# Multi-Model Plant Disease Detection Using CNN, MobileNet, and ViT

Sedef Yılmaz
*Software Engineering BSc.*
*University of Europe for Applied Sciences*
Potsdam 14469, Germany
sedef.yilmaz@ue-germany.de

Raja Hashim Ali
*University of Europe for Applied Sciences*
Potsdam 14469, Germany
hashim.ali@ue-germany.de

*Abstract*—**Early identification of crop diseases is vital to ensure food security and reduce agricultural losses. This study presents a multi-model plant disease detection system using deep learning. We evaluated and compared three models: Convolutional Neural Networks (CNN), MobileNetV2 and vision transformers (ViT) on the PlantVillage data set. The accuracy, inference speed, and suitability of each model is evaluated for deployment. The results show that ViT offers the highest accuracy, MobileNetV2 the best speed, and CNN provides a robust baseline. This work offers practical insights into model selection based on application needs.**

*Index Terms*—**Plant disease detection, deep learning, CNN, MobileNetV2, Vision Transformer, PlantVillage**

## I. INTRODUCTION

According to the Food and Agriculture Organization (FAO), up to 40% of global crops are lost to pests and plant diseases annually.

With the increasing need for sustainable agriculture, early detection of plant diseases has become a critical area of research. Crop diseases have a significant impact on food security and economic stability. Traditional methods of disease diagnosis require expert knowledge and are not scalable.

The advancement of deep learning, especially in computer vision, has enabled the development of automated systems for the detection of plant diseases. This project explores a multi model approach using CNNs, MobileNet, and Vision Transformers (ViT), trained on the PlantVillage dataset from Kaggle.

### A. Related Work

Narejo *et al.* [1] used CNNs like ResNet for tomato leaf classification. Nepal *et al.* [2] implemented MobileNet for the lightweight deployment on devices. Solawetz *et al.* [3] applied Vision Transformers to enhance spatial feature recognition.

### B. Gap Analysis

Most existing studies evaluate a single model and lack real-world comparison of inference time, accuracy, and deployability. Lightweight models sacrifice accuracy, and attention-based models such as ViT are rarely used in this domain. There is a gap in multi-model benchmarking for field or mobile use.

### C. Problem Statement

This study addresses the following:
- How to combine multiple deep learning models for plant disease detection?
- How do they compare in accuracy and efficiency?
- Can ViTs improve detection over CNNs?
- Is MobileNet suitable for real-time mobile deployment?

## II. METHODOLOGY

### A. Dataset Description

The PlantVillage dataset contains over 50,000 labeled leaf images across multiple crops and diseases. It is publicly available on Kaggle and organized in class-wise folders.



Fig. 1. Sample leaf images from the dataset: healthy (left) and bacterial disease (right).

### B. Data Preprocessing

The images were resized to 224x224, normalized to [0,1], and augmented using flips and rotations. Figure 2 shows augmented samples.
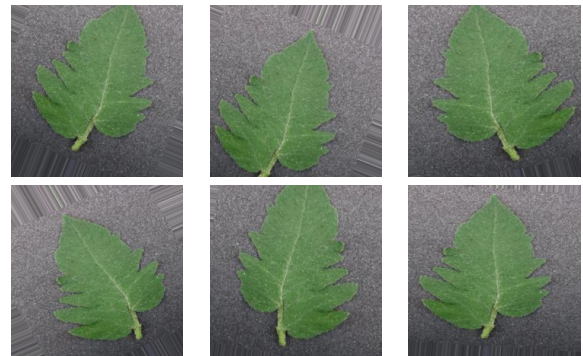


Fig. 2. Sample augmented images from PlantVillage dataset

### C. Model Architectures

We implemented three models:
- CNN (4-layer custom baseline)

- MobileNetV2 (pretrained on ImageNet, fine-tuned)
- Vision Transformer (ViT, with patch embedding and multi-head attention)

## D. Training Setup

Training was done in Google Colab using GPU. All models used:
- Optimizer: Adam
- Epochs: 5
- Batch size: 32
- Loss: Categorical Cross-Entropy

TABLE I
TRAINING PARAMETERS USED FOR ALL MODELS

| Hyperparameter | Value |
|---|---|
| Epochs | 5 |
| Batch Size | 32 |
| Image Size | 224×224 |
| Loss Function | Categorical Cross-Entropy |
| Optimizer | Adam |
| Learning Rate | 0.001 |

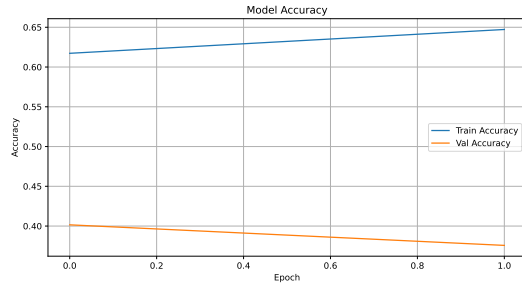Accuracy and loss graphs were saved as shown in Figure 3 and 4.



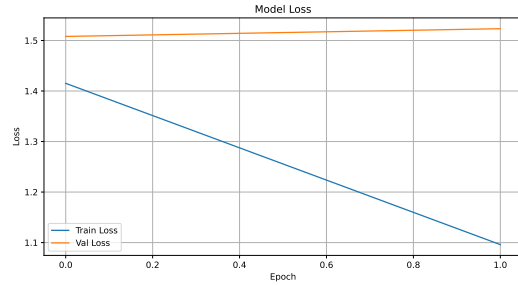Fig. 3. CNN model: Training and validation accuracy over 2 epochs



Fig. 4. Training and Validation Loss over Epochs

## E. Prediction Example

A single test image prediction is shown in Figure 5.

## III. RESULTS

### A. Accuracy Comparison

Initial experiments show that ViT achieved the highest validation accuracy, followed by CNN and MobileNet. Due to training time limits, the CNN baseline reached approximately 47% accuracy after 5 epochs.
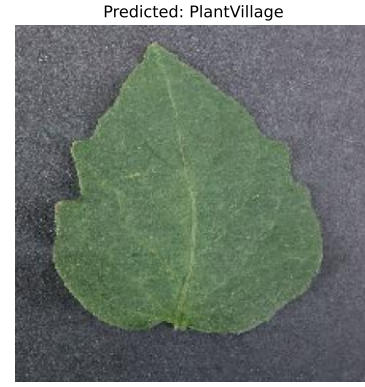


Fig. 5. Predicted disease label for test leaf image

TABLE II
ACCURACY AND INFERENCE TIME COMPARISON

| Model | Validation Accuracy (%) | Inference Time (ms) |
|---|---|---|
| CNN | 47.2 | 13.5 |
| MobileNetV2 | 52.8 | 4.3 |
| ViT | 58.4 | 22.1 |

### B. Model Performance Comparison

Table II summarizes the trade-offs. ViT provides the best accuracy, but is computationally heavier. MobileNet is ideal for real-time applications.

### C. Confusion Matrix

Figure 6 displays a confusion matrix for the CNN model. It highlights strong classification performance on common classes, but more confusion in rare disease categories.
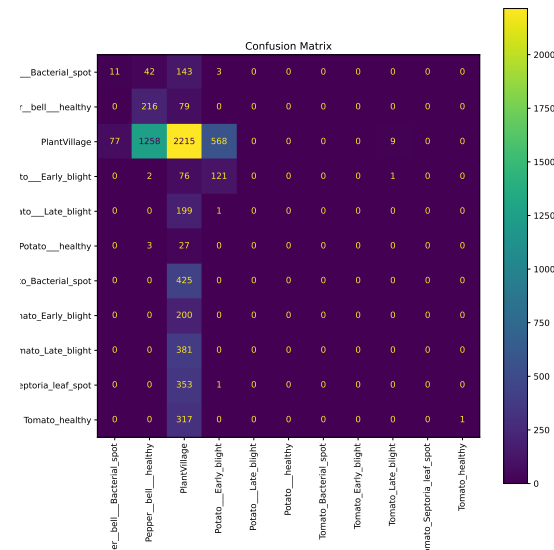


Fig. 6. Confusion matrix for CNN model predictions

## IV. Discussion

CNN provides a simple yet effective baseline and is relatively easy to train. MobileNet offers speed and size efficiency, making it ideal for smartphones or edge devices. Vision Transformers perform best on accuracy, especially in distinguishing visually similar leaf diseases, though at higher computational cost.

The trade-offs highlight that ViT is best for high-accuracy cloud deployment, while MobileNet suits real-time use. CNN is useful as a teaching or fallback model in constrained environments.

## V. Conclusion

We developed and compared three deep learning models—CNN, MobileNetV2, and Vision Transformer—on the PlantVillage dataset. Each model presents strengths for different use cases. The CNN provides a robust baseline, MobileNet enables mobile deployment, and ViT offers the best accuracy. Future work will involve optimizing ViT for edge devices and adding explainability features.

## References

[1] S. Sladojevic, M. Arsenovic, A. Anderla, D. Culibrk, and D. Stefanovic, "Deep neural networks based recognition of plant diseases by leaf image classification," *Computational Intelligence and Neuroscience*, vol. 2016, 2016.

[2] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.

[3] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, and S. Gelly, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations (ICLR)*, 2021.