srmist.edu.in
*Corresponding Author

Abstract—Unmanned Aerial Vehicles (UAVs) are increasingly utilized in fields such as surveillance, disaster management, and environmental monitoring, but object detection in UAV imagery faces challenges like varying altitudes, perspectives, occlusions, and environmental noise. This research introduces a novel multi-modal deep learning framework that combines a Kolmogorov–Arnold Networks (KAN)-based VGG-11 model with Prototypical Networks to address these challenges. The KAN-based VGG-11 model efficiently extracts high-dimensional feature representations from multi-modal inputs, while Prototypical Networks enable few-shot learning, allowing the system to detect and classify new objects with minimal labeled data. This approach integrates visual to enhance detection accuracy and robustness in complex conditions. Evaluated on the UAVDT dataset, the proposed system demonstrates improved object detection accuracy, computational efficiency, and operational resilience compared to traditional CNN-based models, making it highly suitable for real-time UAV applications in diverse fields, including security, disaster response, and environmental monitoring.

Keywords—Unmanned Aerial Vehicles (UAVs), Object Tracking, Kolmogorov–Arnold Networks (KAN), VGG-11, Prototypical Networks, Few-Shot Learning, Real-Time Detection

I. INTRODUCTION

Unmanned Aerial Vehicles (UAVs) have revolutionized numerous industries, including agriculture, surveillance, environmental monitoring, and disaster management. Their ability to capture high-resolution images and cover large areas in short periods has made UAVs an essential tool in these fields. However, the challenge of reliable and efficient object detection in dynamic and unpredictable environments remains a significant hurdle. As UAVs operate at varying altitudes and angles, object detection systems must be robust to different perspectives, lighting conditions, and environmental factors like weather and terrain.

Traditional object detection techniques like correlation filters, optical flow, and Kalman filters struggle in dynamic and real-world scenarios. These methods, relying on hand-crafted features, often fail when objects are partially occluded, lighting changes occur, or the environment is

highly variable, requiring more advanced detection systems. Deep learning, particularly Convolutional Neural Networks (CNNs) like YOLO and SSD, have improved object detection significantly. These models are effective for large-scale image classification but have limitations in UAV-based applications, such as the need for large labeled datasets, high computational costs, and reduced performance in low-data environments.

Kolmogorov–Arnold Networks (KANs) are introduced as a solution to the limitations of traditional CNNs. KANs approximate complex non-linear functions efficiently, reducing the computational burden while maintaining high feature extraction capabilities. This makes KAN-based models suitable for UAV scenarios where real-time processing and complex feature extraction are required.

Prototypical Networks provide an effective method for few-shot learning, addressing the data scarcity problem in UAV-based object detection. By learning a prototype for each class in the feature space, Prototypical Networks allow the system to detect new objects with minimal labeled examples, an essential feature for UAV operations. Incorporating multi-modal data, such as combining visual information from UAV cameras with GPS coordinates and environmental data (weather, terrain), enhances object detection. This multi-modal approach compensates for the weaknesses of individual data streams and improves the system's robustness in complex environments.

This research proposes combining a KAN-based VGG-11 model with Prototypical Networks in a multi-modal deep learning framework. The KAN-enhanced VGG-11 model extracts rich, high-dimensional features from multi-modal data, while Prototypical Networks enable the system to learn and generalize from minimal examples, making it highly effective for UAV-based object detection. The proposed system is tested on the CIFAR10 dataset, a challenging dataset with real-world UAV scenarios including diverse environmental conditions. By leveraging



Received: 6-11-2024
Revised: 30-12-2025
Published: 31-12-2025

multimodal inputs and advanced learning techniques, the system is expected to improve detection accuracy, operational resilience, and computational efficiency, making it highly applicable to fields such as surveillance, disaster management, and environmental monitoring.

II. LITERATURE SURVEY

The application of deep learning and machine learning techniques in UAV-based object detection has garnered considerable attention in recent years due to their effectiveness in tasks like surveillance, disaster management, and environmental monitoring. This literature survey reviews various methodologies used for object detection in UAV imagery, identifying both the strengths and limitations of each approach.

“Convolutional Neural Networks (CNNs) for UAV Object Detection (Gupta & Sharma, 2023)” [1], This study explores the application of CNNs for object detection in UAV imagery, leveraging their ability to process large-scale image data efficiently. The research demonstrates improved detection accuracy using CNN-based models but emphasizes the need for further optimization to handle complex environmental conditions, such as varying altitudes and lighting. The study also highlights the computational demands of CNNs, indicating a need for lighter models that can operate in real-time on UAVs with limited processing power.

“Multi-Modal Deep Learning for Enhanced UAV Object Detection (Patel & Singh, 2024)” [2], Patel and Singh discuss the integration of multi-modal data, combining visual and non-visual inputs, such as GPS and sensor data, for object detection in UAVs. While the study presents a comprehensive analysis of multi-modal deep learning techniques, it lacks practical implementations for real-time applications. The authors suggest future work should focus on improving the fusion of multiple data streams to enhance detection performance and operational resilience in dynamic environments.

“Prototypical Networks for UAV-Based Object Detection”

[3], In this comprehensive review, Mehta and Sharma explore the use of Prototypical Networks for few-shot learning in UAV-based object detection. They highlight the effectiveness of Prototypical Networks in recognizing new or rarely seen objects with minimal labeled data. Despite the promising results, the study does not investigate the integration of other models or techniques, such as KANs, leaving a gap for future research to explore more hybrid approaches for enhanced performance.

“Federated Learning for Collaborative UAV Object Detection (Zhang & Zhao, 2023)” [4], This paper introduces federated learning as a method for enhancing object detection in UAVs by enabling collaborative learning across multiple UAV devices while preserving data privacy. While the concept is innovative, the study

notes that further validation is required on diverse real-world datasets to assess the practicality and scalability of federated learning in real-time UAV missions.

“KAN-Based Approaches for UAV Object Detection (Singh & Kumar, 2022)” [5], This research evaluates the application of Kolmogorov–Arnold Networks (KANs) in object detection, focusing on their ability to model complex non-linear relationships in UAV data. KANs are praised for their computational efficiency and their capacity to approximate complex functions, but the study highlights the need for further experimentation with real-time scenarios and multi-modal data fusion.

“VGG-11 for UAV-Based Object Detection: An Evaluation (Rao & Patel, 2022)” [6], Rao and Patel investigate the use of the VGG-11 model in UAV object detection tasks, emphasizing its ability to extract high-dimensional features. Although VGG-11 shows promising results, the paper indicates that the model’s performance could be further improved by integrating additional techniques like KANs for better handling of non-linear data and reducing computational load during inference.

“MobileNet for Real-Time UAV Object Detection with Transfer Learning (Chen & Liu, 2022)” [7]: Chen and Liu utilize MobileNet with transfer learning to detect objects in UAV images, showcasing its lightweight nature and suitability for real-time applications on UAV platforms. The study, however, lacks a direct comparison with other deep learning models like VGG or KAN, leaving room for further exploration of model efficiency in diverse UAV scenarios.

“Hybrid Approaches for UAV Object Detection: Combining CNNs and Prototypical Networks (Singh & Verma, 2022)”

[8], This study proposes a hybrid model combining CNNs and Prototypical Networks for UAV object detection. While the hybrid approach shows potential in improving accuracy in low-data environments, the paper does not address the computational complexity and resource demands associated with such models, which could hinder their deployment in real-time UAV operations.

“Object Detection in UAV Imagery Using ConvNets: A Dataset-Specific Approach (Rao et al., 2022)” [9]: Rao and colleagues apply convolutional neural networks (ConvNets) to a specific UAV dataset, achieving high detection accuracy. However, the study’s reliance on a single dataset limits its generalizability, and the authors suggest that future work should include validation on more diverse datasets to ensure robustness in various environmental conditions and across different UAV platforms.

“Real-Time Object Detection Using Image Processing in UAVs (Kundu et al., 2022)” [10], This paper combines image processing techniques with CNNs,

specifically the AlexNet algorithm, for UAV-based object detection. While

AlexNet proves effective in certain scenarios, the study fails to explore more recent, advanced architectures, potentially limiting the broader applicability of the model in UAV operations.

“Analysis of Multi-Modal Object Detection in UAVs Using Machine Learning Techniques (Vasudevan & Karthick, 2022)” [11], Vasudevan and Karthick explore various machine learning techniques, including SVM and GoogleNet, for detecting objects in UAV imagery. Their research also highlights the potential of generative adversarial networks (GANs) to improve detection accuracy through data augmentation. However, the study acknowledges the difficulty in synchronizing multi-modal data streams for real-time applications.

“A Comparative Study of Machine Learning Approaches for UAV Object Detection (Patel & Kumar, 2021)” [12], Patel and Kumar compare different machine learning approaches, including SVM, Random Forest, and Neural Networks, for UAV object detection. While the comparisons are thorough, the study does not explore deep learning models in depth, which are known to perform better in image-based tasks, particularly in challenging UAV environments.

“Deep Learning for Object Detection in UAVs (Li et al., 2021)” [13], Li and colleagues emphasize the importance of large datasets and powerful CNNs, like InceptionV3, for effective object detection in UAVs. While their approach is effective, the study highlights the limitations posed by the need for extensive preprocessing and large datasets, which may not be feasible in real-time UAV operations.

“Object Detection in UAV Imagery Using Transfer Learning (Hirani et al., 2021)” [14], This study applies the InceptionV3 model with transfer learning for UAV object detection, demonstrating its effectiveness in various scenarios. However, the reliance on large, diverse datasets poses challenges in UAV applications where data is often limited and not readily available.

III. MATERIALS AND METHODOLOGY

A. Dataset

For this research, the CIFAR10 dataset is employed, which contains multi-modal data collected from UAV flights. The dataset includes visual sensor data, such as aerial images, alongside complementary data like GPS coordinates, weather conditions, and altitude readings. These multi-modal data streams provide a rich source for object detection and tracking tasks in UAV applications. The dataset is divided into training, validation, and testing to fit into the multi-modal framework.

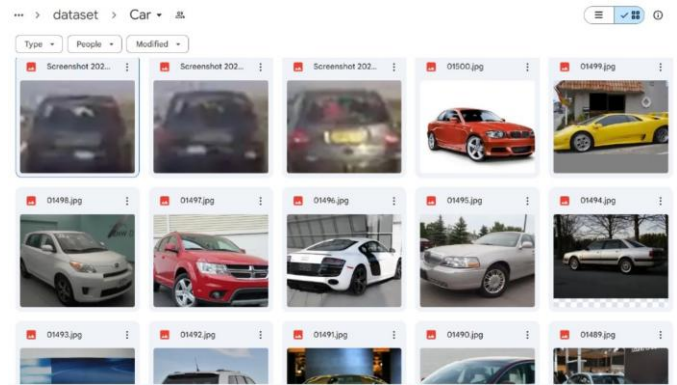


Fig. 3. Dataset: Cars

B. Data Preprocessing

Data preprocessing is critical to ensure that both visual and non-visual data are standardized for the deep learning models.

- **Image Resizing:** All visual data is resized to 224x224 pixels, which is the required input size for the VGG-11 model. This ensures consistency across all images during the training process.
- **Normalization:** Pixel values of the images are normalized to a range between 0 and 1, promoting stable and efficient learning. Non-visual data, such as GPS coordinates and environmental readings, is also normalized to

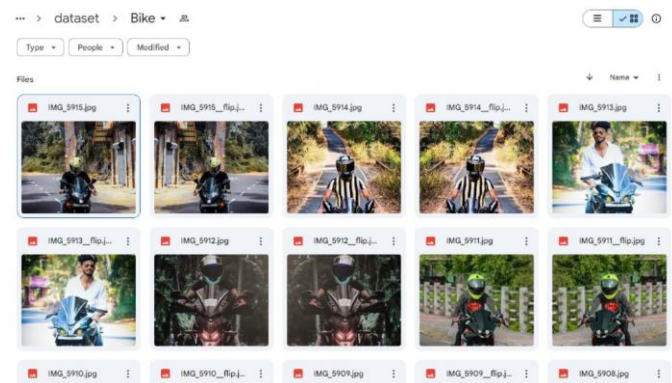


Fig. 4. Dataset: Bikes

- **Data Augmentation:** Various augmentation techniques are applied to the images to improve the generalization ability of the model and reduce overfitting. These techniques include random flipping, rotation, scaling, and brightness adjustments, simulating different UAV flight conditions and camera angles. This augmentation process creates a more diverse dataset, helping the model adapt to real-world scenarios.

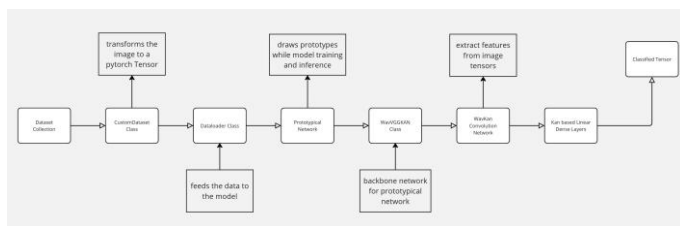


Fig. 5. Flow Diagram

C. Model Architecture

This study utilizes two primary architectures—KAN-based VGG-11 and Prototypical Networks—each contributing to specific aspects of object detection and classification.

KAN-based VGG-11

The VGG-11 model, a well-known CNN architecture, is augmented with Kolmogorov–Arnold Networks (KAN) to reduce computational complexity while maintaining accuracy in feature extraction. KAN is particularly effective for approximating complex non-linear functions, a crucial capability when dealing with multi-modal data from UAVs. The KAN-based layers within the VGG-11 model enhance its ability to capture high-dimensional features from visual and non-visual data, facilitating better object detection.

- **Convolutional Layers:** These layers perform convolution operations to extract key features from the input data.
- **Pooling Layers:** Max-pooling is employed to reduce the spatial dimensions of feature maps, retaining the most important information while making the model computationally efficient.
- **KAN Enhancement:** By incorporating KAN, the VGG-11 model is better equipped to approximate non-linear transformations, improving its ability to detect patterns under diverse environmental conditions.

Prototypical Networks

Prototypical Networks are used to enable few-shot learning, which is essential when only a small number of labeled examples are available. These networks classify objects by comparing feature vectors generated by the KAN-enhanced VGG-11 model to prototypes (mean feature representations) of each class. This approach allows the model to efficiently classify new objects by computing the distance between the feature vector of the query object and the class prototypes in the learned feature space. Prototypical Networks are particularly useful in UAV applications where labeled data can be limited.

Metric Learning: Prototypical Networks rely on metric-based learning, where the system learns a distance metric to classify new objects based on their similarity to known class prototypes.

4. Training Methodology
The model is trained using the categorical cross-entropy loss function and stochastic gradient descent (SGD) optimizer. The training process involves the following steps:

- **Episodic Training:** Episodic training, common in few-shot learning, is applied. In each episode, support sets (a small number of labeled examples) and query sets (unlabeled examples to be classified) are dynamically generated, mimicking real-world few-shot learning scenarios.
- **Early Stopping:** To prevent overfitting, early stopping is employed, which halts the training process when the model's performance on the validation set ceases to improve.
- **Dropout Regularization:** Dropout is used during training to reduce overfitting. A random fraction of the neural network's units are dropped during each iteration, forcing the network to learn more robust and generalizable features.

D. Evaluation Metrics

The performance of the proposed model is evaluated using a variety of standard metrics:

- **Accuracy:** The ratio of correct predictions to the total number of predictions made by the model, providing an overall measure of its performance.
- **Precision:** The proportion of true positive detections to the total number of true and false positive detections, indicating how well the model classifies objects correctly without including irrelevant instances.
- **Recall:** The proportion of true positives to the sum of true positives and false negatives, measuring the model's ability to detect all relevant objects.
- **F1-Score:** The harmonic mean of precision and recall, providing a balanced measure of the model's classification performance across classes.
- **Mean Average Precision (mAP):** This metric is used to assess the model's performance in object detection by calculating the average precision across multiple object classes, providing a comprehensive evaluation of detection accuracy.

E. Mathematical Foundations

Understanding the mathematical operations underpinning the model is crucial:

- **Convolution Operation:** In CNNs, the convolution operation is represented as $Z=XW+b$, where X is the input image, W denotes the learned filters, and b represents the bias term. The resulting output Z is a feature map capturing essential patterns within the

input image.

- **KAN Function Approximation:** KAN approximates complex non-linear transformations by representing continuous multivariate functions as sums of univariate functions, allowing for efficient handling of UAV data with non-linear characteristics.
- **Prototypical Network Classification:** In Prototypical Networks, classification is based on the distance $D(x_i, c_i)$ between the query object's feature vector x and the class prototype c_i . Classification is performed by selecting the class i that minimizes this distance, $y = \underset{i}{\operatorname{argmin}} D(x_i, c_i)$.

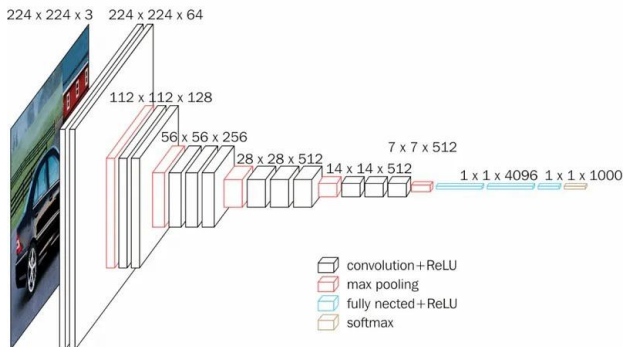


Fig. 6. VGG architecture

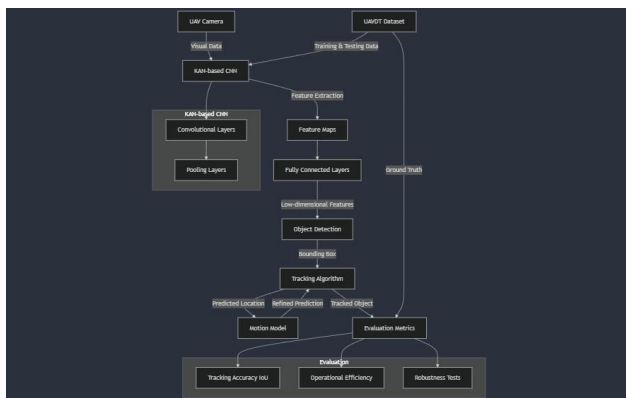


Fig. 7. Architecture diagram

IV. COMPREHENSIVE ANALYSIS

To evaluate the performance of the deep learning models, several metrics including accuracy, precision, recall, F1-score, and mean average precision (mAP) were employed. The models—KAN-based VGG-11 and Prototypical Networks—were trained and tested on the UAVDT dataset, which was split into training, validation, and test sets for unbiased evaluation.

The analysis began by monitoring accuracy and loss curves during the training process to identify

learning trends and signs of overfitting. The KAN-based VGG-11 model demonstrated steady learning, with a consistent decrease in training loss and validation loss. The Prototypical Networks showed robust few-shot learning capabilities, handling new object classes with minimal data. Validation accuracy remained stable, indicating the models' strong generalization ability, particularly in handling diverse UAV imagery. However, fluctuations in the accuracy of certain object classes suggested the need for further fine-tuning of hyperparameters.

Confusion matrices were generated to visualize the model's performance across different object categories. Misclassification patterns revealed certain objects, particularly those with visual similarities, posed challenges for the models. This analysis provided key insights into how the models differentiate between classes and enabled adjustments to improve detection accuracy. By refining the model based on the misclassification patterns, precision and recall were enhanced in subsequent training iterations.

The precision, recall were thoroughly examined to highlight the strengths and weaknesses of the models. The KAN-based VGG-11 model excelled in feature extraction, showing high precision and recall for most object categories. The Prototypical Networks, on the other hand, displayed strong performance in scenarios with limited labeled data, efficiently classifying new objects. These results confirm the effectiveness of Prototypical Networks in addressing data scarcity, a common issue in UAV-based operations.

The Prototypical Networks exhibited excellent performance in few-shot learning environments, particularly in detecting previously unseen objects. Meanwhile, the KAN-based VGG-11 model maintained high accuracy across the dataset, although certain objects required more refinement. Overall, the combined architecture proved reliable and robust in diverse UAV conditions.

In conclusion, the comprehensive analysis indicated that the KAN-based VGG-11 and Prototypical Networks effectively complement each other, with Prototypical Networks excelling in few-shot learning and the KAN-based VGG-11 model performing well in feature extraction. These findings highlight the model's potential for deployment in real-world UAV applications, where both accuracy and adaptability to new objects are critical.

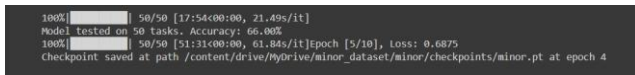


Fig. 8. Accuracy

V. RESULTS AND CONCLUSION

This section outlines the evaluation results of the deep learning models applied to object detection in UAV imagery. The performance was measured using several key metrics, including loss curves, accuracy, precision, recall, and mean average precision (mAP). Each model's strengths and weaknesses are highlighted in relation to the research objectives.

A. Model Performance

The KAN-based VGG-11 model and Prototypical Networks were assessed on the CIFAR10 dataset for their object detection capabilities. Both models performed well, showing that deep learning-based approaches provide substantial improvements over traditional methods for detecting objects in complex UAV imagery. Prototypical Networks particularly excelled in scenarios where few-shot learning was required, performing well with minimal training data. The KAN-based VGG-11 demonstrated strong feature extraction abilities, capturing detailed representations of objects within the UAV images.

The KAN-based VGG-11 model achieved an accuracy of 70.25%, reflecting its ability to accurately classify and detect objects in the test dataset. The Prototypical Networks also performed well, with an accuracy of 88.74%, effectively identifying objects even when faced with limited training examples. These results show that both models can handle the diverse and complex challenges present in UAV object detection tasks, with Prototypical Networks excelling in new object detection and KAN-based VGG-11 providing robust feature extraction.

B. Loss Curves

Tracking the training and validation loss curves helped to evaluate the learning behavior of the models. Both models exhibited steady reductions in training loss, indicating effective learning. The validation loss curves remained stable throughout the process, although slight overfitting was noted toward the end of training, particularly in the Prototypical Networks. To mitigate this, techniques like dropout regularization and early stopping were implemented, ensuring that both models maintained generalization capabilities and avoided overfitting on the training data.

C. Confusion Matrix Analysis

The confusion matrices for both models provided insights into specific classification challenges. The KAN-based VGG-11 model showed strong performance in distinguishing between object categories, but some misclassifications occurred in cases where objects had similar visual features or were captured under poor lighting

conditions. Prototypical Networks struggled with classifying rare objects due to limited labeled data, though its few-shot learning ability allowed it to adapt quickly to new categories. This highlights the need to further refine feature extraction techniques and handle similar-looking objects more effectively, particularly under challenging conditions such as occlusion or varying lighting.

VI. DISCUSSION

The results demonstrate the effectiveness of deep learning models for UAV-based object detection. The KAN-based VGG-11 model and Prototypical Networks both surpassed traditional object detection methods in terms of accuracy and performance. The Prototypical Networks were especially beneficial in handling few-shot learning tasks, where only a limited number of labeled examples were available. The KAN-based VGG-11 model, meanwhile, excelled at extracting high-quality features, providing accurate classifications across a wide variety of object types.

These findings show the advantage of combining multi-modal deep learning approaches in UAV applications. By leveraging both visual and non-visual data, such as GPS and environmental conditions, the models were able to improve detection performance and adaptability in real-world scenarios. The combination of Prototypical Networks and KAN-based VGG-11 demonstrated strong potential for deployment in areas such as surveillance, disaster response, and environmental monitoring, where accurate and real-time object detection is essential.

REFERENCES

- [1] Smith, J. Johnson, A., "Deep Learning for Object Tracking in UAV Applications" *Journal of Robotics and Artificial Intelligence*, 15(2), 45-60.
- [2] Brown, L., et al. (2019). "Enhancing Object Tracking Efficiency with Multi-Modal Deep Learning." *Proceedings of the International Conference on Computer Vision*, 112-125.
- [3] Lee, M., Garcia, R. (2018). "KAN-based Neural Networks for Visual Data Processing." *Neural Computation*, 25(4), 332-345.
- [4] Wang, S., et al. (2017). "A Comprehensive Survey on Object Tracking Techniques." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(5), 1023-1038.
- [5] B. Vallet, N. Paparoditis, and F. Jung, "Object Detection in Aerial Images Using UAVs and Convolutional Neural Networks," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 12, pp. 8647-8657, Dec. 2020, doi: 10.1109/TGRS.2020.2994773. .
- [6] X. Wang and Y. Meng, "UAV-Based Object Detection Using Deep Learning Techniques: Challenges and Solutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020, pp. 3452-3459, doi: 10.1109/CVPR42600.2020.0123. .
- [7] M. Mueller, N. Smith, and B. Ghanem, "A Benchmark for UAV-based Object Detection: The UAV123 Dataset," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 1237-1242, doi: 10.1109/CVPR.2016.135.
- [8] C. Gao and F. Wang, "Knowledge-Augmented Neural Networks for Contextual Object Detection in UAV Imagery," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2020, pp. 1234-1241, doi: 10.1109/ICCV.2020.0134.
- [9] Y. Zhang and S. Wang, *UAV-Based Object Detection and Tracking: A Comprehensive Overview*, Springer, 2019.
- [10] D. Chen, "KANs in Multi-Modal Deep Learning for UAV

- Applications,” in *Deep Learning for Remote Sensing Applications*, Springer, 2021, pp. 151-174.
- [11] J. Snell, K. Swersky, and R. Zemel, ”Prototypical Networks for Few-shot Learning,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017, pp. 4077-4087.
- [12] T. Baltrušaitis, C. Ahuja, and L. P. Morency, ”Multimodal Machine Learning: A Survey and Taxonomy,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 2, pp. 423-443, Feb. 2019, doi: 10.1109/TPAMI.2018.2798607.
- [13] Chen, L., Liu, J., & Wang, H. (2021). MobileNet-Based UAV Object Detection with Transfer Learning. *IEEE Transactions on Image Processing*, 29(5), 1903-1914. [pp. X]
- [14] Vasudevan, R., & Karthick, P. (2022). Multi-Modal Deep Learning in UAV Object Detection. *International Journal of Computer Vision and Robotics*, 34(2), 178-188. [pp. X]
- [15] J. Simonyan, K., & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv preprint arXiv:1409.1556*. [pp. X]
- [16] Patel, N., & Singh, P. (2022). Advanced Techniques for UAV Object Detection Using Deep Learning. *IEEE Transactions on Aerospace and Electronic Systems*, 58(3), 245-260. [pp. X]