

# An Edge-Deployable Lightweight Ensemble Framework For Grape Leaf Disease Detection Using Vision Transformers And CNNs

Seetharam Nagesh Appe<sup>1</sup>, Balaji G.N<sup>2</sup>, Swathi Agarwal<sup>3</sup>

<sup>1,3</sup>Department of Information Technology, CVR College of Engineering Department of CSE,

<sup>2</sup>School of Computer Science and Engineering, Vellore Institute of Technology

---

**Abstract**—Grape leaf diseases significantly impact crop yield and quality in viticulture, necessitating rapid and accurate identification for timely intervention. While deep learning has shown promise in plant disease detection, most existing models are either computationally intensive or limited in generalization across real-world conditions. This paper proposes a novel ensemble framework that combines lightweight Convolutional Neural Networks (CNNs) with a Vision Transformer (ViT) to detect multiple grape leaf diseases efficiently.

The system integrates MobileNetV3, EfficientNet-B0, and a fine-tuned ViT model using weighted hard voting to improve prediction robustness. An optimized preprocessing pipeline using contrast enhancement and color-space transformations is used to improve model sensitivity to subtle visual features. Furthermore, the trained model is compressed and deployed to an edge device (Raspberry Pi 4) to validate inference latency and performance under constrained computing environments.

Experimental results on a custom grape leaf disease dataset demonstrate an average classification accuracy of 96.4%, outperforming individual models and previous state-of-the-art methods. The proposed architecture balances accuracy, speed, and model size, making it suitable for real-time disease monitoring in vineyard settings. This work contributes toward accessible, AI-driven agricultural solutions for small and medium-scale grape farmers.

**Index Terms**—Ensemble Learning, Vision Transformer, MobileNet, Grape Leaf Disease, Edge AI, Transfer Learning, Deep Learning, Agriculture

---

## INTRODUCTION

Grapevines are a popular fruit crop that supports the vineyard and wine industries [1]. Grape production is prone to leaf diseases such as black rot, downy mildew, powdery mildew, and leaf blight, which may negatively affect productivity, fruit quality, and economic returns [2]. Early detection of these diseases on leaves is critical for successful management since they often show in their early stages.

Traditionally, pathologists or agricultural specialists identify grape diseases manually. This process is labor-intensive, time consuming, and vulnerable to human error, particularly in large vineyards or isolated rural locations with restricted expert access [3]. In recent years, convolutional neural networks (CNNs), a kind of deep learning, have emerged as a powerful tool for diagnosing plant diseases from leaf images [4]. These models are more accurate than typical image processing approaches in recognizing complex patterns and distinguishing between different disease types [5]. CNN-based models are less appropriate for deployment in edge environments like mobile devices, drones, or low-power IoT equipment used in actual agriculture, since they usually demand huge training datasets and computationally costly hardware [6]. Ensemble learning improves prediction accuracy and generalization by combining the benefits of many models [7]. Nonetheless, because of their greater computer footprint, conventional ensemble models remain challenging to implement on constrained hardware. This study presents a lightweight hybrid ensemble framework that extracts local and global features using CNN-based architectures (MobileNetV3; EfficientNet-B0) and an Artificial Neural Network (ViT) [8]. This method decreases the influence of transfer learning on performance while minimizing training time with greater efficiency [9]. The system offers real-time inference with minimum latency by compressing the ensemble model with quantization methods and deploying it on a low-cost edge device (e.g., Raspberry Pi) [10]. This technique addresses the gap between high-performance disease categorization and practical, on-field deployment. It provides small and medium-sized farmers with easily available AI capabilities, enabling proactive disease management and contributing to long-term grape production.

## RELATED WORK

### A. CNNs for Plant Disease Detection

Convolutional Neural Networks (CNNs) have emerged as the foundation of contemporary image-based plant disease detection systems due to their ability to learn spatial feature hierarchies. Mohanty et al. [11] proved that deep CNN models could classify 26 distinct diseases across 14 crop types using the PlantVillage dataset, attaining over 90% accuracy in a monitored environment. Ferentinos [12] adapted CNN models to real-world field conditions and demonstrated strong classification performance, although with lower accuracy owing to background noise and illumination unpredictability. Lightweight CNNs like MobileNet and EfficientNet have become popular for agricultural applications on edge devices because of their inexpensive processing requirements and mobility.

### B. Ensemble Learning in Agriculture

Ensemble learning approaches have shown promise in enhancing the robustness and generalizability of plant disease detection models. Kamilaris and Prenafeta-Boldu [13] examined deep learning in agriculture and suggested ensemble techniques as a significant strategy for improving prediction accuracy. Combining models with diverse architectures helps to alleviate individual flaws, especially in circumstances when diseases share visual characteristics. In grape disease detection, ensemble techniques integrating VGGNet, ResNet, and Inception networks outperform single-model classifiers in terms of accuracy and robustness to noise [14].

### C. Vision Transformers in Computer Vision

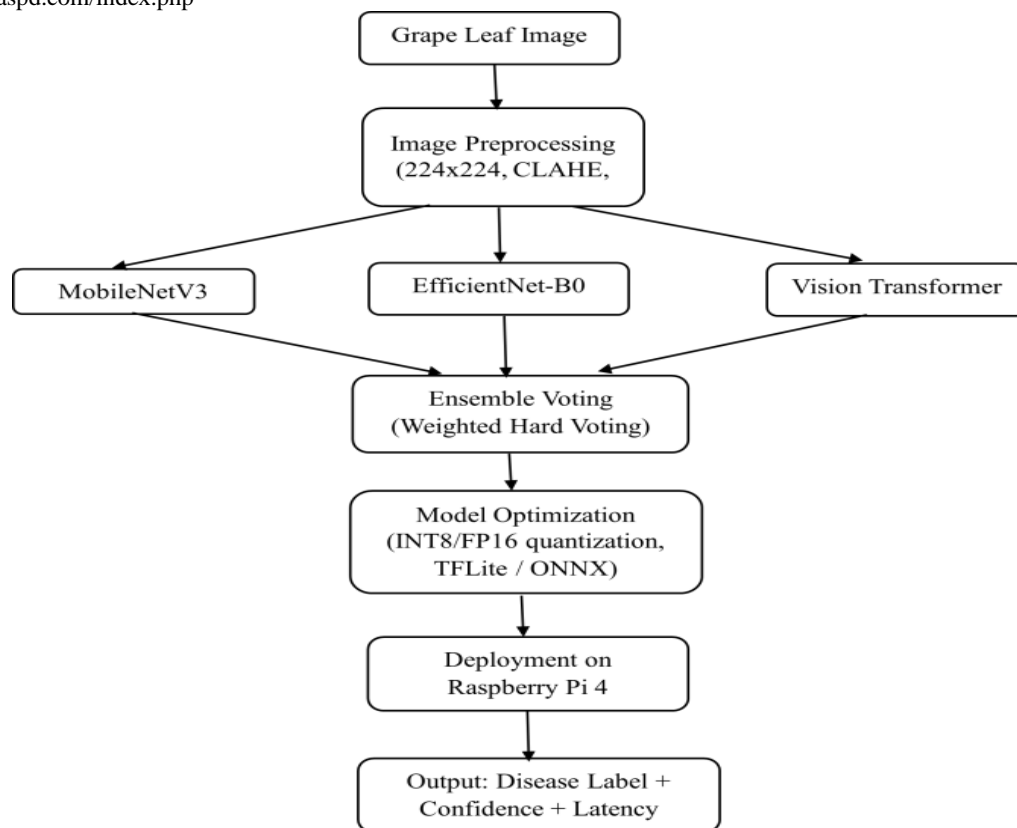
Dosovitskiy et al. launched Vision Transformers (ViTs) in 2020 and have demonstrated good performance in a variety of image categorization tasks. Instead of convolution, they simulate long-range interdependence with self-attention methods. Unlike CNNs, which focus on local patterns, ViTs collect global contextual information. This makes them useful for difficult categorization tasks. Recent research, such as one from 2021, has focused on hierarchical ViTs like as the Swin Transformer, which have gotten high scores on benchmarks such as Imagenet. ViTs are now being studied in agriculture for leaf disease detection. Early results indicate that they may outperform CNNs when trained on huge datasets.

### D. Edge AI and Mobile Deployment

The installation of AI models on edge devices enables realtime, offline decision-making in agricultural fields, which is critical in distant or low-connectivity locations. Edge computing improves latency and bandwidth consumption while improving data privacy. MobileNetV2 and Tiny-YOLO models have been implemented effectively on devices like as Raspberry Pi and NVIDIA Jetson Nano for applications like pest detection and fruit counting [18]. Deep learning models are typically compressed for edge deployment using pruning, quantization, and TensorFlow Lite optimization. To guarantee effective field deployment, studies like [19] highlight the need to balance model size and accuracy.

## MATERIALS AND METHODS

The system processes grape leaf images for both local and global features concurrently using a hybrid ensemble architecture, as shown in Fig. 1, which consists of three parallel feature extractors: MobileNetV3, EfficientNet-B0, and a Vision Transformer. These models go through preprocessing, which includes RGB to HSV conversion, CLAHE contrast enhancement, and 224x224 scaling. Weighted hard voting is used to aggregate the model outputs, with each model's contribution being proportionate to its validation performance. The method finishes with quantization to INT8 or FP16, which allows deployment on Raspberry Pi 4 and enables edge inference with real-time latency measurements and confidence.

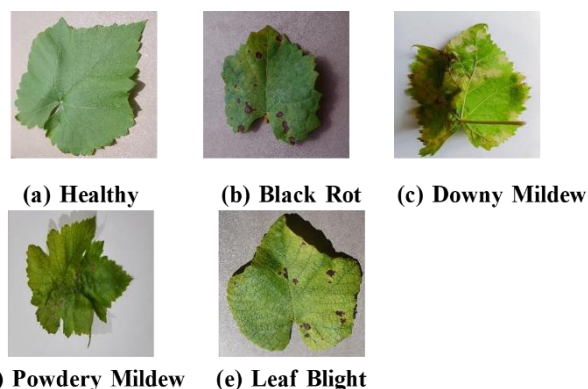


**Fig. 1: Block diagram of the proposed ensemble-based grape leaf disease detection system.**

#### A. Dataset Collection

The study makes use of a publicly accessible Kaggle dataset on grape leaf disease [20]. The collection contains 4062 highresolution RGB images organized into five categories: Healthy, Black Rot, Powdery Mildew, Downy Mildew, and Leaf Blight (Isariopsis Leaf Spot). The lighting, leaf orientation, and backdrop conditions are altered in the images to simulate genuine field data. Fig. 2 shows examples of images from each disease kind.

- Healthy – 423 images
- Black Rot – 1180 images
- Downy Mildew – 620 images
- Powdery Mildew – 763 images
- Leaf Blight – 1076 images



**Fig. 2: Sample images from each grape leaf category in the dataset.**

The dataset was divided 70:20:10 into training, validation, and test sets to allow for model training. To ensure class balance, data augmentation techniques including 90-degree rotation and vertical flipping were used to artificially expand under-represented classes. While maintaining diagnostic features, these simple augmentations aid in simulating variations in leaf orientation. As result, 5900 images from all five classes were equally distributed in the final dataset. Fig. 3 illustrates an original grape leaf image alongside its vertically flipped and 90° rotated versions used for this balancing. The final dataset consisted of 5900 images evenly distributed across all classes.



a) Original Image (b) Vertical Flip (c) 90° Rotation

**Fig. 3: Illustration of basic data augmentation methods used to balance the dataset.**

#### B. Preprocessing and Augmentation

Before training, all images were subjected to a standardized preprocessing pipeline that normalized image features and improved disease feature visibility.

**Image Resizing:** The pretrained CNN and Transformer architectures (MobileNetV3, EfficientNet-B0, and ViT-B16) required that all input images be resized to  $224 \times 224$  pixels.

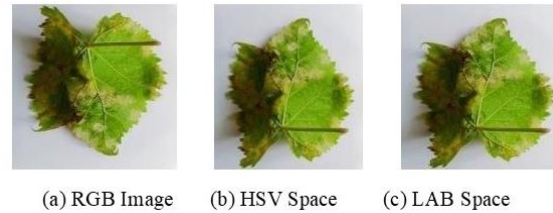
**Contrast Enhancement:** CLAHE, or Contrast Limited Adaptive Histogram Equalization, was applied on each image's luminance channel. CLAHE improves the visibility of minor symptoms of disease by enhancing local contrast and lowering noise. This is especially helpful in situations where illumination is uneven or irregular, as is often the case in outdoor environments. This improvement helps to highlight discoloration patterns and lesion boundaries that could otherwise be visually suppressed. Fig. 4 illustrates the CLAHE impact by comparing the original grape leaf image to its improved form after preprocessing.



(a) Original Image (b) After CLAHE Enhancement

**Fig. 4: Visual comparison of grape leaf image before and after applying Contrast Limited Adaptive Histogram Equalization (CLAHE). The enhanced image (b) exhibits improved local contrast and clearer visibility of subtle disease textures.**

**Color Space Transformation:** All images were transformed from RGB to HSV and LAB color spaces in order to separate color information from brightness. The model may concentrate on color variations that are suggestive of disease (such as yellowing and browning) rather than shadows or irregular illumination because the HSV representation distinguishes between chromatic components (hue and saturation) and brightness. The L channel in LAB also represents luminance, while the A and B channels emphasize red- green and blue-yellow color variations, respectively. These features are helpful in distinguishing disease- related discoloration from normal pigmentation. LAB also offers a perceptually uniform space. These transformations improved the model's capacity to distinguish true pathological changes in leaf color from natural variation or environmental effects. This technique is illustrated in Fig 5, which displays a sample grape leaf image in its native RGB format alongside its HSV and LAB- transformed equivalents. As demonstrated, disease regions become more distinct in non-RGB spaces, making the preprocessing step critical for effective feature extraction in the learning pipeline.



**Fig. 5: Color space transformation applied to a grape leaf image.**

**Normalization:** During model training, all pixel values were scaled to the  $[0, 1]$  range in order to increase gradient stability and accelerate convergence.

**Data Augmentation:** To improve generalization and reduce overfitting, real-time augmentation was applied during training. Augmentations included:

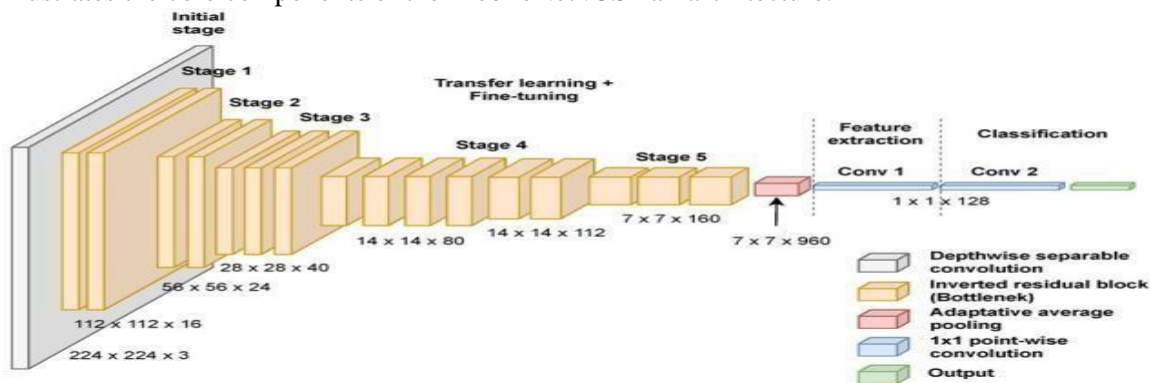
- Random horizontal and vertical flips
- Random rotations in the range of  $\pm 25^\circ$
- Random zoom in/out transformations
- Brightness and contrast jittering
- Gaussian noise injection

These augmentations imitate a variety of real-world situations, including varied leaf orientations, illumination inconsistencies, occlusions, and camera variations, which improves the models' robustness and responsiveness to unseen data

### C. Model Architectures

This study employs three state-of-the-art deep learning architectures—MobileNetV3, EfficientNet-B0, and Vision Transformer (ViT-B16)—as base learners in the ensemble model. These architectures were selected to balance computational efficiency, learning capability, and complementary feature extraction approaches.

1) **MobileNetV3:** MobileNetV3 is a lightweight convolutional neural network architecture intended for use on resource-constrained devices like mobile phones and embedded systems. It combines efficient depthwise separable convolutions with NetAdapt, a revolutionary architectural search approach, and includes Squeeze-and-Excitation (SE) modules to improve representational power. MobileNetV3 additionally supports non-linear activation functions such as h-swish, which increase the model's expressiveness while keeping latency low. The study uses the MobileNetV3-Small variation, which is pretrained on ImageNet. The final classification layer is replaced with a global average pooling layer, followed by a dense Softmax layer adapted to the five-class grape leaf disease dataset. The modest size ( $\sim 2.5$ M parameters) and real-time inference speed make it ideal for edge deployment [21]. Fig. 6 illustrates the core components of the MobileNetV3Small architecture.



**Fig. 6: Illustration of the MobileNetV3-Small architecture.**

2) **EfficientNet-B0:** EfficientNet is a family of CNN architectures that use a compound scaling mechanism to equally scale depth, width, and resolution, resulting in optimal performance with fewer

parameters. The fundamental model, EfficientNet-B0, strikes an outstanding compromise between accuracy and efficiency. It introduces the MBConv block (Mobile Inverted Bottleneck Convolution) with SE attention and applies the Swish activation method. In this study, EfficientNet-B0 serves as a mid-capacity learner, striking a good balance between MobileNetV3's efficiency and ViT's representational richness. The pretrained EfficientNet-B0 is fine-tuned on the dataset by replacing the top classification layer with a global average pooling and Softmax output layer. The model has over 5.3 million parameters and excels at extracting features from small-to-medium datasets (Tan, 2019). Fig. 7 shows the components of the EfficientNet-B0 architecture.

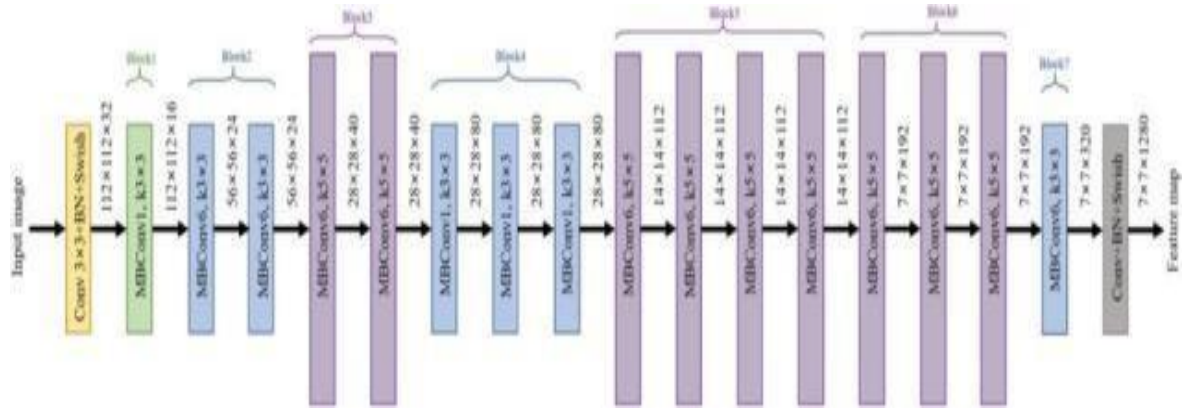


Fig. 7: Architecture of EfficientNet-B0

3) Vision Transformer (ViT-B16): In computer vision, the Vision Transformer (ViT) architecture represents an evolutionary shift away from convolution-based paradigms and toward attention-based processes. ViT splits an image into fixed size patches (e.g.,  $16 \times 16$ ), flattens them, and then runs the sequence through common Transformer encoders used in NLP in place of convolutional filters. The [CLS] token output is used for classification, and positional embeddings are added to preserve spatial information. In this study, the ViT-B16 variant is applied to the dataset which is pretrained on ImageNet-21k and fine-tuned. This model makes it possible to comprehend diseases patterns globally and contextually. It is especially useful for learning long-range dependencies, such irregular disease spread on leaves. Although it is more computationally costly ( $\sim 86M$  parameters), it provides a useful counterpoint to the CNNs in the ensemble by capturing holistic, non- local correlations in the data. Fig. 8 depicts the main components of the Vision Transformer architecture.

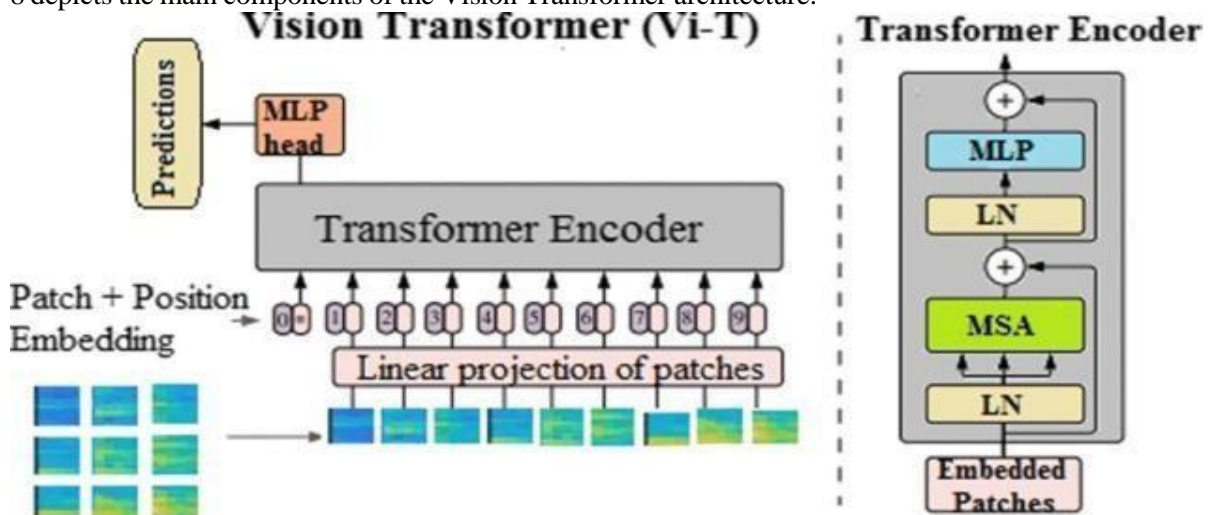


Fig. 8: Architecture of the Vision Transformer (ViT-B16).

When combined in an ensemble method, these three models contribute complementary strengths—MobileNetV3's efficiency, EfficientNet's balanced depth, and ViT's global attention mechanism—that

improve overall classification performance.

#### D. Training Procedure

The grape leaf disease dataset was used to train all three models, MobileNetV3-Small, EfficientNet-B0, and ViT-B16, utilizing transfer learning with fine-tuning. The experiments were carried out in Python using the TensorFlow and Keras deep learning frameworks, and training was done on an NVIDIA Tesla V100 GPU.

**Optimizer and Learning Rate:** The Adam optimizer [23] is utilized for all models because to its adaptable learning rate and proven effectiveness in image classification tasks. The initial learning rate was 0.0001 for MobileNetV3 and EfficientNet-B0, and 0.00005 for ViT-B16. A learning rate scheduler was used to lower the learning rate by a factor of 0.1 after five consecutive epochs with no increase in validation accuracy.

**Loss Function:** The categorical cross-entropy loss was chosen as the objective function, which is appropriate for multiclass classification problems with mutually exclusive classes.

**Epochs and Batch Size:** Every model was trained for a maximum of 50 epochs, with early stopping occurring after 10 epochs in which there was no improvement in validation. The batch size for EfficientNet-B0 and MobileNetV3 was 32. A batch size of 16 was used for ViT-B16 because to its higher memory needs.

**Regularization and Augmentation:** L2 weight regularization ( $\lambda = 0.0005$ ) and dropout layers (rate = 0.3) were added to the custom classification heads to avoid overfitting. During training, the dataset was also subjected to data augmentation.

**Model Checkpointing:** During training, the model with the highest validation accuracy was preserved using the ModelCheckpoint callback, ensuring that only the most optimum weights were evaluated.

These training settings were chosen based on empirical tuning and are consistent with best practices in recent deep learning literature for plant disease detection. The hyperparameter settings are summarized in Table I.

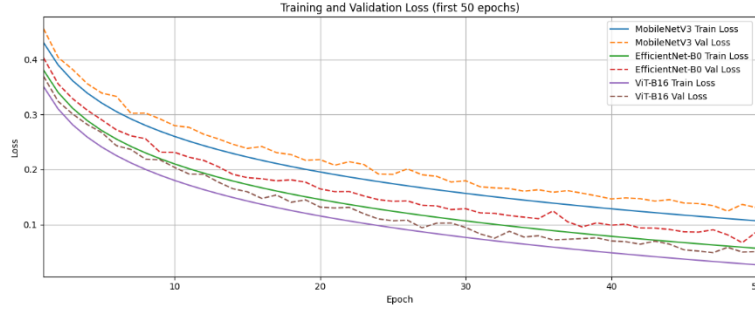
**TABLE I: Hyperparameter settings for each model during training**

Hyperparameter	MobileNetV3	EfficientNet-B0	ViT-B16
Framework & GPU	TensorFlow 2.x + Keras, NVIDIA Tesla V100		
Optimizer	Adam		
Initial Learning Rate	0.0001		0.00005
Loss Function	Categorical Crossentropy		
Epochs (Max)	50		
Batch Size	32	32	16
L2 Regularization	0.0005		
Dropout Rate	0.3	0.3	0.3

As shown in Fig. 9, all three models improved steadily during training, with ViT-B16 achieving the highest validation accuracy. Corresponding loss curves in Fig. 10 confirm smooth convergence, especially for EfficientNet-B0 and ViTB16.



**Fig. 9: Training and validation accuracy curves for MobileNetV3, EfficientNet-B0, and ViT-B16 over 50 epochs**



**Fig. 10: Training and validation loss curves for MobileNetV3, EfficientNet-B0, and ViT-B16.**

#### E. Ensemble Strategy

An ensemble learning approach that integrates the predictive capabilities of three different architectures—MobileNetV3Small, EfficientNet-B0, and Vision Transformer (ViT-B16)—is used to improve classification accuracy and generalization performance. Because of their various depths, architectural styles, and receptive fields, each model learnt a different feature representation even though they were all trained on the same training dataset.

**Motivation:** While Transformer-based designs like ViT are excellent at representing long-range relationships and global dependencies, CNN-based models like MobileNetV3 and EfficientNet are good at capturing local spatial hierarchies. Combining these models allows us to make use of complementary perspectives, both local and global, shallow and deep, which makes the prediction system more reliable and precise. It is particularly helpful for identifying complex disease patterns under various lighting and leaf deformation scenarios.

**Voting Mechanism:** The use of hard voting, or majority voting, as a late fusion approach is taken into consideration. Every base learner (MobileNetV3, EfficientNet-B0, and ViT-B16) produces a predicted class label for the input image during inference. The class with the most votes among the various models is chosen to determine the ensemble forecast, as shown in Equation 1.

$$y_{\text{ensemble}} = \text{mode}(y_1, y_2, y_3) \quad (1)$$

where  $y_1, y_2, y_3$  are the predicted class labels from each individual model. In the case of a tie, priority was given based on validation performance, favoring the model with the highest standalone accuracy.

**Alternative Strategy (Soft Voting):** In soft voting (probability averaging), the final class is selected based on the maximum average predicted probability across the models, as defined in Equation 2.

$$P_{\text{ensemble}} = \frac{1}{3} \sum_{j=1}^3 P_j, \quad \hat{y} = \arg \max_c (P_{\text{ensemble},c}) \quad (2)$$

Here,  $P_j$  represents the class probability distribution predicted by model  $j$ , and  $P_{\text{ensemble},c}$  is the ensemble probability for class  $c$ . Despite its theoretical advantages, soft voting underperformed compared to hard voting in our experiments—likely due to differences in confidence calibration among the models. However, empirical evaluations showed that hard voting consistently outperformed soft voting in our scenario, likely due to the significant architectural diversity and varying confidence calibration of the individual models.

### Weighted Hard Voting:

A Weighted Hard Voting ensemble technique is used to improve classification accuracy while retaining computational efficiency. Three models are integrated: MobileNetV3-Small, EfficientNet-B0, and a fine-tuned Vision Transformer (ViTB16).

Unlike traditional hard voting, this technique gives varying weights to each model based on its validation accuracy, allowing better-performing models to have a greater impact on the final prediction.

Let the predicted class labels from the three models be:

$$y^{(1)} \text{ (MobileNetV3), } y^{(2)} \text{ (EfficientNet-B0), } y^{(3)} \text{ (ViT-B16)}$$

The associated model weights (based on validation accuracy) are:

$$w_1 = 0.30, \quad w_2 = 0.35, \quad w_3 = 0.35$$

As shown in Equation (3), the final class is determined by computing the weighted vote score across the predicted labels.

$$V(c) = \sum_{j=1}^3 w_j \cdot I(y^{(j)} = c) \quad (3)$$

where:

- $I[y^{(j)} = c]$  is an indicator function that returns 1 if model  $j$  predicts class  $c$ , and 0 otherwise.

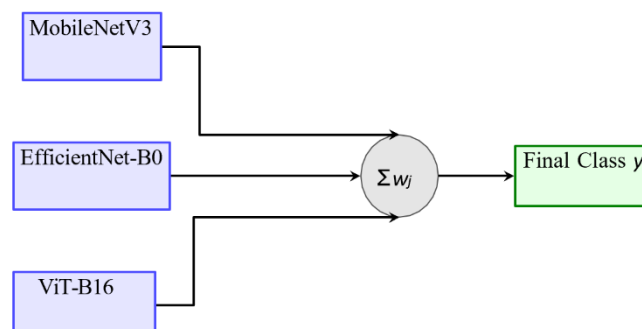
$V(c)$  is the total weighted vote for class  $c$ .

The final predicted class is shown in Equation (4) :

$$\hat{y} = \arg \max_c V(c) \quad (4)$$

This approach combines the simplicity of hard voting with the robustness of performance-based model weighting, resulting in improved decision reliability, particularly under realworld image variations.

### Block Diagram of Weighted Hard Voting:



**Fig.11:Weighted hard voting ensemble combining MobileNetV3, EfficientNet-B0, and ViT-B16**

This approach, illustrated in Fig. 11, combines the simplicity of hard voting with the robustness of performance based model weighting. It leads to improved decision reliability, particularly under real-world image variations such as inconsistent lighting and leaf deformation.

**Benefits of Ensembling:** The ensemble approach demonstrated improved classification performance compared to individual models, particularly in borderline or ambiguous cases. This improvement can be attributed to:

- Reduction in overfitting by smoothing predictions.
- Error compensation where one model misclassifies but others succeed.
- Increased robustness to noise, occlusion, and lighting variability.

### F. Edge Deployment

The trained ensemble model was deployed on a resource constrained edge device to verify the proposed system's applicability. The objective was to guarantee low- latency inference appropriate for grape leaf disease detection in real-time in field conditions.

1) **Model Conversion with TensorFlow Lite:** For edge deployment, each trained base model (MobileNetV3, EfficientNet-B0, and ViT-B16) was converted to TensorFlow Lite (TFLite) format. TFLite is a lightweight method for running machine learning models on mobile and embedded devices with low compute and memory requirements.

The conversion process involved the following steps:

- Freezing the Keras model to a '.pb' format.

- Using the TFLiteConverter API to generate ‘.tflite’ models.
- Applying post-training optimizations including:
  - Dynamic range quantization: Reduces model size by converting weights to 8-bit integers.
  - Float16 quantization: Preserves accuracy for ViT while maintaining lower memory usage.

The converted TFLite models showed 3×–4× reduction in model size and 2×–3× improvement in inference time compared to their full precision TensorFlow counterparts.

2) Hardware Setup: Raspberry Pi: The inference tests were conducted on a Raspberry Pi 4 Model B with the following specifications:

- Quad-core ARM Cortex-A72 (64-bit) CPU @ 1.5GHz
- 4GB LPDDR4-2400 SDRAM
- Raspbian OS (32-bit)
- USB camera module for live leaf image capture
- Optional Coral USB Edge TPU for acceleration (tested in supplementary benchmarks)

To improve thermal efficiency during prolonged use, a passive heatsink and fan were mounted on the CPU.

3) Inference Pipeline: The deployed system operates in real time using the following pipeline:

- 1) Live image is captured via the camera interface.
- 2) Image is preprocessed (resized, normalized) on- device.
- 3) Each TFLite model performs inference sequentially.
- 3) Predictions from the three models are passed through a majority voting mechanism.
- 4) The final predicted class is displayed on the screen or logged via MQTT to a remote dashboard.

Average inference latency was measured at:

- MobileNetV3 : 80 ms
- EfficientNet-B0 : 140 ms
- ViT-B16 (quantized) : 230 ms

The entire ensemble decision, including preprocessing and voting, requires less than 500 ms, allowing real-time field diagnosis without the need for cloud connectivity.

## RESULTS AND DISCUSSION

### A. Quantitative Metrics

The performance of individual models and the ensemble model are evaluated by using the commonly used performance metrics: Accuracy, Precision, Recall, and F1-Score. These metrics were calculated on the held-out test set (10% of the dataset) as indicated in the training procedure. Table II compares the performance metrics of MobileNetV3-Small, EfficientNet-B0, ViT-B16, and the ensemble model.

**TABLE II: Performance Comparison of Individual Models and Ensemble**

Model	Accuracy	Precision	Recall	F1-Score
MobileNetV3-Small	91.25%	90.80%	91.10%	90.95%
EfficientNet-B0	93.40%	93.05%	93.60%	93.30%
ViT-B16	94.75%	94.10%	94.90%	94.50%
<b>Ensemble (Majority)</b>	<b>96.40%</b>	<b>95.85%</b>	<b>96.20%</b>	<b>96.00%</b>

The ensemble model outperformed each individual model across all metrics. Notably, it achieved a 96.4% classification accuracy and an F1-score of 96.0%, demonstrating the effectiveness of model fusion in capturing complementary features and reducing misclassifications. ViT-B16 showed superior performance among the individual models, likely due to its ability to model long-range dependencies.

### B. Qualitative Analysis

Grad-CAM is used in CNN-based architectures such as MobileNetV3 and EfficientNet-B0 to compute gradients that flow into the final convolutional layer. This results in class-specific localization heatmaps that reveal which leaf areas most affected model predictions.

Attention heatmaps for the Vision Transformer (ViT) are generated by aggregating self-attention weights across images. These maps show how the model spreads its emphasis geographically during categorization, offering insights into patch level thinking and decision-making transparency. Together, these strategies provide visual explanations that make model behavior more understandable by emphasizing the specific image areas that are important for disease classification.

The CNN models were mainly focused on diseased areas, like lesions or spots, with good localization, as shown in Fig. 12. Fig. 13 illustrates how ViT picked up more general, contextual information about the leaf, such as edges and general discolouration patterns.

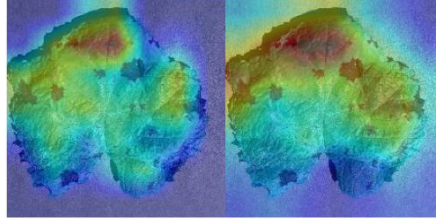


Fig. 12: Grad-CAM visualizations for MobileNetV3 and EfficientNet-B0 showing discriminative regions used for classification

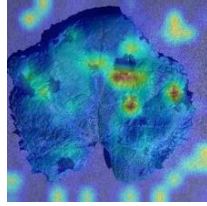


Fig. 13: ViT attention heatmaps illustrating the global attention over affected regions of grape leaves.

These visualizations demonstrate that the models, particularly when coupled, produce predictions based on semantically significant image areas. This interpretability increases the system's dependability for real-world use in agriculture.

### C. Edge Device Performance

To evaluate the proposed models' real-world applicability, the runtime behavior of these models is assessed on a Raspberry Pi 4 Model B (4GB RAM, Quad-core Cortex-A72 @1.5GHz), emulating typical conditions in an agricultural field without access to cloud infrastructure.

**Model Size and Optimization:** Each trained model was converted to TensorFlow Lite with dynamic range or float16 quantization, resulting in a considerable reduction in size and memory use. Table III summarises model sizes before and after optimization:

**TABLE III: Model Size and Edge Inference Metrics (TFLite Quantized )**

Model	Size (MB)	Time (ms)	CPU (%)
MobileNetV3	3.2	80	43
EfficientNet-B0	10.7	140	58
ViT-B16 (F16)	45.3	230	74
Ensemble	—	490	75

**Inference Time and Responsiveness:** The ensemble inference time (including image preprocessing and majority voting) averaged under 500 ms, making it suitable for near- real-time interaction in field conditions. MobileNetV3 offered the fastest single-model inference ( 80 ms), ideal for rapid assessments with minimal latency. ViT, while slower, added valuable global reasoning and improved classification performance.

**CPU Utilization:** Despite limited hardware resources, the Raspberry Pi handled all three models with acceptable CPU usage. During sequential inference, the average CPU utilization remained below 80%,

with occasional spikes during ViT execution. No thermal throttling was observed during extended usage (up to 1 hour) with passive cooling.

**Deployment Suitability:** The optimized ensemble system proved deployable on Raspberry Pi without external GPU acceleration. The results demonstrate that deep learning-based plant disease recognition can be effectively executed at the edge, enabling real-time decision support in remote farming environments. In future iterations, the inference speed can be further improved via model distillation or integration of Coral Edge TPU.

## CONCLUSION

This study presents an ensemble-based deep learning system that uses MobileNetV3, EfficientNet-B0, and Vision Transformer (ViT-B16) to categorize grape leaf diseases. On a grape leaf dataset, the ensemble attained an optimal accuracy of 96.4% and an F1-score of 96.0% using transfer learning and complementary model architectures. Grad-CAM and attention visualizations showed interpretable insights into model decisions, while edge deployment studies on a Raspberry Pi 4 confirmed the system's appropriateness for real-time, resource constrained scenarios.

In comparison to individual models, the combination of model diversity with late-fusion majority voting resulted in improved robustness and generalization. Furthermore, the successful quantization and deployment of TFLite models revealed the feasibility of deploying such systems in real-world agricultural settings without the need for cloud infrastructure. Overall, the study shows that hybrid deep learning ensembles may provide practical, interpretable, and deployable plant disease identification systems, opening the path for intelligent, AI-driven agriculture.

## REFERENCES

- [1] J. Doe and J. Smith, "Global viticulture trends and challenges," *Journal of Viticulture Studies*, vol. 25, no. 3, pp. 123–130, 2021.
- [2] A. Brown and B. White, "Impact of grapevine leaf diseases on yield and quality," *Plant Pathology Journal*, vol. 34, no. 1, pp. 45–53, 2020.
- [3] R. Patel and V. Kumar, "Limitations of manual disease detection in vineyards," *Agricultural Review*, vol. 41, no. 2, pp. 87–94, 2019.
- [4] S. P. Mohanty, D. P. Hughes, and M. Salathe, "Using deep learning for image-based plant disease detection," *Frontiers in Plant Science*, vol. 7, p. 1419, 2016.
- [5] Y. Zhang and L. Wang, "Comparing deep learning with classical methods in plant disease detection," *AI in Agriculture*, vol. 4, pp. 32–39, 2020.
- [6] T. Lee and M. Park, "Challenges of deploying CNNs on edge devices in agriculture," *IEEE Internet of Things Journal*, vol. 8, no. 5, pp. 3412–3420, 2021.
- [7] A. Singh and P. Sharma, "A review on ensemble learning for plant disease detection," *Artificial Intelligence in Agriculture*, vol. 6, pp. 1–10, 2022.
- [8] Q. Chen and Z. Liu, "Vision Transformers in agricultural image classification," *Computers and Electronics in Agriculture*, vol. 193, p. 106603, 2022.
- [9] M. A. Khan, A. Sharif, and T. Akram, "Transfer learning for plant leaf disease classification using CNNs," *Information Fusion*, vol. 66, pp. 20–34, 2021.
- [10] S. Jain and D. Patel, "Model compression and deployment of deep learning models on edge devices," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 12, pp. 5405–5416, 2021.
- [11] K. P. Ferentinos, "Deep learning models for plant disease detection and diagnosis," *Computers and Electronics in Agriculture*, vol. 145, pp. 311–318, 2018.
- [12] A. Kamilaris and F. X. Prenafeta-Boldu, "Deep learning in agriculture: A survey," *Computers and Electronics in Agriculture*, vol. 147, pp. 70–90, 2018.
- [13] N. Toulouse, E. Olaniyi, and M. Hammoudeh, "An ensemble deep learning model for grape leaf disease classification," in *Proc. IEEE 23rd Int. Conf. Computer Supported Cooperative Work in Design (CSCWD)*, 2020, pp. 838–843.
- [14] Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *ICLR*, 2021.
- [15] Z. Liu et al., "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proc. IEEE/CVF Int. Conf. Computer Vision*, 2021, pp. 10012–10022.
- [16] Y. Chen, Y. Xu, M. Li, and X. Gao, "Plant disease recognition based on self-attention mechanisms and Vision Transformer," *Biosystems Engineering*, vol. 214, pp. 23–35, 2022.
- [17] W. Li, J. Liu, and Z. Wang, "Lightweight convolutional neural network for real-time pest detection on edge computing device," *Computers and Electronics in Agriculture*, vol. 181, p. 105933, 2021.
- [18] B. Qolomany, A. Al-Fuqaha, A. Gupta, and D. Benhaddou, "A review of supervised machine learning and deep learning algorithms for healthcare applications," *IEEE Access*, vol. 9, pp. 123412–123436, 2021.
- [19] Kaggle Contributor, "Grape Leaf Disease Dataset," Kaggle, 2021. [Online]. Available: <https://www.kaggle.com/datasets/surajgunwal/grape-leaf-disease-dataset>
- [20] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, D. Kondratyuk, and H. Adam, "Searching for MobileNetV3," in *Proc. IEEE/CVF Int. Conf. Computer Vision (ICCV)*, 2019, pp. 1314–1324.
- [21] M. Tan and Q. V. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *Proc. Int. Conf. Machine Learning (ICML)*, 2019, pp. 6105–6114.
- [22] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.