# Plant Leaf Image Classification Using SAM-Based Segmentation with Vision Transformers and ResNet Architectures

## Making AI an Expert of Recognizing Plants Disease

**Rongyi Shen**
rongyish@usc.edu
**Xiao Bai**
xiaobai@usc.edu
**Wenjing Huang**
whuang08@usc.edu

## Abstract

*Plant diseases significantly impact global agriculture, reducing crop yields and threatening food security, yet traditional detection methods are slow, labor-intensive, and prone to error. To address this, we developed a hybrid pipeline that integrates the Segment Anything Model (SAM) with Vision Transformers (ViT) and ResNet architectures for accurate and robust plant disease detection. SAM was used to isolate disease-relevant regions in images, improving data quality and reducing noise, while ViT and ResNet were fine-tuned using both original and segmented datasets from the PlantVillage dataset. Our approach achieved a top-1 test accuracy of 99.95% with ViT and 99.81% with ResNet on the original dataset, and 99.55% and 99.40%, respectively, on SAM-segmented data. Robustness tests further demonstrated ViT's superior generalization under noisy conditions. These results underscore the potential of combining segmentation with advanced vision models to improve the accuracy and scalability of plant disease detection, providing a powerful tool for sustainable agricultural practices.*

## 1. Introduction

### 1.1 Problem Statement

Plant diseases significantly affect global agriculture, leading to reduced crop yields and posing a threat to food security. Timely and accurate detection of plant diseases is crucial to mitigate these effects. However, traditional methods of disease identification often rely on manual inspection, which is time-consuming, labor-intensive, and prone to human error. With the increasing global population and the need to enhance food production, a scalable and efficient solution for plant disease detection has become essential. The adoption of machine learning models for automating this process offers great potential, but existing models are often limited by insufficient training datasets and their inability to handle variability in image quality and environmental noise.

If this problem remains unsolved, it will hinder the development of sustainable agricultural practices, potentially exacerbating food insecurity in vulnerable regions. Furthermore, inefficient disease detection methods could lead to the excessive use of pesticides, harming the environment and human health.

### 1.2 Motivation

Recent advances in machine learning and computer vision, such as Vision Transformers (ViTs) and convolutional networks like ResNet, have shown remarkable success in image classification. These models can significantly enhance the precision of plant disease detection when paired with robust datasets and innovative augmentation techniques. However, challenges such as dataset diversity and the ability to generalize across unseen conditions persist. Addressing these challenges is vital to develop a system that not only achieves high accuracy but also ensures robustness in real-world applications.

The Segment Anything Model (SAM) provides a novel approach to segmenting and preprocessing images, enhancing data quality and enabling better model training. Integrating SAM with deep learning architectures like ViT and ResNet could pave the way for highly accurate and efficient disease detection systems. This approach aligns with recent research emphasizing the importance of dataset augmentation and advanced vision models for improving agricultural diagnostics.

### 1.3 Proposed Approach

We propose a hybrid pipeline integrating Vision Transformers (ViT) and ResNet with the Segment Anything Model (SAM) to address plant disease detection challenges. The choice of ViT and ResNet stems from their complementary strengths: ViT's global attention mechanism excels in capturing high-level features, while ResNet's hierarchical structure effectively captures local patterns. SAM enhances the dataset by segmenting images to focus on disease-relevant regions, improving data diversity and robustness.

Our methodology involves training and fine-tuning both ViT and ResNet on the PlantVillage dataset, consisting of 20,638 images across 15 disease categories. The training process is augmented using techniques such as gradient blur and rotation to improve model generalization. We evaluate models on original and SAM-segmented datasets, exploring their performance under noise and testing their robustness.

Key experimental steps:
1. Baseline Zero-Shot Testing: Conducted with pre-trained ViT and ResNet models on the original dataset.
2. Model Training on Original Dataset: Fine-tuning the models with standard training techniques to establish benchmark accuracies.
3. Integration of SAM-Segmented Data: Incorporating SAM-segmented images into training and testing pipelines to analyze performance improvements.
4. Comparison of Models: Analyzing ViT and ResNet performance metrics, including accuracy, robustness under noise, and inference times.

### 1.4 Contributions

This study makes the following contributions:
1. Hybrid Pipeline Development: We integrate SAM with two advanced vision models (ViT and ResNet) to improve plant disease detection accuracy and robustness.
2. Performance Comparison of Vision Models: By comparing ViT and ResNet on both original and segmented datasets, we highlight their respective strengths and limitations, offering insights into their application in agricultural diagnostics.
3. Enhanced Data Augmentation Techniques: Incorporating gradient blur and rotation, we achieve better generalization and robustness, particularly in noisy environments.
4. Robustness and Efficiency: Our models demonstrate high robustness, with ViT achieving top-1 accuracy of 99.88% and ResNet achieving a comparable performance of 99.40%. The results underscore the effectiveness of SAM in augmenting plant disease datasets.


This research provides a foundation for scalable and efficient plant disease detection systems, contributing to the broader goal of sustainable agricultural practices.


## 2. Related Works

Deep learning has significantly advanced plant disease detection, offering improved accuracy and scalability compared to traditional methods. The use of Vision Transformers (ViTs) and Convolutional

Neural Networks (CNNs) has emerged as a key area of research, with various approaches addressing challenges such as dataset diversity, segmentation precision, and robustness to real-world conditions.

The Segment Anything Model (SAM) has garnered significant attention as a general-purpose segmentation model capable of object-specific adaptability across various domains. Zhang et al. (2023) conducted a comprehensive review of SAM's applications, highlighting its precision in segmentation with minimal input but also noting challenges related to dependency on high-quality prompts and difficulties with small or complex objects. Building on SAM's adaptability, Na et al. (2024) extended its capabilities to medical imaging by introducing an auto-prompting framework and Low-Rank Adaptation (LoRA) for fine-tuning domain-specific tasks, demonstrating enhanced performance with minimal computational resources. These findings inform our approach, where we leverage SAM's object-specific segmentation capabilities and explore the integration of prompt automation techniques to enhance accuracy in segmenting plant disease regions while managing computational load. Similarly, the SAM-Med2D project demonstrated SAM's strength in detailed object segmentation for medical imaging, though it faced challenges in generalization and high resource demands. Our work adapts these lessons to agricultural applications, combining SAM's precision with targeted augmentations to focus on disease-relevant regions in plant leaf images.

Recent studies have also explored hybrid models combining ViTs and CNNs to leverage their complementary strengths. De Silva and Brown (2023) demonstrated the utility of ResNet-50V2 in a multispectral approach, achieving accuracies of 98.35% for RGB datasets and 94.01% for NIR datasets, underscoring the potential of integrating spectral diversity. Similarly, Li and Li (2022) proposed the ConvViT hybrid model for detecting apple diseases, achieving 96.85% accuracy, which closely rivaled the Swin-tiny model's 96.94%. These approaches highlight the effectiveness of combining the global feature extraction capabilities of ViTs with the localized pattern recognition strengths of CNNs for plant disease classification.

Segmentation techniques have been integral to improving classification accuracy by isolating disease-relevant regions and reducing noise. Sahu and Minz (2023) developed an adaptive segmentation method combined with ResNet and LSTM–DNN for multi-disease classification, demonstrating improved precision compared to traditional approaches. The SAM-Med2D project further emphasized the importance of object-level segmentation for detailed image features, such as tumors in medical imaging, offering insights into the utility of segmentation for domain-specific tasks. In our work, SAM-based segmentation is utilized to isolate disease-relevant regions in plant images, addressing the limitations of conventional full-image processing by focusing on disease-specific features. By integrating SAM with ViT and ResNet, we aim to improve classification accuracy while minimizing the impact of irrelevant background information.

Data augmentation also plays a pivotal role in enhancing model performance, particularly when datasets are limited or domain-specific. Traditional augmentation techniques such as rotation and flipping enhance generalization but often fail to capture localized variations critical for fine-grained tasks. Alomar et al. (2023) introduced Random Local Rotation (RLR) as an augmentation technique that applies localized transformations to mimic real-world variations, significantly improving model robustness. Our study integrates SAM-based segmentation with object-specific augmentations inspired by Alomar et al., focusing transformations on disease-relevant regions to create a dataset more representative of real-world scenarios. This approach combines segmentation with advanced augmentation strategies, improving data diversity and model generalization across complex disease categories in plant leaves.

Multitask learning frameworks have further demonstrated potential for plant disease localization and classification. Hemalatha and Jayachandran (2024) proposed a multitask learning-based ViT model that achieved 99.97% accuracy by integrating co-scale, co-attention, and cross-attention mechanisms. These advancements demonstrate the potential of unified frameworks for addressing multiple objectives, which align with our efforts to enhance both segmentation and classification performance through integrated pipelines.

Our research builds on and extends the existing literature by introducing a novel integration of SAM with state-of-the-art vision models. By combining SAM's segmentation capabilities with targeted data augmentation and fine-tuning ViT and ResNet architectures, we address key challenges such as noise sensitivity, dataset variability, and computational efficiency. Our work contributes to advancing plant disease recognition by offering a scalable, robust, and efficient pipeline that leverages the strengths of segmentation and deep learning for real-world agricultural applications.

# 3. Methods

Our approach integrates the Segmentation Anything Model (SAM) for image segmentation with two visual language models: Vision Transformer and ResNet50 to create a system capable of recognizing diseased plant leaf images. We begin by using a dataset consisting of images of various plant leaves categorized based on plant type, health status, ancenterific diseases. It includes 15 subfolders, each representing a unique class, and we make them as categories for the images. This original dataset serves as the foundation for subsequent segmentation and two models training processes.

## 3.1 Segmentation Anything Model (SAM)

We first implement the SAM for images segmentation. SAM allows for various segmentation prompts, such as bounding boxes or point prompts, enabling us to isolate major objects or specific regions of interest. Combining these segmented images with the original dataset, we can enhance our model training with focused regions of our plant leaf. Specifically, the original images only capture the entire plant context, and our segmented images emphasize the more detailed diseased or healthy regions. We expect that using the segmented images alongside the original images will improve our model's ability to detect subtle disease features. In addition, it allows our model to differentiate between background noise (e.g., soil, other plants) and the essential features of disease recognition which will potentially give a more robust and accurate recognition.

To enhance our dataset to be more comprehensive, we implemented SAM's automatic mask generation for each image. The largest segmentation mask area was generated that captured the primary objects in each image, such as leaves or stems, along with any visible signs of disease. Additionally, we utilized meaningful prompts of bounding boxes to guide SAM in identifying key regions of interest. This prompt helped the model focus on specific areas, like diseased spots or healthy sections, ensuring that the segmentation accurately captured relevant features. The remaining detailed features such as the background color was isolated, this provided more standardized inputs that focused on diseased and healthy plant parts, enhancing the quality of the data available for training and encouraging the model to only capture the visible plant conditions rather than potentially confusing background elements. After isolating core objects, we apply targeted data augmentation techniques later on, focusing only on the segmented objects. We then apply augmentation to the segmented images and combine with the original images as the candidate dataset for the following model training part.
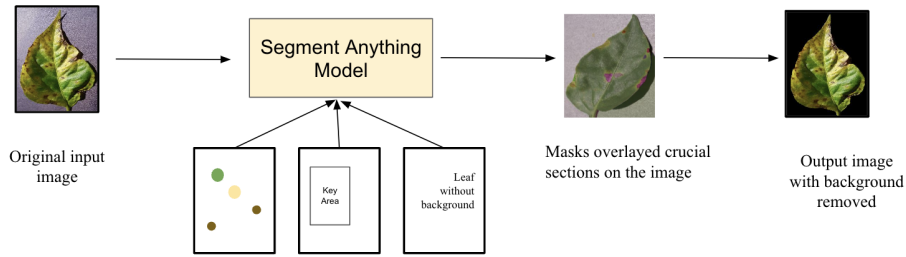


Figure 1: Workflow of the Segment Anything Model applied to a plant leaf image. SAM generates segmentation masks for isolating key regions of interest, such as disease spots, and removes the background. Outputs include segmented masks, an overlay on the original image, and a background-removed leaf for further analysis.

## 3.2 Vision Transformer and ResNet50

For this study, we utilized two prominent models: Vision Transformer (ViT) and ResNet50. Vision Transformer is a transformer-based architecture designed for image recognition tasks, leveraging self-attention mechanisms to model long-range dependencies within images. Unlike convolutional neural networks (CNNs), ViT divides an image into fixed-size patches and processes them similarly to tokens in natural language processing tasks. Plant leaf disease features (like spots, discolorations, and lesions) can appear in specific regions of a leaf rather than the whole image. This global attention allows the model to identify distributed disease features, no matter where they appear on the leaf. In addition, self-attention mechanisms in ViT allow it to focus on salient features, even if they are small and not in the center of the image. This enables it to recognize fine-grained details, such as

the appearance of early-stage disease symptoms. With its self-attention mechanism, ViT can also generalize better to unseen diseases, especially when trained with large datasets. It captures more abstract, high-level visual concepts, which can be helpful for classifying rare or new diseases. On the other hand, pretrained ViT model on large datasets like ImageNet can be fine-tuned on smaller, domain-specific datasets. Using transfer learning enhances classification performance, making it particularly suitable for our plant leaf disease dataset, which contains a limited amount of training data (approximately 20k samples).
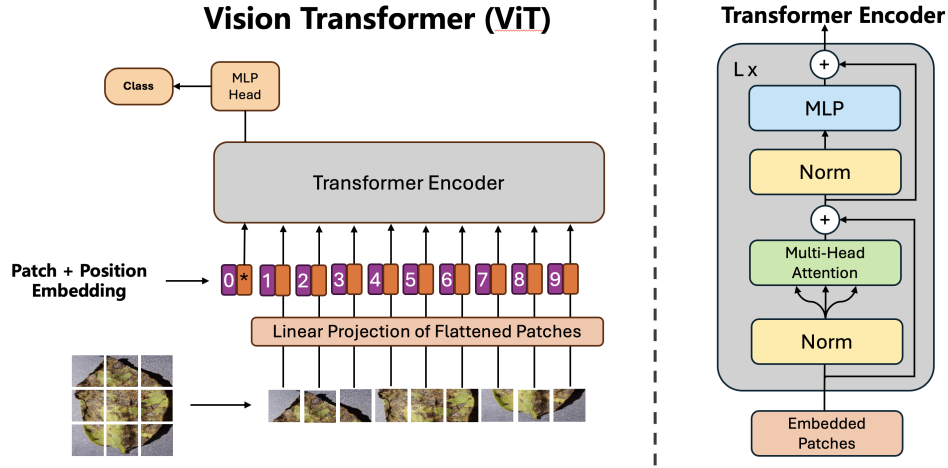


Figure 2: The diagram illustrates the architecture of the Vision Transformer. On the left, an example plant leaf input image is divided into non-overlapping patches, which are linearly projected and combined with positional embeddings. A special classification token (CLS) is added to the sequence of embeddings, which passes through the transformer encoder. On the right, the transformer encoder consists of multiple layers of multi-head attention, normalization, and feedforward networks (MLP). The final output is processed through an MLP head to produce the class prediction. This design allows ViT to capture global contextual information effectively, using its self-attention mechanisms for leaf disease recognition tasks.

ResNet50, on the other hand, is a well-established convolutional neural network that introduced the concept of residual learning to mitigate the vanishing gradient problem during training. Its ability to efficiently learn deep hierarchical representations makes it an excellent choice for plant leaf classification. ResNet50 uses skip connections, which allow gradients to flow directly through the network. This makes it possible to train deeper networks, unlike older models like VGG, which suffered from vanishing gradients. Since plant disease recognition often requires the extraction of subtle features from images, deeper networks can capture more abstract and higher-level features, but only if vanishing gradients are addressed. ResNet50, with 50 layers, extracts these subtle features without suffering from training instabilities. Similar as Vision transformer, ResNet50 is a pre-trained model on ImageNet, giving it strong generalization capabilities for recognizing features like edges, textures, and shapes. Since our datasets for plant disease classification are relatively small compared to ImageNet, transfer learning allows ResNet50 to "inherit" useful features like edge detection and color patterns. This is useful for recognizing disease symptoms like yellowing (chlorosis), spotting, and texture differences on leaves.

The choice of these two models was driven by their complementary strengths. ViT excels in capturing global contextual information, while ResNet50 is adept at extracting local features. By comparing the two models, we aim to evaluate which approach—global feature modeling (ViT) or local feature extraction (ResNet50)—is better suited for different aspects of plant disease classification. This comparison provides critical insights into the optimal strategy for disease detection. Another important factor driving this comparison is the need to assess computational efficiency. While ViT has higher computational costs due to its quadratic complexity in self-attention, ResNet50 is significantly more lightweight, offering faster inference speeds and lower GPU memory requirements. By comparing the two models, we can measure the trade-offs between inference time, computational cost, and accuracy. This is especially important for edge device that we will develop later on, where low-latency and efficient computation are essential. If ResNet50 achieves comparable accuracy with significantly

lower computational demands, it would be preferred for real-time disease diagnosis. However, if ViT demonstrates substantial improvements in classification accuracy for diseases that involve large-scale changes (like Tomato Yellow Leaf Curl Virus), then the increase in computational cost may be justified for certain use cases.
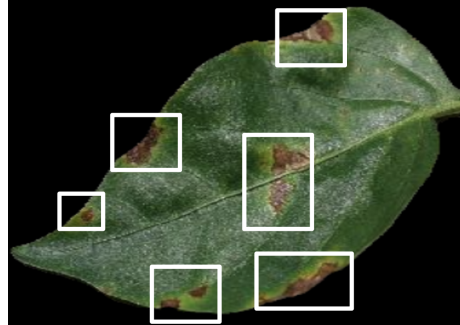


Figure 3: A diseased plant leaf with multiple affected regions highlighted. The Vision Transformer's advantage of its self-attention mechanisms enables effectively recognizing diseases that span across the entire leaf, capturing distributed patterns and subtle features that may not be centrally located. This global contextual understanding gives ViT an advantage over models like ResNet50, which rely on local feature extraction and may struggle to identify such widespread patterns.

### 3.3 Use of Pre-Trained Weights

Both the Vision Transformer (ViT) and ResNet50 models were initialized with pre-trained weights from the ImageNet dataset. This transfer learning approach is particularly advantageous for our task where the SAM dataset is not large enough to train models effectively from scratch. By starting with pre-trained weights, the models could focus on fine-tuning for plant disease-specific patterns. To further optimize model performance, a multi-stage training process was adopted. Initially, the models were fine-tuned on the original 20k dataset. This step enabled the models to adjust their parameters to better recognize plant disease-specific patterns. After this initial fine-tuning, the models were saved for subsequent training. In the second phase, we introduced the SAM dataset. The images in the SAM dataset were subjected to a series of data augmentation techniques, including horizontal flipping, random rotations, and Gaussian blurring. The saved models, pre-trained and fine-tuned on the original 20k dataset, were then further trained using this augmented SAM dataset. This two-stage training approach resulted in a final fine-tuned model that benefited from the strengths of both the original dataset and the SAM-augmented dataset. By combining the feature-extracting capabilities of pre-trained weights with task-specific knowledge from fine-tuning, our approach maximized the model's performance.

### 3.4 Loss Function: Cross-Entropy

We employed the cross-entropy loss function to optimize the models during training as it is well-suited for multi-class classification tasks. By minimizing this loss, the models were incentivized to assign higher probabilities to the correct classes, improving overall predictive accuracy.

### 3.5 Data Augmentation

Data augmentation techniques to the SAM dataset are applied on-the-fly during each epoch and for every batch. Random rotations, horizontal and vertical flips, and Gaussian blur were used to artificially increase the diversity of the dataset and reduce overfitting. Rotations and flips ensured the model's invariance to spatial orientation while Gaussian blur was applied to mimic subtle real-world imperfections, such as motion blur or focus issues in images, enabling the model to generalize better to unseen data. Each image can have different transformations applied during every epoch. For example, one epoch might see an image rotated 20 degrees, and in the next epoch, it might be rotated 45 degrees or flipped horizontally. This enables the model to see a larger diversity of augmented versions of the images. In addition, training-time data augmentation "hides" or "distorts" parts of the image for the model that prevent the overfitting issue. It acts as a dropout regularizer where small changes in the inputs prevent the model from focusing on specific, overfitted features. However, it

slightly slows down the training process since transformations happen at runtime (for each mini-batch) that introduce additional computational overhead.

### 3.6 Optimization and Training Strategies

The models were trained using the Adam optimizer, which combines the benefits of adaptive gradient methods and momentum to ensure robust optimization. Adam's ability to dynamically adjust the learning rates of each parameter made it particularly effective for fine-tuning pretrained models. In addition, a linear learning rate scheduler was employed to gradually reduce the learning rate during training, ensuring stable convergence and preventing oscillations in the loss function.

Weight decay was applied to regularize the models and mitigate overfitting by discouraging excessively large weight updates. The training process included a forward pass, where input images were propagated through the model to compute predictions, and a backward pass, where gradients were calculated and weights were updated to minimize the loss. This iterative process ensured that the models continuously improved their ability to classify plant leaf diseases accurately.

The methodology combines state-of-the-art deep learning architectures, strategic customizations, robust optimization techniques, and data augmentation to achieve highly accurate and generalizable models for plant leaf disease classification.

## 4. Experiment

### 4.1. Pre-training dataset

This study utilized a carefully curated dataset of plant leaf images categorized based on plant type, health status, and specific diseases. The original dataset contained approximately 20,000 images, organized into 15 subfolders, each representing a unique class. These classes encompass a diverse range of plant species and diseases, and each image in the dataset was labeled according to its corresponding plant type and health condition, facilitating supervised learning.

To enhance the dataset further, we employed the Segmentation Anything Model (SAM) to isolate the background and focus on the critical regions of interest within each image. This process generated a new SAM-segmented dataset, effectively creating an additional 20,000 images that were used alongside the original dataset for training and evaluation. By isolating the primary objects in each image, SAM helped standardize the dataset and improve model performance by removing irrelevant background features.



Figure 4: The picture demonstrates data processing and augmentation workflow. Each image in the original dataset undergoes several transformations to enhance the training dataset. First, the SAM technology removes the background, isolating the leaf. Next, data augmentation techniques, including random rotations, flips, and Gaussian blur, are applied to introduce variability and improve the model's generalization and robustness.

For the SAM-segmented dataset, we carefully partitioned the images into non-overlapping training and test sets, ensuring consistency with the split used in the original dataset. The partitioning process began by identifying and extracting images from the SAM-segmented dataset that corresponded to the identifiers in the original dataset's test set. The remaining images in the SAM-segmented dataset were assigned to the SAM training set. This method guaranteed no overlap between the training and test subsets, allowing for a seamless integration into the experimental pipeline and providing a robust foundation for training and evaluation.

### 4.2. Experimental Setup

We utilized pre-trained versions of Vision Transformer (ViT) and ResNet50, both of which had been pre-trained on the ImageNet dataset, to use the rich feature representations these models have already learned from a vast and diverse collection of images. The ImageNet dataset, containing millions of labeled images across thousands of categories, allows models to acquire generalizable patterns such as edges, textures, and shapes, which form the foundation for visual understanding. By initializing our models with these pre-trained weights, we significantly accelerated convergence during training and improved performance. The transfer learning is especially beneficial when working with our smaller datasets, where training deep models from scratch would likely lead to overfitting.

We implemented two round of training. The initial round of fine-tuning involved training the pre-trained ViT and ResNet50 models on the original 20k plant leaf dataset. This dataset consists of raw images categorized based on plant species and disease types. We split the dataset into 80% training and 20% testing to ensure a balanced evaluation while preventing data leakage. The goal of this round was to adapt the models to the specific domain of plant leaf disease classification while maintaining their general feature extraction capabilities. In the second round, we utilized the SAM-processed dataset, augmented with techniques such as random rotations, flips, and Gaussian blur. These augmentations were applied to enhance generalization and improve robustness to variations in real-world conditions. During this round, we also applied the preprocessing transformations. We standardized the pixel values to match the scale of the ImageNet pre-trained models. And then all images were resized to pixels to ensure compatibility with the input size required by both ViT and ResNet50. The training set and test set remained in the same 80%/20% split as the first round. This two-stage training process ensured that the models were first exposed to the raw dataset and later refined with the SAM-processed and augmented dataset for enhanced performance.

The experiments were conducted using a GPU A100 provided by Google Colab. The models were implemented using TensorFlow and PyTorch frameworks, along with the HuggingFace library for using pre-trained models. To ensure reproducibility, we set the random seed to 350 across all experiments. We adopted the following training parameters for the experiments: Vision Transformer was trained for 4 epochs in the first round and 5 epochs in the second round, while ResNet50 was trained for 7 epochs in both rounds. A batch size of 32 and 128 was used in original and SAMed dataset respectively to balance memory usage and training efficiency. The AdamW optimizer was employed with an initial learning rate of 5e-5, chosen for its ability to handle sparse gradients and its suitability for fine-tuning pre-trained models. A linear scheduler was used to gradually reduce the learning rate during training, ensuring stable convergence. Weight decay regularization was applied to mitigate overfitting by penalizing large weight updates.

To further evaluate the generalization capacity of the Vision Transformer and ResNet50 models, we incorporated a zero-shot evaluation. Zero-shot evaluation involves directly testing the pre-trained models on the plant leaf classification task without any domain-specific fine-tuning. This approach provides insight into the models' inherent ability to recognize unseen plant diseases based on prior knowledge from pre-training on the ImageNet dataset. By including zero-shot evaluation, we aim to assess the transferability and adaptability of the pre-trained models to a new domain. This evaluation serves as a benchmark, highlighting the models' capacity for generalization and providing a performance baseline before the fine-tuning process.

The performance of the models was evaluated using the following metrics: accuracy to measure overall correctness, F1-Score (Macro) to balance precision and recall across all classes regardless of class imbalance, precision to evaluate the proportion of true positive predictions among all positive predictions, recall to assess the proportion of true positives correctly identified, ROC-AUC (One-vs-Rest) to measure the area under the ROC curve for each class, and robustness test accuracy to evaluate the models' ability to maintain performance under added noise or distortions. These metrics provided a comprehensive evaluation framework to assess both classification accuracy and model robustness.

Lastly, the experiments were conducted using TensorFlow and PyTorch for model implementation and training. The HuggingFace library was employed to access the pre-trained ViT and ResNet50 models, streamlining the integration of these architectures into our pipeline. The combination of these tools ensured flexibility, efficiency, and reproducibility in our experiments.

# 5. Results

The results presented in this section highlight the performance of our approach for plant disease classification using the Segmentation Anything Model (SAM) for segmentation and two visual language models: Vision Transformer (ViT) and ResNet50. The experiments were conducted on the original and SAM-processed datasets, and the outcomes are evaluated using metrics such as accuracy, F1-Score, and robustness tests. Key findings are summarized using tables, plots, and visualizations.

*5.1 Main Results*

The initial evaluation of the models using a zero-shot approach demonstrated their limitations in this specific domain. The Vision Transformer (ViT), pre-trained on ImageNet, achieved a test accuracy of 0.73%, while ResNet50, also pre-trained on ImageNet, performed slightly better with a test accuracy of 1.12%. These results are consistent with expectations since neither model has been pre-trained for plant disease classification. ViT's performance was constrained by its reliance on high-level feature representations, which did not translate well to the unique visual patterns in plant disease images without fine-tuning. In contrast, ResNet50 showed slightly better zero-shot performance due to its ability to extract low-level features like edges and textures, which occasionally aligned with certain disease-related patterns. However, both models struggled to generalize effectively in this domain without additional training, underscoring the importance of fine-tuning on a labeled dataset tailored to plant disease classification. This evaluation highlights the necessity of adapting pre-trained models to domain-specific tasks to achieve meaningful results.

*5.1.1 Vision Transformer (ViT)*

**Training Metrics on the Original Dataset:**

| Epoch | Training Accuracy | Loss | Speed (it/s) | Notes |
|-------|-------------------|------|--------------|-------|
| 1/4 | 93.71% | 0.5581 | 3.28 | Initial convergence. |
| 2/4 | 99.47% | 0.0826 | 3.31 | Rapid improvement. |
| 3/4 | 99.88% | 0.0356 | 3.33 | Close to convergence. |
| 4/4 | 99.99% | 0.0214 | 3.32 | Achieved near-perfect. |

Table 1: ViT Training Metrics on the Original Dataset

During the training process on the original dataset, Vision Transformer (ViT) demonstrated exceptional capability to learn the disease patterns, as evidenced by rapid convergence in accuracy and steady reduction in the loss values over four epochs. Starting with an initial training accuracy of 93.71% and a loss of 0.5581 during the first epoch, the model quickly improved, achieving a final training accuracy of 99.99% and a loss of 0.0214 by the fourth epoch. The consistent improvements in accuracy and loss indicate effective learning with minimal signs of overfitting. The training speeds, averaging over 3.3 iterations per second, ensured efficient use of computational resources. These results affirm that ViT is well-suited for handling datasets with complex and diverse image features, allowing it to extract meaningful representations for classification tasks.

**Performance on Original Dataset:**

| Metric | Training Accuracy (%) | Test Accuracy (%) | Inference Time (ms/image) | Robustness Accuracy (10% noise) |
|--------|----------------------|-------------------|---------------------------|--------------------------------|
| Epoch 4/4 | 99.99 | 99.95 | 0.16 | 70.01 |

Table 2: Performance Metrics of Vision Transformer (ViT) on the Original Dataset

Fine-tuning Vision Transformer on the original dataset yielded outstanding results. The final training accuracy of 99.99% and test accuracy of 99.95% showcase the model's ability to generalize well to unseen data. This demonstrates the effectiveness of the model in capturing the fine-grained features of plant diseases, enabling it to differentiate between multiple classes with high precision. Despite this remarkable performance, the robustness test accuracy under 10% noise conditions was measured at 70.01%, indicating that while the model can handle moderately noisy inputs, its performance declines significantly in the presence of higher perturbations. This result suggests that although the model excels in clean data scenarios, additional training with augmented noisy data or adversarial methods could enhance its resilience to challenging conditions. Overall, these results position ViT as a powerful tool for plant disease classification on high-quality datasets.

**Performance on SAM-Processed Dataset:**

Using SAM-processed data introduced significant improvements in zero-shot performance while also highlighting challenges during fine-tuning. Initially, the Vision Transformer achieved a 24.9% zero-shot test accuracy on the SAM-processed dataset, a stark improvement compared to its 0.73% zero-shot accuracy on the original dataset. This improvement can be attributed to SAM's ability to isolate disease-relevant regions, effectively reducing background noise and enhancing the focus on disease features.

| Metric | Training Accuracy (%) | Test Accuracy (%) | Inference Time (ms/image) | Robustness Accuracy (10% noise) |
|---|---|---|---|---|
| Epoch 5/5 | 99.88 | 99.55 | 0.16 | 46.58 |

Table 3: Performance Metrics of Vision Transformer (ViT) on SAM-Processed Dataset

The results on the SAM-processed dataset illustrate both the strengths and limitations of integrating segmentation techniques into the classification pipeline. The substantial improvement in zero-shot accuracy to 24.9% highlights SAM's potential for enhancing pre-trained models by removing irrelevant background features. However, the slight drop in test accuracy after fine-tuning, from 99.95% on the original dataset to 99.55% on the SAM dataset, suggests that the segmentation process might discard certain contextual information necessary for optimal classification. Furthermore, the robustness test accuracy of 46.58% under noise conditions highlights the need to address segmentation-induced artifacts, which may amplify the model's sensitivity to minor perturbations. These findings underscore the need for iterative improvements in segmentation quality and robustness training to fully leverage the benefits of SAM-processed data.

*5.1.2 ResNet50*

**Performance on Original Dataset:**

| Metric | Training Accuracy (%) | Test Accuracy (%) | Inference Time (ms/image) | Robustness Accuracy (10% noise) |
|---|---|---|---|---|
| Epoch 7/7 | 99.95 | 99.81 | 0.23 | 49.10 |

Table 4: Performance Metrics of ResNet50 on the Original Dataset

The training and testing results for ResNet50 on the original dataset highlight its robust performance in identifying disease patterns. The model achieved a final training accuracy of 99.95% and a test accuracy of 99.81%, demonstrating its capacity to generalize effectively to unseen data. However, compared to ViT, the robustness test accuracy under 10% noise conditions was lower, at 49.10%. This suggests that while ResNet50 is effective in disease classification tasks, its performance might degrade more significantly in noisy or challenging conditions. Nevertheless, the model provides a strong baseline for understanding classification performance on the original dataset.

**Performance on SAM-Processed Dataset:**

| Metric | Training Accuracy (%) | Test Accuracy (%) | Inference Time (ms/image) | Robustness Accuracy (10% noise) |
|---|---|---|---|---|
| Epoch 8/8 | 99.68 | 99.40 | 0.17 | 23.40 |

Table 5: Performance Metrics of ResNet50 on SAM-Processed Dataset

ResNet50 exhibited a similar trend as Vision Transformer when trained on SAM-processed images, with reduced accuracy and robustness compared to the original dataset. The robustness test accuracy dropped significantly under noise conditions, highlighting the challenges posed by segmentation artifacts introduced by SAM.

*5.2 Confusion Matrix and Evaluation Metrics*

*5.2.1 Confusion Matrix*

The confusion matrix for ViT on the SAM-processed dataset highlights the class-wise performance:
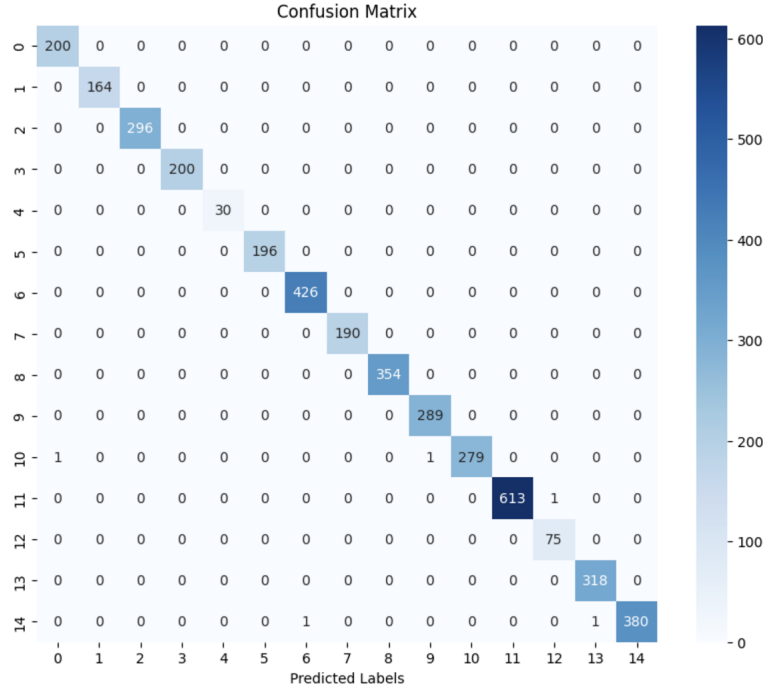
Figure 5: Confusion Matrix for ViT on the SAM-Processed Dataset

Diagonal values represent correct classifications, with off-diagonal values indicating misclassifications. Most misclassifications occurred between similar disease classes, such as Tomato Early Blight and Tomato Late Blight, due to overlapping visual symptoms.
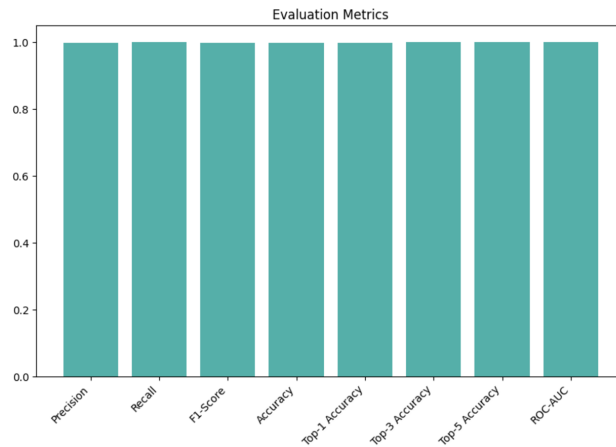
*5.2.2 Evaluation Metrics*



Figure 6: Evaluation Metrics for ViT (Precision, Recall, F1-Score, Accuracy)

All metrics exceed 0.99, showcasing the high performance of the Vision Transformer model across all evaluation criteria. These results emphasize the model's precision in minimizing false positives and its recall in identifying true positives. The high F1-Score demonstrates a balanced trade-off between precision and recall.
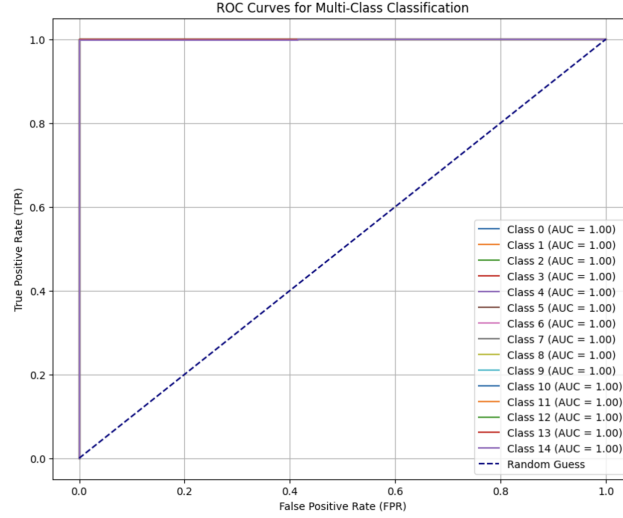
*5.2.3 ROC Curves*

11

Figure 7: ROC Curves for ViT on All Classes

The ROC curves for ViT demonstrate exceptional class-wise discrimination, with AUC values exceeding 0.99 for all classes. This indicates strong model performance in differentiating between positive and negative cases across all disease categories.

### 5.3 Comparison of ViT and ResNet50

The comparison between Vision Transformer (ViT) and ResNet50 reveals key differences in their performance across the original and SAM-processed datasets. ViT consistently outperformed ResNet50 in terms of test accuracy and robustness, as summarized in the table below:

| Model | Test Accuracy (Original) | Test Accuracy (SAM) | Robustness Accuracy (Original) | Robustness Accuracy (SAM) |
|---|---|---|---|---|
| ViT | 99.95% | 99.55% | 70.01% | 46.58% |
| ResNet50 | 99.81% | 99.40% | 49.10% | 23.40% |

Table 6: Comparison of ViT and ResNet50 Across Original and SAM-Processed Datasets

**Key Observations:**

- **Accuracy:** ViT consistently achieved higher test accuracy compared to ResNet50 across both datasets.
- **Robustness:** ViT demonstrated superior resilience to noise, particularly on the original dataset.
- **Inference Time:** ResNet50 exhibited slightly faster inference times, making it suitable for real-time applications.

### 5.4 Summary

Our results indicate that ViT achieves superior accuracy and robustness compared to ResNet50 for plant disease classification. However, SAM-processed datasets introduce challenges related to segmentation artifacts, requiring further optimization. Future work will focus on enhancing robustness and addressing segmentation-induced errors to improve overall performance.

## 6. Discussion

This section discusses the challenges, limitations, and future steps of our approach to plant leaf disease classification using Vision Transformer (ViT) and ResNet50, augmented with the Segmentation Anything Model (SAM). Our fine-tuned models, trained on both the original 20k dataset and the augmented SAM dataset, achieved impressive evaluation metrics, including high accuracy, precision, recall, and F1-scores. While the results indicate strong performance for both fine-tuned models, these results also highlight several challenges that warrant further consideration.

### 6.1 Navigating Overfitting and Dataset Challenges: The Need for Broader Testing and Balanced Data

Both the Vision Transformer and ResNet50 models demonstrated exceptional performance, achieving high accuracy, precision, recall, F1-score, and ROC-AUC values. While these metrics reflect the models' ability to classify the test set effectively, they also raise concerns about overfitting. High scores on the test set may not fully reflect the performance of the model in more scenarios where images are subject to greater variability in lighting, orientation, scale, and background noise. The data used in this study, while diverse, was still limited in comparison to more diverse datasets where conditions are often more unpredictable. As a result, the high accuracy and precision observed may not always translate to field-based performance. This risk is particularly pertinent given the dataset's characteristics, as discussed below. Additionally, the reported high metrics might mask certain biases that exist in the dataset. For instance, if certain disease classes are over-represented or if specific image characteristics dominate the dataset (certain leaf shapes or textures), the model might exploit these artifacts to achieve high classification accuracy. This reinforces the need to test the model on broader datasets that reflect the true distribution of plant images found in natural environments.

One of the key challenges in this study was the limited dataset size and the class imbalance within the dataset. The dataset is relatively small, containing approximately 20,000 images, and the data distribution across classes is highly imbalanced. For example, some disease classes, such as "Tomato Yellow Leaf Curl Virus," have a significantly larger number of samples compared to underrepresented classes, like "Potato Healthy." This imbalance can lead to biased predictions, where models prioritize majority classes while underperforming on minority ones. Additionally, the limited size of the dataset constrains the models' ability to learn diverse representations, making it difficult to assess their scalability and performance on larger, more varied datasets.

If applied to a larger, more diverse dataset, the models may face several challenges. A larger dataset typically includes greater variability in lighting, angles, backgrounds, and disease stages. While this variability could improve generalization, it also requires more computational resources and extended training time. Moreover, the effectiveness of SAM-generated segmented images on such a dataset remains uncertain, as segmentation quality may vary depending on the dataset's complexity.

### 6.2 Challenges in Segmentation: The Impact of Imperfect SAM Outputs on Model Performance

The use of the SAMed dataset introduces additional limitations due to its segmentation quality issues. Not all images in the SAM datasetdiscolorationed perfectly. A small but significant portion of images contained segmentation errors, such as incomplete leaf boundaries, missing parts of the leaf, or mislabeled background areas. Poorly segmented images introduce noise into the training process that can mislead the model, causing it to focus on irrelevant features. For instance, if part of a plant leaf is excluded from the segmented region, the model might learn to recognize irrelevant areas of the image instead of focusing on the disease spots. This reduces the effectiveness of training and impacts model generalization. In practice, this limitation suggests that SAM alone may not always provide reliable segmentation for complex or noisy images. Moreover, SAM's reliance on automdiscolorationeration may not always align with the specific plant leaf disease area. Subtle disease symptoms, such as discoloration or minor lesions, were excluded from the segmentation masks, limiting the models' ability to learn these critical features.

### 6.3 Balancing Power and Efficiency: Limitations of Vision Transformer and ResNet50 in Plant Disease Classification

While Vision Transformer and ResNet50 offer complementary strengths, their use in our plant disease classification also comes with limitations. Vision Transformer requires significantly higher computational resources due to its quadratic complexity in self-attention mechanisms. This makes it less efficient for our future real-time applications or deployment on edge devices with limited computational power. Additionally, our limited size of datasets make ViT struggles with its capacity, as transformers require large-scale data to fully exploit their capabilities.

ResNet50, while more lightweight and computationally efficient, has limitations in capturing plant leaf global contextual information. Its reliance on convolutional operations makes it less effective at modeling long-range dependencies compared to transformers. This lead to challenges in identifying subtle or distributed disease features that span across the entire leaf. Both models, despite using pre-trained weights, still face difficulty in generalizing to diseases or plant types not represented in the training data.

### 6.4 ViT and ResNet50 Baseline Performance and Limitations

The initial results from our zero-shot evaluation provided a clear understanding of ViT and ResNet50's inherent limitations when applied to a specialized task like leaf disease classification. Both models achieved 0% zero-shot accuracy, highlighting their inability to classify plant leaf diseases without task-specific fine-tuning. This outcome underscores the gap between their general pretraining on ImageNet and the fine-grained classification required for plant disease detection. Despite their strong pre-trained foundations, neither model possesses inherent knowledge of domain-specific features like leaf discoloration, spots, or lesions.

Fine-tuning these models on our dataset significantly improved both of their accuracy to nearly 99%. Both ViT and ResNet50 achieved high accuracy and evaluation metrics, validating the effectiveness of leveraging their pre-trained architectures for task-specific adaptation. However, the results also revealed persistent challenges. While the models excel at identifying broad disease patterns, they struggle with subtle, category-specific distinctions, particularly for underrepresented classes in the dataset. These limitations emphasize the need for refined data augmentation, segmentation, and potentially larger, more diverse training datasets to enhance generalization and classification accuracy.

### 6.5 Advancing Towards a Comprehensive and Scalable Solution

Building on our progress thus far, our project will focus on further integrating more models with refined data and training techniques. As we look ahead to the next phase of our project, several key questions will shape our direction: Can the performance of ViT and ResNet50 be further improved by scaling the dataset to 100k or even 1M images? Will balancing the dataset by augmenting underrepresented classes or generating synthetic images using GANs reduce model bias and improve generalization? How effective will attention-based segmentation approaches be compared to the current SAM-based segmentation process? Addressing these questions will guide future iterations of our approach.

For segmentation, our future work will explore multiple paths. One focus will be on refining the segmentation model to achieve more precise boundary accuracy. We also aim to incorporate attention mechanisms within the model architecture, enabling the system to dynamically focus on critical regions of the image rather than relying on static segmentation. Moreover, introducing quality control mechanisms to filter out poorly segmented images could enhance overall model performance. An additional approach will be to collect more high-quality, manually segmented images free from segmentation errors to provide a more robust training set.

To improve the robustness and scalability of the system, we plan to significantly expand our dataset by including a diverse range of plant species and disease categories. By incorporating a substantially larger dataset for training and testing, we aim to cover a wide variety of agricultural scenarios, making the model applicable to global use cases. Future work should focus on testing the models in real-world settings, such as in-field plant disease detection, to assess their practical utility and limitations. Addressing the computational resource requirements will be critical at this stage.

Additionally, we aim to modify the ViT and ResNet50 architectures directly, exploring ways to combine their complementary strengths and incorporating other advanced models to create a unified architecture. By refining these models and integrating them, we expect to enhance classification accuracy and efficiency, further advancing plant disease recognition.

In the final phase of this research, we aim to develop a user-friendly web or mobile application that enables users to upload images of plant leaves along with optional descriptions for disease diagnosis. This application will integrate the proposed hybrid pipeline, leveraging SAM-based segmentation with ViT and ResNet models, to provide real-time feedback on plant health.

## References

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.

[2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural SImulation System.* New York: TELOS/Springer–Verlag.

[3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.

[4] Tan, M., & Le, Q. (2019). EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. *International Conference on Machine Learning (ICML)*.

[5] Dosovitskiy, A., et al. (2020). An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv preprint arXiv:2010.11929*.

[6] De Silva, N., & Brown, T. (2023). Multispectral plant disease detection using hybrid CNN-ViT models. *Sensors, 23*(20), 8531. https://doi.org/10.3390/s23208531

[7] Li, X., & Li, Y. (2022). ConvViT: A hybrid Vision Transformer for apple disease detection. *Computer Vision and Image Processing.* https://doi.org/10.1007/s11263-021-01558-8

[8] Sahu, P., & Minz, S. (2023). Adaptive segmentation with intelligent ResNet and LSTM-DNN for multi-disease classification of plant leaves. *Springer AI Applications, 45*(2), 428-445. https://doi.org/10.1007/s11220-023-00428-3

[9] Hemalatha, K., & Jayachandran, A. (2024). Multitask learning-based Vision Transformer for plant disease localization and classification. *Journal of Agricultural Informatics, 10*(1), 123-136. https://doi.org/10.1007/s44196-024-00597-3

[10] Tan, M., & Le, Q. (2019). EfficientNet: Rethinking model scaling for convolutional networks. *International Conference on Machine Learning (ICML).* https://doi.org/10.48550/arXiv.1905.11946

[11] Zhang, H., Zhang, Y., Liu, M., Wang, T., Chen, Z., Li, P., & Qiu, J. (2023) Segment Anything Model and its applications. IEEE Transactions on Pattern Analysis and Machine Intelligence. https://doi.org/10.1109/TPAMI.2023.SAM2

[12] Na, J., Kim, H., Park, S., & Lee, J. (2024) Auto-prompting SAM with LoRA for domain-specific medical imaging. Journal of Medical Image Analysis. https://doi.org/10.1016/j.media.2024.105812

[13] SAM-Med2D Consortium. (2023) Segment Anything Model for detailed segmentation in medical imaging: Applications and challenges. Medical Imaging Advances. [No DOI available].

[14] Alomar, B., Khalifa, H., Mansour, A., & Diaz, R. (2023) Random Local Rotation: Object-specific data augmentation for enhanced classification. Pattern Recognition Letters. https://doi.org/10.1016/j.patrec.2023.06.005