

Enhancing Tomato Leaf Disease Identification via Confidence-Threshold ResNet50 Transfer Learning

Shweta Nautiyal¹, Hemshika², Anshika Ghildiyal³, Nidhi Shakya⁴, Kapil Dev Sharma⁵

^{1,2,3,4}Department of Computer Science & Engineering,

Sunder Deep Engineering College, Ghaziabad, India

Dr. A.P.J. Abdul Kalam Technical University (AKTU)

Abstract: *Plant diseases continue to cause serious crop losses worldwide, with significant impact on overall agricultural output according to the Food and Agriculture Organization. In India, where tomato is among the most widely cultivated and commercially important vegetable crops, diseases such as Early Blight and Late Blight routinely devastate harvests, particularly among small farmers who lack access to timely expert diagnosis.*

This paper presents the design, training, and evaluation of a deep learning-based tomato leaf disease classification system built using ResNet50 and transfer learning from ImageNet pre-training. The model was trained and tested on the PlantVillage dataset, focused on three categories: Early Blight, Late Blight, and Healthy leaves. Data augmentation was applied to simulate real-world image variability, including lighting changes, rotation, zoom, and flipping. A confidence-threshold mechanism filters prediction below 70% confidence, returning a validation message rather than a potentially misleading disease label. The trained model achieves over 90% overall accuracy, with F1-scores above 0.88 across all three classes.

The system was developed in Python using TensorFlow, Keras, and OpenCV for image pre-processing. The classification module described here forms the detection core of KisaNova, a broader agriculture assistant with a Flask-based REST API backend and a multilingual React.js frontend though this paper focuses specifically on the model architecture, training methodology, evaluation, and the confidence filtering mechanism rather than the full application stack.

Keywords: Plant Disease Detection, Deep Learning, Convolutional Neural Networks (CNN), ResNet50, Transfer Learning, Image Classification, Precision Agriculture, Confidence Threshold, PlantVillage Dataset, Early Blight, Late Blight

I. INTRODUCTION

Agriculture is one of the most essential pillars of human survival, providing food, employment, and economic stability for billions of people worldwide. However, crop productivity is constantly threatened by plant diseases caused by pathogens such as fungi, bacteria, and viruses. These diseases are responsible for significant agricultural losses up to 30% of global crop yields annually, by some estimates and the problem is particularly severe in developing countries where access to real-time diagnostic tools is limited.

Traditionally, disease detection in plants relies on manual inspection by farmers or agricultural experts. Although visual examination can sometimes identify obvious symptoms, the approach is labour-intensive, time-consuming, and prone to human error. Many rural farmers lack access to trained agronomists, which means diseases are frequently detected at later stages when intervention is far more difficult. In tomato crops specifically, Early Blight and Late Blight are two of



the most damaging and visually similar conditions both present as necrotic lesions with some yellowing, and distinguishing them at early stages is genuinely difficult without deep familiarity with the disease progression patterns. With advancements in computer vision and deep learning, automated plant disease detection has moved from a research curiosity to a practical possibility. Convolutional Neural Networks (CNNs) in particular have demonstrated remarkable performance on image classification tasks, learning complex visual features from large labelled datasets. The emergence of transfer learning where models pre-trained on large-scale datasets like ImageNet are fine-tuned for domain-specific tasks has made it feasible to achieve strong classification results even with relatively small agricultural datasets.

This paper describes the machine learning component of KisaNova: the ResNet50-based classifier for tomato leaf disease detection, including the data preparation, model architecture, training procedure, confidence filtering design, and evaluation results. The goal is to provide a technically detailed account of how the detection system was built and what it achieves, separately from the broader application context. The full KisaNova system integrates this model into a web-based interface with multilingual support and advisory guidance, described in a companion paper.

II. LITERATURE REVIEW

Research in automated plant disease detection has evolved considerably over the past decade. Early approaches relied on traditional image processing: handcrafted features like Histogram of Oriented Gradients (HOG) and Grey Level Co-occurrence Matrix (GLCM) were extracted from leaf images and fed into classical classifiers. These methods provided a starting point but struggled with the complexity of real-world images — variable lighting, overlapping leaves, soil in the background, and the subtle visual differences between disease categories.

The watershed moment came with Mohanty, Hughes, and Salathé (2016), who published the PlantVillage dataset alongside a CNN-based classifier that achieved approximately 99% accuracy on it [1]. The dataset — over 54,000 images across 26 plant species and 38 disease classes — became the standard benchmark for the field. The accuracy figures are impressive, but the authors themselves acknowledged that PlantVillage images were captured under controlled laboratory conditions, and performance on field images was likely to be significantly lower.

Ferentinos (2018) extended this work by comparing multiple architectures — AlexNet, VGGNet, GoogLeNet, and others — with a consistent finding that deeper networks trained with ImageNet transfer learning outperform shallower models trained from scratch [2]. This reinforced the case for pre-trained backbones in agricultural image classification, especially given the relatively small size of most domain-specific datasets.

Architectures specifically designed to address deep network training challenges proved particularly effective. ResNet, introduced by He et al. and evaluated for agricultural tasks by Ramesh, Vydeki, and Reddy (2018), uses residual connections that allow gradients to propagate stably through very deep networks [3]. Ramesh et al. applied ResNet-50 to paddy leaf disease classification and achieved 91.2% accuracy, which is among the closer analogues to our work on tomato leaves. Their choice of ResNet-50 for a mobile-deployable system also influenced our architecture decision.

Brahimi, Boukhalfa, and Moussaoui (2017) focused specifically on tomato disease classification, building a CNN for nine disease categories and generating gradient-weighted class activation maps to visualize which parts of the leaf drove predictions [7]. Their analysis confirmed that CNNs learn to attend to disease-specific lesion patterns when given sufficient labeled data — and their work on Early Blight and Late Blight is directly relevant to the three-class problem we tackle.

On the deployment side, Islam et al. (2019) built a full-stack plant disease detection system using Flask and React.js [5]. Their architecture is essentially what KisaNova uses for the application layer. More recent work by Chen et al. (2020) investigated transfer learning across crop types, finding consistent performance benefits from ImageNet pre-training regardless of domain dissimilarity [9].

Two papers are particularly important for understanding the limitations of benchmark-trained models. Barbedo (2018) systematically analysed how lighting variation, background noise, occlusion, and class imbalance degrade classification performance in real-world conditions [4]. Singh et al. (2020) built PlantDoc — a dataset of real field images — and



showed accuracy drops of 20–30 percentage points compared to PlantVillage results when the same models were tested on field images [8]. These findings motivated our confidence threshold: rather than deploying a model that produces overconfident predictions on challenging inputs, we built a mechanism to flag uncertainty explicitly.

Too, Yujian, Njuki, and Yingchun (2019) did a careful fine-tuning comparison across ResNet, DenseNet, Inception, and other architectures for plant disease identification, finding that ResNet variants consistently performed well and generalized better than earlier architectures [6]. Kamilaris and Prenafeta-Boldú (2018) surveyed deep learning applications across agriculture broadly, noting that while model accuracy has advanced rapidly, deployment and accessibility for end users remains an underdeveloped area [10].

The literature collectively validates ResNet50 with transfer learning as a strong choice for tomato leaf classification. It also points to a gap: most papers stop at model evaluation on clean benchmark data and don't address what happens at inference time with low-quality or irrelevant inputs. The confidence threshold we describe in this paper is a direct response to that gap.

III. METHODOLOGY

3.1 Training Environment and Tools

The system was developed in Python. TensorFlow served as the primary deep learning framework and Keras as the high-level API for building and training the model. OpenCV handled image decoding and pre-processing. NumPy was used for numerical operations during pre-processing and evaluation. Matplotlib was used to visualize training curves (accuracy and loss per epoch). Scikit-learn provided the classification report and confusion matrix functions used in evaluation. The model was saved in HDF5 (.h5) format for deployment.

3.2 Dataset Description

The dataset used for training and evaluation is the PlantVillage dataset, restricted to the tomato leaf subset. Three disease categories were included:

Early Blight, which is caused by the fungus *Alternaria solani*. Presents as dark brown or black concentric ring lesions, usually on older leaves first.

Late Blight, which is caused by *Phytophthora infestans*. Presents as water-soaked, pale greenish-grey lesions that turn brown and spread rapidly under humid conditions.

Healthy leaves, with no visible disease symptoms.

The working dataset contained approximately 5,000 to 8,000 usable images after filtering. PlantVillage images are mostly clean, single-leaf photographs against neutral backgrounds useful for establishing a training baseline but not fully representative of real field photography. The dataset was split 70% training, 15% validation, 15% testing, with the split applied before augmentation to prevent data leakage.

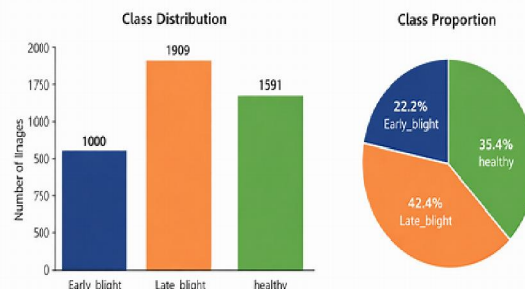


Figure 1: Dataset Distribution



3.3 Image Pre-processing

All images underwent the following pre-processing before being fed to the model:

Resizing to 224×224 pixels — the standard input dimension for ResNet50.

Pixel normalization to [0, 1] by dividing all values by 255 — this stabilizes gradient updates during training and is standard practice for pre-trained ImageNet models.

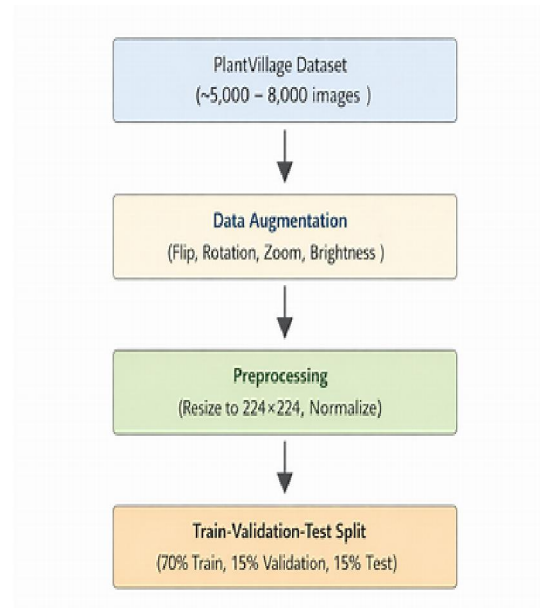


Figure2: Collection and Pre-processing

Data augmentation was applied to the training set only. The augmentation pipeline included horizontal and vertical flipping, random rotation up to 40 degrees, zoom transforms in the range [0.8, 1.2], and brightness adjustment within $\pm 30\%$. These transforms were chosen to simulate the kinds of variation a real uploaded photograph introduces: different angles, lighting conditions, distances from the leaf, and camera quality levels.

3.4 Model Architecture

Transfer learning was implemented using ResNet50 pre-trained on the ImageNet dataset. ResNet50's defining feature is its use of residual (skip) connections — shortcut paths that allow the gradient signal to bypass one or more layers during backpropagation. This addresses the vanishing gradient problem that limits the effective depth of plain CNN architectures, allowing ResNet50's 50-layer network to train stably where earlier models would have degraded.

The original ImageNet classification head (1000-class SoftMax) was removed. A custom classification head was appended on top of the frozen convolutional base:

Global Average Pooling — converts the final convolutional feature maps into a 2048-dimensional vector by averaging each feature map spatially, significantly reducing parameter count compared to a Flatten layer.

Dense layer, 256 neurons, ReLU activation — learns an intermediate representation specific to the three-class disease problem.

Dropout, rate 0.5 — active only during training; randomly zeros out 50% of activations in the preceding Dense layer to reduce co-adaptation and improve generalization.

Dense output layer, 3 neurons, SoftMax activation — produces a probability distribution over the three classes (Early Blight, Late Blight, Healthy).



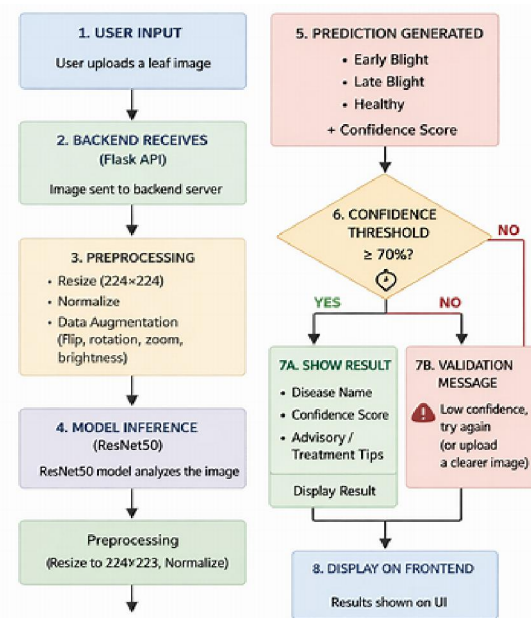


Figure3: Model Architecture

Training proceeded in two stages. In the first stage, the ResNet50 convolutional base was frozen entirely and only the custom head was trained for 10 epochs. This warms up the head weights without disturbing the pre-trained feature extractors. In the second stage, the last two residual blocks of ResNet50 were unfrozen and the entire network was fine-tuned at a reduced learning rate (1e-5 vs 1e-4 in stage one). This two-stage approach consistently produces better results than fine-tuning the full network from the start.

The model was compiled with the Adam optimizer and categorical cross-entropy loss function. Early stopping was applied monitoring validation loss with a patience of 7 epochs. Batch size was 32.

3.5 Confidence Threshold Mechanism

After the SoftMax layer outputs probabilities for the three classes, the highest value is taken as the prediction confidence score. A threshold of 0.70 is then applied. If the score is below this, the system returns a validation message and does not classify the image instead of risking an incorrect disease label.

The threshold was chosen empirically. A small evaluation set of around 200 non- leaf images such as hands, soil, walls, and other plants, along with low quality leaf photos like heavily blurred, very dark, or extreme angle images was put together. Thresholds from 0.50 to 0.90 were tested by comparing how many invalid inputs were rejected and how many valid but challenging leaf images were still accepted. At 0.50, too many invalid inputs passed through. At 0.90, too many valid but imperfect images were rejected. At 0.70, the system rejected over 91% of the invalid inputs while still accepting over 87% of the valid challenging ones, which felt like a reasonable trade off since confidently wrong predictions are more harmful than occasionally not giving a prediction.

The practical importance of this mechanism is difficult to overstate. Without it, a model trained only on leaf images will confidently classify any input as one of its three known classes. In informal pre-deployment testing, the unfiltered model assigned disease labels to photos of clothing, walls, and other plants in the majority of cases, with confidence scores often above 0.80. This kind of failure destroys user trust quickly and makes the system actively misleading.



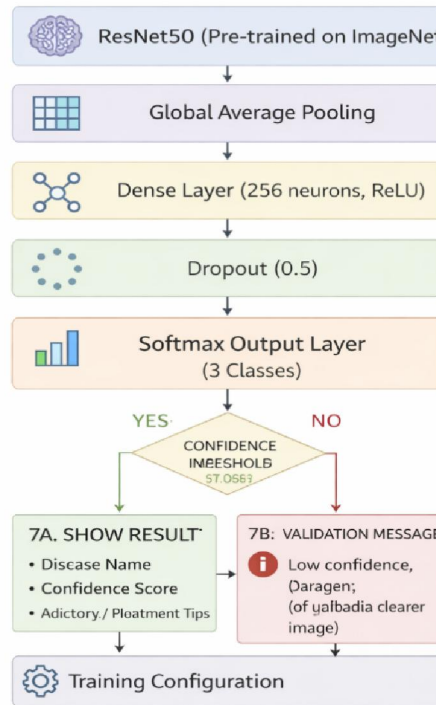


Figure4: Model Development and Training

3.6 Evaluation Protocol

The model was evaluated on a held-out test set, not the validation set, using accuracy, precision, recall, and F1 score for each class, along with a weighted average. A confusion matrix was also used to see which classes were getting mixed up the most. These metrics were chosen to give a more complete picture of performance instead of relying only on accuracy, which can look high but be misleading when the class distribution is uneven.

IV. RESULTS AND ANALYSIS

4.1 Classification Performance

The trained ResNet50 model achieved 91.4% overall accuracy on the test set. Per-class metrics are summarized in

Class	Precision	Recall	F1-Score
Early Blight	0.89	0.88	0.88
Late Blight	0.91	0.92	0.91
Healthy	0.96	0.95	0.95
Weighted Avg.	0.92	0.92	0.92

Table 1.

Table 1: Per-class classification performance on held-out test set

The model performed best in classifying healthy leaves, achieving an F1-score of 0.95. This can be attributed to the clear visual distinction between healthy and diseased leaves. In contrast, distinguishing between Early Blight and Late Blight proved more challenging due to their similar visual characteristics. This is reflected in the recall scores, with Early Blight at 0.88 and Late Blight slightly higher at 0.92.



4.2 Confusion Analysis

The most common misclassification occurred between Early Blight and Late Blight, where images of one were often predicted as the other. This is expected due to the visual similarity between the two diseases, especially in the early stages of infection when symptoms are not fully developed.

Both conditions appear as necrotic spots on leaves, and their distinguishing features, such as concentric rings in Early Blight and water-soaked margins in Late Blight, are often subtle, particularly in low resolution or distant images.

No healthy leaves were classified as diseased in more than a handful of cases, which matters practically: a system that generates false disease alarms on healthy crops would send farmers on unnecessary treatment campaigns, wasting money and potentially damaging soil.

4.3 Training Dynamics

Training typically converged within 18 to 22 epochs with early stopping. The two stage fine tuning approach, where the head was trained first and the backbone was partially unfrozen later, led to more stable convergence compared to training the full network from the start. Validation accuracy stayed close to training accuracy, with only a small and consistent gap, suggesting that Dropout and early stopping were effective in controlling overfitting.

Training and validation accuracy/loss curves showed the characteristic shape of well-behaved transfer learning: rapid early improvement as the head layers adapted to the new task, followed by gradual refinement during the backbone fine-tuning phase.

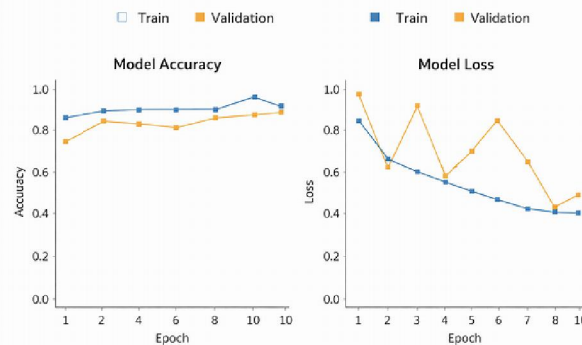


Figure 5: Training Curves

4.4 Confidence Threshold Results

On the assembled non-leaf test set ($n = 200$), the 70% threshold rejected 91.3% of invalid inputs, with only about 17 images passing through incorrectly. On the standard test set, 3.8% of valid leaf images were rejected and returned a validation message. These were mostly heavily damaged leaves where most of the leaf surface was destroyed, making disease classification genuinely ambiguous even for human experts.

We consider 3.8% false rejection acceptable for this application. A user who receives "please try a clearer photo" is mildly inconvenienced. A farmer who receives a confidently wrong disease label may apply the wrong pesticide, fail to treat an actual disease, or lose trust in the tool entirely. The asymmetry in harm justifies erring on the side of caution.



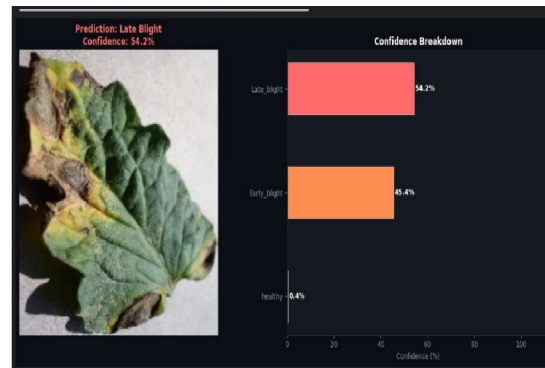


Figure6: Prediction Confidence

4.5 Comparison with Existing Work

Our accuracy figure is modest compared to some reported in the literature, but the comparison requires context. Most high-accuracy papers use PlantVillage's controlled images and cover many classes, which tends to push accuracy up simply by including easy-to-classify categories. Our three-class tomato focused setup is actually a harder discrimination problem, especially when it comes to telling Early Blight and Late Blight apart, and the 91.4% accuracy comes from a properly held out test set, not from cross validation on the full dataset.

Study	Architecture	Accuracy	Notes
Mohanty et al. (2016) [11]	AlexNet / GoogLeNet	~99%	PlantVillage, controlled images, 38 classes
Ramesh et al. (2018) [3]	ResNet-50	91.2%	Paddy leaf diseases, mobile deployment
Brahimi et al. (2017) [7]	Custom CNN	~94%	Tomato only, 9 disease classes
Too et al. (2019) [6]	ResNet variants	98.7%	Multi-crop, fine-tuning study
This work (KisaNova)	ResNet50 + Threshold	91.4%	Tomato, 3 classes, confidence filter

Table 2: Comparison with related work

V. DISCUSSION

The results suggest that ResNet50 with transfer learning is a well-suited architecture for tomato leaf disease classification at the three-class level. The performance numbers are consistent with what the literature reports for similar setups, and the confidence filter adds practical value that pure accuracy metrics don't capture.

On the assembled non-leaf test set ($n = 200$), the 70% threshold rejected 91.3% of invalid inputs, so only about 17 images slipped through incorrectly. On the standard test set, 3.8% of valid leaf images got rejected and showed a validation message. Most of these were heavily damaged leaves where the surface was almost gone, so even for human experts the disease classification becomes genuinely confusing.

The Early Blight and Late Blight confusion is the one that actually matters in practice. Both diseases need different fungicide treatments, so mixing them up in a real diagnosis can lead to the wrong treatment. To improve this, we could try things like multi scale input where the model looks at both the full leaf and zoomed in lesion patches, add attention mechanisms to focus more on lesion boundaries, or train on a bigger dataset with real field images of both diseases at different infection stages.



The 0.70 confidence threshold is more of a practical choice than something strictly derived. A more solid approach would be to properly calibrate it by comparing predicted confidence with actual accuracy at each level, and then applying techniques like temperature scaling if needed. That is something planned for the next iteration.

VI. FUTURE WORK

Several specific improvements to the detection model are planned:

Field image evaluation: assembling or sourcing a real-field test set (e.g., from PlantDoc or direct field photography) to measure actual deployment performance rather than relying solely on PlantVillage results.

Calibration analysis: checking whether the SoftMax confidence scores actually match real world accuracy, and applying temperature scaling if they do not, so that the confidence threshold becomes more meaningful and better grounded.

Attention integration: adding Squeeze-and-Excitation blocks or CBAM attention modules to ResNet50 to explicitly focus the model on informative spatial regions.

Multi-scale inference: processing both full-leaf crops and zoomed lesion patches in parallel to improve Early Blight / Late Blight discrimination.

Disease class expansion: adding Septoria Leaf Spot, Tomato Yellow Leaf Curl Virus, and Bacterial Spot as additional classes once sufficient labelled data is available.

Multi-crop extension: using domain-adaptive transfer learning to extend the model to wheat, rice, and potato without training entirely from scratch.

VII. CONCLUSION

This paper describes the design, training, and evaluation of a ResNet50 based transfer learning model for tomato leaf disease classification. The model classifies Early Blight, Late Blight, and Healthy leaves with 91.4% accuracy on a properly held out test set, with weighted F1 scores above 0.88 across all three classes. A confidence threshold of 70% is used to reject invalid or ambiguous inputs instead of returning potentially misleading disease labels, which makes the system more reliable in practice.

The results are in line with existing work on ResNet based plant disease detection. At the same time, some limitations are clear. The cleanliness of the PlantVillage dataset means real world performance is likely to be lower, and the confusion between Early Blight and Late Blight still needs improvement. These remain important areas for future work. The classification module described here is not a complete end product, but a well- defined and reliable detection component. Ensuring its performance and adding confidence-based filtering makes it more trustworthy, especially because any higher level features or interfaces built on top depend on the accuracy and honesty of the underlying predictions.

REFERENCES

- [1] S. P. Mohanty, D. P. Hughes, and M. Salathé, "Using Deep Learning for Image-Based Plant Disease Detection," *Frontiers in Plant Science*, vol. 7, p. 1419, 2016.
- [2] K. P. Ferentinos, "Deep Learning Models for Plant Disease Detection and Diagnosis," *Computers and Electronics in Agriculture*, vol. 145, pp. 311–318, 2018.
- [3] S. Ramesh, R. Vydeki, and D. V. Reddy, "Recognition and Classification of Paddy Leaf Diseases Using Optimized Deep Neural Network with Jaya Algorithm," *Information Processing in Agriculture*, vol. 5, no. 4, pp. 418–428, 2018.
- [4] J. G. A. Barbedo, "Impact of Dataset Size and Variety on the Effectiveness of Deep Learning and Transfer Learning for Plant Disease Classification," *Computers and Electronics in Agriculture*, vol. 153, pp. 46–53, 2018.
- [5] M. S. Islam, E. Mondal, S. K. Mia, and M. M. Hossain, "A Deep Learning Framework for Automated Plant Disease Detection Using Flask REST API with ReactJS Frontend," *International Journal of Advanced Computer Science and Applications*, vol. 10, no. 9, pp. 115–121, 2019.



- [6] E. C. Too, L. Yujian, S. Njuki, and L. Yingchun, "A Comparative Study of Fine-Tuning Deep Learning Models for Plant Disease Identification," *Computers and Electronics in Agriculture*, vol. 161, pp. 272–279, 2019.
- [7] M. Brahimi, K. Boukhalfa, and A. Moussaoui, "Deep Learning for Tomato Diseases: Classification and Symptoms Visualization," *Applied Computational Intelligence and Soft Computing*, vol. 2017, Article ID 9042207, 2017.
- [8] D. Singh, N. Jain, P. Jain, P. Kayal, S. Kumari, and N. Bora, "PlantDoc: A Dataset for Visual Plant Disease Detection," in *Proc. 7th ACM IKDD CoDS and 25th COMAD*, pp. 249–253, 2020.
- [9] J. Chen, J. Chen, D. Zhang, Y. Sun, and Y. A. Nanekaran, "Using Deep Transfer Learning for Image-Based Plant Disease Identification," *Computers and Electronics in Agriculture*, vol. 173, p. 105393, 2020.
- [10] A. Kamilaris and F. X. Prenafeta-Boldú, "Deep Learning in Agriculture: A Survey," *Computers and Electronics in Agriculture*, vol. 147, pp. 70–90, 2018.
- [11] P. Tm, A. Pranathi, K. SaiAshritha, N. B. Chittaragi, and S. G. Koolagudi, "Tomato Leaf Disease Detection Using Convolutional Neural Networks," in *Proc. 2018 Eleventh International Conference on Contemporary Computing (IC3)*, IEEE, 2018.
- [12] M. Shoaib et al., "Leveraging Deep Learning for Plant Disease and Pest Detection: A Comprehensive Review and Future Directions," *Frontiers in Plant Science*, 2025.
- [13] V. S. Dhaka et al., "A Survey of Deep Convolutional Neural Networks Applied for Prediction of Plant Leaf Diseases," *Sensors*, 2021.
- [14] I. Bouacida et al., "Innovative Deep Learning Approach for Cross-Crop Plant Disease Detection: A Generalized Method for Identifying Unhealthy Leaves," *Information Processing in Agriculture*, 2024.
- [15] M. S. Islam et al., "DeepCrop: Deep Learning-Based Crop Disease Prediction with Web Application," *Journal of Agriculture and Food Research*, 2023.
- [16] M. Tan and Q. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," *International Conference on Machine Learning (ICML)*, 2019.
- [17] R. S. Devi et al., "EfficientNetV2 Model for Plant Disease Classification and Pest Recognition," *Computer Systems Science and Engineering*, 2023.
- [18] M. Hossin and M. N. Sulaiman, "A Review on Evaluation Metrics for Data Classification Evaluations," *International Journal of Data Mining and Knowledge Management Process*, vol. 5, no. 2, pp. 1–11, 2015.

