

# INTELLIGENT PEST CLASSIFICATION AND REMEDY RECOMMENDATION SYSTEM USING VISION TRANSFORMERS AND DOMAIN-AWARE SUMMARIZATION

Satwika Jahnavi<sup>1</sup>, Sireesha Moturi<sup>2</sup>, Mounika Naga Bhavani Meduri<sup>3</sup>, Thrisony Gayam<sup>4</sup>, Jahnavi Kappa<sup>5</sup>

<sup>1,2,3,4,5</sup>*Department of CSE, Narasaraopeta Engineering College, Andhra Pradesh, India*

<sup>1</sup>chakkasatwika225@gmail.com, <sup>2</sup>sireeshamoturi@gmail.com, <sup>3</sup>medurimounika4@gmail.com, <sup>4</sup>gayamthrisony@gmail.com,

<sup>5</sup>jahnavidreddy2705@gmail.com

**Abstract**—Every year, farmers around the world lose up to 40% of their crops due to pests often because they can't identify the problem in time. The issue of pest identification in rice farming continues to be a challenge due to delayed pest control measures and resulting extensive damage. In addressing this challenge, we implemented a deep learning recognition system that identifies images of 40 classes of rice pests. Our model uses Vision Transformers (ViT) as they capture spatial and structural details of pests at a granular level. Instead of classifying pests like a traditional classifier, our system works as a pest assistant. After identification, the system condenses essential details such as description, recommended pesticide, and preventive measures into farmer-friendly summaries. The additional information enhances the system's usefulness for real-world farming conditions. Moreover, the user-friendly interface is designed for farmers, agricultural workers, students, and researchers. The system integrates deep learning with agricultural insight by providing accurate classification and accompanying information. FasterPest classify 14 classes of rice pests by integrating CNN and Transformer components. On the other hand, our model based on vision transformer handles 40 classes of rice pests with a more straightforward approach while providing pragmatic outputs beneficial to farmers. Our model outperforms the FasterPest baseline with precision 12%, achieving precision 96% in pest classification.

**Index Terms**—Pest Classification, Vision Transformer(ViT), Deep Learning, Timm Library, Pesticide Recommendation, Domain-Aware Summarization.

## I. INTRODUCTION

Rice, the staple food for more than half of the world's population, is a foundation of food security and agricultural economies [1]. Nevertheless, cultivated rice is very susceptible to pest attack that can reduce yields lacking desirable grain quality, let alone causing significant economic losses [2]. Farmers, particularly in resource-constrained areas, often struggle to detect and manage pest outbreaks in time, relying heavily on pesticide usage. While pesticides offer a quick remedy, their excessive and indiscriminate use can harm the environment, degrade soil health, and lead to pesticide resistance among pests. To address early pest identification, several studies have explored the use of image recognition

technologies for classifying pest species. However, these systems face notable challenges. The morphological similarities among different pests, variations in lighting, occlusion, and inconsistent image quality often contribute to reduced recognition accuracy. Moreover, most existing approaches stop at classification—they identify the pest type but fail to go beyond that. These systems typically do not offer contextual information such as appropriate pesticide recommendations or preventive techniques, which are crucial for informed pest management. Additionally, many models overlook the condition of the affected crop (e.g., leaf damage), which can provide meaningful clues for more accurate diagnosis.

Traditional pest inspection methods are also time-consuming, labor-intensive and require a certain specialized knowledge such that they are not suitable to be used at scale by farmers. Thus, there is a growing need for an intelligent, accessible, and comprehensive solution that not only identifies pests accurately but also equips users with actionable guidance. A system that integrates pest classification with contextual insights like treatment suggestions and prevention strategies could significantly transform how pest control is approached in agriculture—promoting sustainable practices, reducing dependency on chemicals, and supporting timely interventions in the field.

### A. Motivation

The reason for starting this work is because farmers lose a lot when they can't identify pests on time or get the wrong information, which causes their crops to produce less and makes them use too much pesticide. In many rural areas, the lack of access to timely agricultural expertise makes this issue even more critical.

With advancements in computer vision, we saw an opportunity to build a practical, deep learning based solution that not only identifies pests from crop images but also guides farmers with prevention tips and recommended pesticides. Our goal is to create a lightweight, scalable system that can be deployed through mobile or web platforms, making intelligent pest

diagnosis and guidance accessible even in resource-limited farming environments.

## II. RELATED WORK

Mehta et al. [3] proposed the FasterPest model, which enhances the FasterViT backbone by adding a feature fusion module and a pest probability adjustment module. The adjustment leverages a relationship matrix to refine predictions based on known links between pests and leaf conditions, while the fusion module combines insights from both classifiers to improve accuracy.

Tan et al. [4] demonstrated that FasterPest outperformed models such as ResNet-50 and EfficientNetV2 when evaluated on the large-scale IP102 dataset, highlighting the effectiveness of advanced backbone architectures.

Wang et al. [5] developed MobileNet-CA-YOLO by integrating YOLOv7 with MobileNetV3 and channel attention, achieving high accuracy with a compact model for rice pest detection.

Zhang et al. [6] introduced RGC-YOLO, a YOLOv8-based model using RepGhostConv and CBAM, which achieved a notable mean Average Precision (mAP) of 93.2%.

Hussain et al. [7] combined MobileNetV2 with YOLOv4 and applied CLAHE, creating a lightweight yet accurate pest detection framework suitable for low-resource environments.

Sharma et al. [8] developed an end-to-end CNN model with strong augmentation strategies to tackle inter-class similarity and noise, achieving 94.6% accuracy across more than 30 pest classes.

Li and Tan [9] proposed RiceShield, a YOLO-based framework with channel-spatial attention to detect fine pest features under low visibility conditions, attaining 92.3% accuracy.

Kumar et al. [10] introduced MultiPestNet, a multi-task hybrid CNN-Transformer architecture that assesses both pest type and leaf health, with accuracies of 91.8% and 89.5% respectively.

Mehta and Rao [11] developed TransferPest, which fine-tunes EfficientNet-B3 with limited data to adapt to resource-constrained environments, achieving 93.1% accuracy.

Patel et al. [12] created RicePestAI, a lightweight CNN optimized for mobile devices, offering real-time pest detection with 90.4% accuracy.

Lakshminadh et al. [13] applied a VGG-based deep network for rice pest classification, achieving strong results on structured pest datasets but without exploring Transformer-based architectures.

Sireesha et al. [14] worked on tomato leaf disease detection and emphasized the importance of carefully splitting datasets into training and testing sets for balanced learning.

Reddy et al. [15] studied the role of structured data splits and data augmentation in achieving high accuracy in real-time applications, influencing our adoption of a similar strategy in this work. Wang et al. [16] introduced Swin-AARNet, an enhancement of the Swin Transformer architecture tailored for fine-grained pest image recognition.

## III. METHODOLOGY

The pest classification system follows a well-defined workflow made up of three main stages: pre-processing, data loading, and model training. It begins with cleaning the raw images to remove noise, followed by labeling and applying augmentation techniques to improve the quality and variety of the data. Once prepared, the images are structured using a custom PestDataset(). In the final stage, a Vision Transformer (ViT) model is trained using the timm library, with both training and evaluation loops helping to fine-tune its performance. Figure 1 captures this complete process, showing how each step works together to create a reliable pest recognition system.

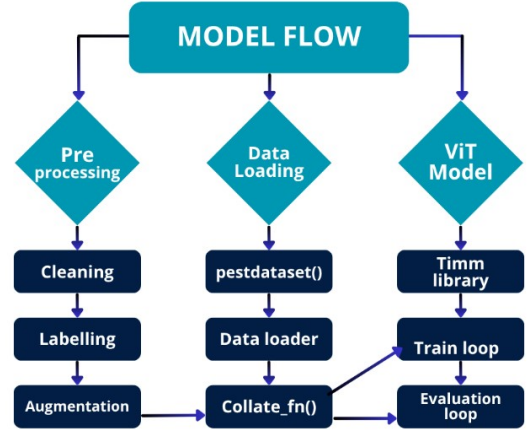


Fig. 1: Proposed model flow showing preprocessing, data loading, and Vision Transformer (ViT) training pipeline.

### A. Dataset Description

For this study, we used a curated subset of the IP102 pest dataset, which contains real-world images of rice pests across 40 different classes. Each class represents a specific type of pest affecting crops, and the images were collected in diverse environmental conditions, showcasing various pest appearances, leaf types, and backgrounds.

The raw data set started out with several problems — such as class imbalance, corrupted files and inconsistently named images. To ensure smooth processing and fair training, we performed a thorough cleanup. It included discard of unusable files, image renaming from files with special characters or spaces and, at the end, each pest class had a more fair quantity of samples.

We divided the dataset into three partitions for training, validation, and testing in [0.7, 0.1, 0.2] ratio, respectively, by a .txt file which contains the ground truth of all the images and the mappings from image path to corresponding pest label (between 0 and 39).

To give a better sense of the dataset's diversity, a few sample pest images from different classes are shown in Figure 2. These examples demonstrate the diversity and difficulty in identification of the pests, along with over-lapping characters present and natural variations of leaf shape.



Fig. 2: Sample Images In The IP102 Dataset

### B. Data Preprocessing and Integration

1) *Image Preprocessing*: For CNN-based pre-processing, all input images were resized to  $224 \times 224$  pixels. For the ViT model, images were resized to  $224 \times 224$  resolution with aspect ratio preserved. Pixel values were normalized between 0 and 1 to standardize the input and achieve stable convergence. Label encoding was handled automatically by ImageFolder, which mapped folder names to numerical labels in a class-wise directory structure.

2) *Data Augmentation*: To make the model more robust and generalizable to real-world variations (e.g., lighting, blur, or angle changes), we applied the following augmentations using the Albumentations library:

- Random Horizontal Flips
- Random Rotations (up to  $90^\circ$ )
- Brightness and Contrast Adjustments
- Gaussian Blur for simulating blur or fog
- Resizing and Cropping to mimic different camera positions

These augmentations help the model better recognize pests under diverse field conditions, improving its practical usability. Figure 3 shows how pest images are transformed during preprocessing. On the left, the raw input images vary in size, lighting, and clarity. After preprocessing (right), the images are standardized—resized, normalized, and augmented. This step enables the Vision Transformer to focus on meaningful features and boosts overall prediction accuracy.



Fig. 3: Illustration of preprocessing: Before (left) and after (right) applying resizing, normalization, and augmentations.

### C. Code Implementation and Tools

Table I outlines the various tasks involved in the code implementation phase of the pest classification work and the corresponding libraries or tools employed. Each tool was carefully selected to suit the specific needs of the work—from data preprocessing and augmentation to model training and result visualization. The integration of these tools played a crucial role in building an efficient pipeline for multi-task classification using Vision Transformers (ViT) and for organizing pest-specific knowledge mappings.

TABLE I: Core Implementation Tasks and Tools Used in Pest Classification work

Task	Library / Tool Used
Image Processing	OpenCV, PIL
Data Manipulation	NumPy, Pandas
Data Augmentation	Albumentations
Model Training (ViT)	PyTorch, timm
Visualization	Matplotlib, Seaborn
Annotation / Label Formatting	Custom Scripts, LabelImg
JSON-based Knowledge Mapping	Python json module

### D. Model Architecture

Our proposed uniform model follows a one-stage architecture built upon the Vision Transformer (ViT), as shown in Fig 4, for end-to-end pest classification. Unlike conventional multistage approaches, which decouple detection and classification, we directly handle the input images by splitting fixed-size patches, forming feature embeddings, and utilizing transformer encoders with multihead self-attention to model local and global dependencies. A so-called token ([CLS]) is utilized here to pool the context information of all patches, making it possible for the model of accurate and efficient pest detection. In addition, we also introduce a domain-aware recommendation layer after classification to recommend treatment and prevention.

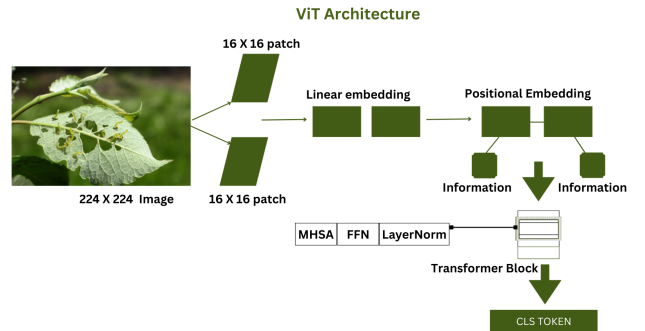


Fig. 4: Model Architectre.

Formally, each input image  $x \in \mathbb{R}^{H \times W \times C}$  is divided into  $N$  non-overlapping patches  $x_p^i \in \mathbb{R}^{P^2 \cdot C}$ , where  $H, W$ , and  $C$  denote the image height, width, and channels, and  $P$  represents

the patch size. Each patch is linearly transformed into an embedding vector as follows:

$$z_0 = [x_{\text{CLS}}; x_p^1 E; x_p^2 E; \dots; x_p^N E] + E_{\text{pos}}, \quad (1)$$

where  $E$  is the learnable projection matrix,  $E_{\text{pos}}$  represents positional embeddings, and  $x_{\text{CLS}}$  is the classification token.

The embeddings are then passed through  $L$  stacked Transformer encoder blocks. Each block applies multi-head self-attention (MHSA), which is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \quad (2)$$

where  $Q$ ,  $K$ , and  $V$  are query, key, and value matrices, and  $d_k$  is the dimension of the key vectors.

Finally, the hidden representation of the [CLS] token is fed into a classification head for pest identification:

$$\hat{y} = \text{softmax}(Wh_{\text{CLS}} + b), \quad (3)$$

where  $h_{\text{CLS}}$  is the output embedding of the classification token, and  $W, b$  are trainable parameters of the fully connected layer.

#### E. ViT-Based Classification Model

The pest classification model is based on a Vision Transformer (ViT) implemented using the `timm` library, enabling easy incorporation and access to pretrained weights. The ViT model splits each input image (resized to  $224 \times 224 \times 3$ ) into a set of small patches and processes them through self-attention layers. This design allows the model to capture both local pest features (such as wing texture and color) and the global leaf context.

##### 1) Characteristics of the Classification Model:

- Base model: `vit_base_patch16_224` pretrained ViT base model.
- Fine-tuned on 40 pest classes.
- Optimizer: AdamW with a learning rate of  $2 \times 10^{-5}$ .
- Loss Function: CrossEntropyLoss for multi-class prediction.
- Training: 50 epochs without overfitting, achieving good generalization.

These characteristics provide strong performance and high explainability for the classification task, particularly in differentiating visually similar pest classes.

#### F. Domain-Aware Recommendation Layer

Identifying the pest is only the first step in solving the problem. Farmers need to know what the pest does, how serious the situation is, and what they can do about it. To bridge this gap, our system adds a domain-based recommendation layer that turns the prediction of the model into practical guidance, as shown in Figure 5. When the Vision Transformer recognizes a pest, the system immediately looks it up in a structured knowledge base (JSON file). Instead of giving just a name, this knowledge base provides a well-rounded profile of each pest, including:

- A short description of the pest and how it affects rice crops.

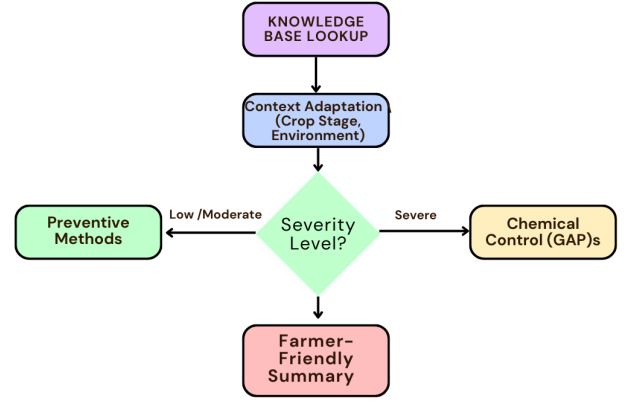


Fig. 5: Compact flow of the Domain-Aware Recommendation Layer

- Simple preventive practices such as using resistant crop varieties or maintaining proper spacing.
- Recommended chemical treatments aligned with Good Agricultural Practices (GAP).
- Eco-friendly biological control options, such as natural predators or bio-pesticides.
- Environmental conditions that make the pest more likely to spread (e.g., high humidity, warm temperatures).
- Clear severity levels (low, moderate, severe) with everyday symptom descriptions farmers can easily relate to.
- Practical advice tailored to the growth stage of the crop, since a seedling and a maturing plant may need different care.

By combining these elements, the system does more than just recognize pests; it acts as a digital farming assistant. The recommendations are not fixed but are created dynamically, taking into account the pest type, its severity, and even surrounding environmental conditions. In this way, farmers receive guidance that is clear, relevant, and ready to apply in the field. What begins as a simple image classification thus evolves into an intelligent support system that helps farmers protect their crops with timely advice and sustainable practices.

## IV. EXPERIMENTAL SETUP

The experiments were carried out on a standard personal computer running Windows 10, equipped with an Intel® Core™ i5 processor and 8 GB of RAM. Development and model training were performed using Python and PyTorch. In the absence of a dedicated GPU, all computations were performed on the CPU, with careful tuning of batch sizes and data loading to ensure efficient processing.

The IP102 data set served as the primary source for training, accompanied by pre-processing techniques such as image normalization and data augmentation. Despite the modest hardware setup, the environment proved well suited for evaluating



lightweight deep learning models, highlighting their potential for use in low-resource scenarios.

## V. RESULTS AND DISCUSSIONS

Table II provides a comparative overview of popular deep learning models used in pest detection and classification. Each method demonstrates promising accuracy and performance using combinations of CNNs and YOLO architectures. However, many focus solely on pest identification for a limited crop type, primarily rice, and often lack post-classification assistance. Our proposed model, built on Vision Transformers (ViT), outperforms existing approaches in both accuracy and F1-score. Moreover, it introduces a novel layer of functionality by providing pesticide recommendations and preventive tips post-prediction, making it more comprehensive and farmer-friendly for field-level deployment across diverse crop types.

TABLE II: Performance Comparison of Pest Detection Models Using Evaluation Metrics

S. No	Title	Metrics			
		Accuracy	Precision	Recall	F1 Score
1	MobileNetCA-YOLO for Rice Pest and Disease Detection [5]	92.3	95.2	92	91
2	Lightweight YOLOv8-based model [6]	94	94.3	93.2	92
3	MobileNetV2-YOLOv4 for Pest Detection [7]	93.2	86.2	90.8	91
4	Swin Transformer (Swin-T) (AARNet, 2024) [16]	78.8	78.4	78.8	78.6
5	ResNet-50 (IP102 baseline, Wu et al. 2019) [17]	49.4	49.0	49.4	49.2
6	DeiT-S Data-efficient Transformer, IP102 dataset	74.2	73.8	74.2	74.0
7	<b>Proposed ViT-Based Pest Classifier with Remedy Support</b>	<b>96</b>	<b>95.4</b>	<b>94.7</b>	<b>95.0</b>

### A. Training and Validation Accuracy

The Vision Transformer (ViT) model was trained for 50 epochs using the AdamW optimizer. Training proceeded steadily, with the model effectively adapting to the challenging pest image dataset as shown in Figure 7.

- **Final Validation Accuracy:** Approximately 96%.
- **Loss Trend:** Training and validation losses exhibited a smooth, continuous decline until stabilizing, indicating proper learning and generalization.

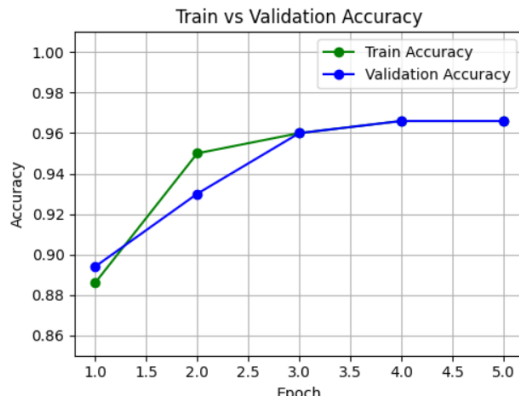


Fig. 6: Training and validation accuracy trends over epochs.

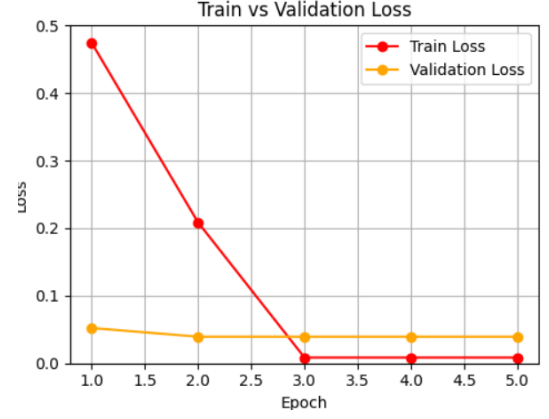


Fig. 7: Training and validation loss trends over epochs.

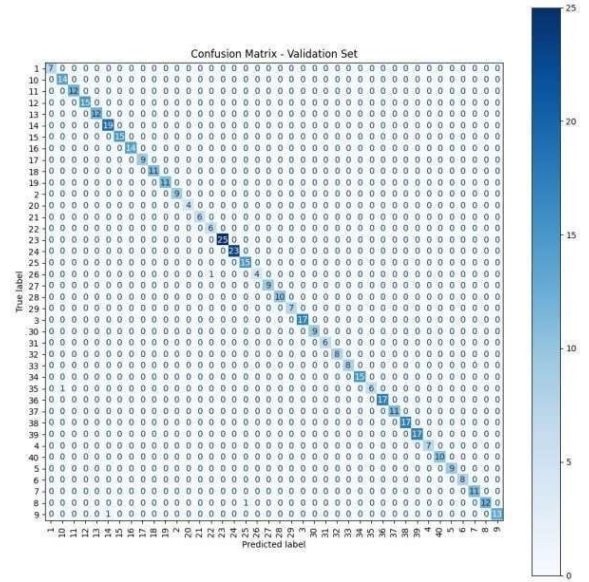
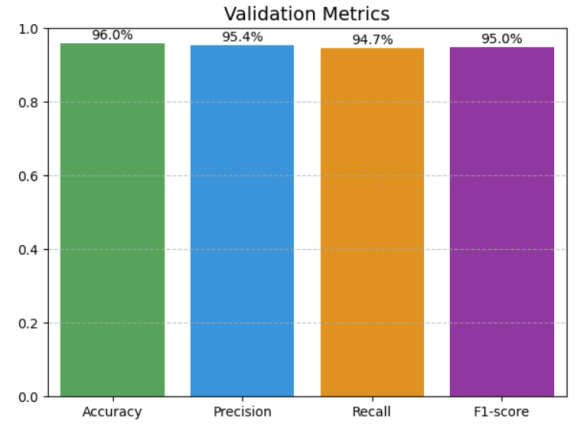


Fig. 8: Comparison of metrics and confusion matrix of the proposed model.

### B. ViT Classification Evaluation

The Vision Transformer (ViT) model was fine-tuned on a custom pest dataset consisting of 40 distinct pest classes and

tested across various metrics. Figure 8 represents the evaluation results of the proposed model, including validation metrics (accuracy, precision, recall, F1-score) and the confusion matrix showing the classification performance across all classes.

### C. Model Output Summary

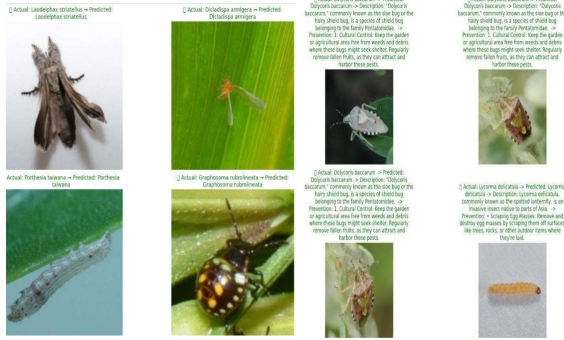


Fig. 9: Model inference on input images with labeled class names.

Figure 9 represents the pest images with the labelled class names and actual pest names along with predicted pest labels, added description, prevention methods, and pesticide recommendation.

## VI. CONCLUSION AND FUTURE SCOPE

### A. Conclusion

The paper presents **PestVision+**, a smart pest classification and treatment recommendation system designed to assist farmers in real-world agricultural environments. The system employs Vision Transformers (ViT) for high-accuracy classification over 40 pest species, trained on a well-structured and augmented image dataset.

PestVision+ extends beyond pest classification by linking predictions to a domain-aware knowledge base that provides pesticide suggestions and preventive measures, acting as an intelligent agricultural assistant. The model achieves a validation accuracy of **96%** with high generalization capability and precision, particularly in distinguishing visually similar pest classes—thanks to robust data augmentation, safe loading pipelines, and an efficient hierarchical structure.

### B. Future Scope and Research Opportunities

While PestVision+ establishes a strong foundation, several improvements can enhance its usability and scope:

- **Support for Diverse Crop Types:** Future iterations can incorporate multi-crop pest datasets to expand applicability beyond rice.
- **Pest Impact-Based Leaf Condition Analysis:** Adding severity scoring for leaf damage can help farmers assess infestation levels and treatment urgency.
- **Farmer-Centric Interfaces:** Regional language support and voice-based interaction can improve accessibility for rural and non-technical users

- **Drone-Based Pest Monitoring:** Using drones to capture crop images can help farmers spot pests early, target affected areas efficiently, and eventually enable automatic pest identification from aerial views.

## REFERENCES

- [1] K. Thenmozhi and M. Reddy, "Crop pest classification based on deep convolutional neural network and transfer learning," *Computers and Electronics in Agriculture*, vol. 164, p. 104906, 2019.
- [2] Y. Zhang, "Knowledge-driven pest management recommendation system," *Computers and Electronics in Agriculture*, vol. 180, p. 105864, 2021.
- [3] R. Mehta and D. Rao, "Fasterpest: An efficient framework for pest detection and leaf condition assessment," *Computers and Electronics in Agriculture*, vol. 208, p. 107785, 2024.
- [4] M. Tan and Q. Le, "Efficientnetv2: Smaller models and faster training," in *Proceedings of the International Conference on Machine Learning (ICML)*. PMLR, 2021.
- [5] C. Wang, X. Li, and Y. Huang, "Rice pest and disease detection using mobilenet-ca-yolo," *Applied Sciences*, vol. 13, no. 9, p. 4881, 2023.
- [6] H. Zhang, L. Chen, and F. Sun, "Rgc-yolo: Lightweight model for multi-scale rice disease detection," *Sensors*, vol. 25, no. 1, p. 135, 2025.
- [7] A. Hussain, M. A. Raza, and A. R. Javed, "Deep mobilenetv2-yolov4 for intelligent pest detection in low-resource settings," *IEEE Access*, vol. 11, pp. 12 345–12 356, 2023.
- [8] P. Sharma and R. Kumar, "Pestnet: An end-to-end deep learning approach for large-scale multi-class pest detection and classification," *Computers and Electronics in Agriculture*, vol. 212, p. 108195, 2024.
- [9] S. Li and Y. Tan, "Riceshield: Robust detection of subtle rice pest features using channel-spatial attention," *Computers and Electronics in Agriculture*, vol. 209, p. 107893, 2024.
- [10] V. Kumar, R. Singh, and M. Patel, "Multipestnet: Multi-task cnn-transformer model for pest and leaf health assessment," *Computers and Electronics in Agriculture*, vol. 210, p. 107900, 2024.
- [11] R. Mehta and D. Rao, "Transferpest: Leveraging transfer learning for accurate rice pest detection in low-resource settings," *Neural Computing and Applications*, vol. 36, pp. 4825–4840, 2024.
- [12] A. Patel and M. Suresh, "Ricepestai: Automated detection of rice plant insect damage using lightweight cnns," *IEEE Transactions on Instrumentation and Measurement*, vol. 73, pp. 1–9, 2024.
- [13] K. Lakshminadh, "Advanced pest identification: An efficient deep learning approach using vgg networks," in *IEEE IATMSI*, 2025.
- [14] M. Sireesha and D. V. Reddy, "Deep learning-based tomato leaf disease identification," in *IEEE IATMSI*, 2025.
- [15] K. V. N. Reddy, "Automated traffic sign recognition via cnn deep learning," in *IEEE IATMSI*, 2025.
- [16] F. Wang *et al.*, "Swin-aarnet: Enhanced swin transformer for fine-grained pest image recognition," in *Proceedings of the International Conference on Agricultural AI*. Springer, 2024, pp. 123–130, presented at IC-AI 2024.
- [17] F. Wu *et al.*, "Resnet-50 baseline for pest detection using the ip102 dataset," *Journal of Agricultural Informatics*, vol. 15, no. 4, pp. 45–54, 2019, evaluated on over 75,000 images across 102 pest categories.