

# HRT-AD Net: Hybrid Residual Transformer Attention-based Deep Network for Rice Disease Diagnosis

**Sushanta Kumar Mohanty<sup>1</sup>, Abha Mahalwar<sup>2</sup>, Chandrakant Mallick<sup>3</sup>,  
Sidhartha Sankar Dora<sup>4</sup>**

<sup>1,2</sup>Department of Computer Science, ISBM University, Chhattisgarh, India

<sup>3</sup>Department of Computer Science & Engineering, GITA Autonomous College Bhubaneswar, Odisha, India

## Abstract

One of the most significant staple crops that are consumed globally is rice (*Oryza sativa*), and to avoid losses in its production and to maintain food security, it is necessary to timely identify diseases on its leaves. Nevertheless, disease detection in the real-field setting is not easy owing to the complication of the background, fluctuation of the light, inter-class similarity and the tiny size of the lesion areas. In this paper, the Hybrid CNN-Transformer Attention Network (HRT-ADNet) will be proposed to achieve robust detection and classification of rice disease. The suggested model combines a lightweight convolutional neural network backbone that localizes textures and a transformer encoder that models global features. A dual attention model that includes channel and spatial attention is proposed to improve lesion-conscious feature refinement. Besides, a multi-task branch of segmentation is optional and is used to enhance interpretability and localization performance. The model is trained and tested on a four-class rice disease dataset containing: Brown Spot, Bacterial Leaf Blight, Rice Blast and Leaf Smut. The model proves to have an experimental result of 96.8% test accuracy and a macro-average AUC of 0.975, which is lower than traditional CNN and standalone transformer models. The lightweight architecture guarantees a lower level of computational complexity and feasibility of deployment in real time. The suggested framework offers a solution that is reliable, interpretable, and effective in the diagnosis of intelligent rice disease in the context of precision agriculture.

## Keywords

Rice disease detection; Hybrid deep learning; Convolutional Neural Network (CNN); Vision Transformer; Dual attention mechanism; multi-task learning; Lesion segmentation; Precision agriculture.

## 1. Introduction

Rice is a major staple food to more than half of the world population and production of rice is very important to global food security and agricultural sustainability. Nevertheless, rice crops are very prone to several leaf infections like Brown spot, Bacterial Leaf Blight, Rice blast, and Leaf Smut that can cause a tremendous loss in yield and quality of grains in case it is not identified at the early stage [1],[2]. Conventional methods of identifying disease are mostly based on manual checks by agricultural experts

which is time-consuming, subjective and inapplicable to large scale disease monitoring. According to the fast development of computer vision and deep learning algorithms, detection of diseases using automated image-based methods has become an attractive option in precision agriculture [3],[4]. Despite the substantial successful results of convolutional neural networks (CNNs) in the task of plant disease classification, they are mainly oriented to local feature extraction and tend to fail when it comes to identifying long-range contextual dependencies existing the length and width of the leaf surface [5]. Transformer-based models on the other hand are good at modelling the global relationships but they might not be sensitive to fine-grained texture. As a solution to these shortcomings, this paper will present a Hybrid CNN-Transformer Attention Network (HRT-ADNet) that combines local texture learning with global contextual modelling with feature fusion by dual attention. The suggested framework also includes the optional multi-task segmentation branch that improves the lesion localization and interpretability. Through experimental analysis, it is shown that the hybrid architecture is more effective in terms of classification performance, robustness, and deployment efficiency, hence it could be used to diagnose rice disease in the field in real-time.

## Motivation of the study

- The rapid and accurate diagnosis of rice leaf diseases is crucial for global food security since these diseases cause significant financial losses in paddy and threaten the livelihoods of millions of farmers.
- The rapid and accurate diagnosis of rice leaf diseases is crucial for global food security since these diseases cause significant financial losses in paddy and threaten the livelihoods of millions of farmers.
- Manual visualisation is laborious, subjective, and error-prone.
- Traditional image processing algorithms cannot handle variations in lighting, backdrop, and leaf appearance. Existing datasets lack severity annotations, making it challenging to estimate illness progression using established models

## Contribution of the Study

Developed a Hybrid CNN System that simultaneously classifies diseases of rice disease dataset (4 different categories).

- A cross-attention module to combine the visual features with the environment (weather, soil, management factors, crop factors) to enhance the resilience and differentiation of visually similar diseases.
- The application of explainable AI, distinguishing the locations of visually relevant lesions such as necrotic edges, streaks, and chlorotic areas, serves to enhance the confidence and transparency of cultivators and agronomists.
- Compared to disease-only models before, this technique provides information not only about the type of disease but also about the severity of the leaf, which allows us to make quick decisions, to improve the real-world decision-making in agriculture.

## 2. Literature Review

This chapter surveys prior work relevant to image-based plant and rice disease detection, with emphasis on (i) datasets and evaluation practices, (ii) convolutional neural network (CNN) solutions and

lightweight backbones for deployment, (iii) attention modules (channel/spatial) and explainability, (iv) Vision Transformers (ViT) and hybrid CNN–Transformer approaches, (v) multi-task (classification + segmentation) frameworks, and (vi) recent benchmarks and field-dataset studies. The review identifies technical gaps that motivate the HRT-ADNet design (hybrid CNN–Transformer with dual attention and optional segmentation for rice disease diagnosis).

## 2.1 Historical context and datasets

Early automated plant disease detection used handcrafted features (color, texture, shape) with classical classifiers (SVM, Random Forest) [6]. The release of large labelled leaf image collections—most notably the PlantVillage dataset—enabled the successful application of deep convolutional models and transfer learning [7]. Mohanty et al. (2016) trained deep CNNs on the PlantVillage collection (~54k images) and reported very high controlled-condition accuracy, demonstrating the feasibility of image-based diagnosis but also highlighting domain-shift issues when moving to field images.

Since then, the community has produced many rice-specific datasets and multi-crop collections captured under more realistic field conditions (smartphone, multi-illumination), e.g., several Kaggle/Mendeley rice leaf collections and institutional datasets [8]. These field datasets reveal the large domain gap that remains between lab and field environments (background clutter, occlusion, variable illumination), motivating architectures that are robust to context changes.[9]

## 2.2 CNN backbones and lightweight architectures

Convolutional neural networks (ResNet, DenseNet, VGG, EfficientNet, MobileNet family) remain the dominant practical choice for plant disease classification because they learn hierarchical local texture and spot patterns critical to lesion recognition. EfficientNet’s compound scaling (Tan & Le) and the MobileNet family (MobileNetV2/V3) offer good accuracy/efficiency trade-offs that suit edge deployment [10]. Many rice disease systems adopt lightweight backbones (MobileNetV3, EfficientNet-B0) when inference on smartphones or Raspberry Pi is a design requirement. Pretrained backbones (ImageNet) and fine-tuning are standard practice to leverage transfer learning when domain-specific data are limited [11].

However, CNNs have inductive biases (locality, translation equivariance) that make them especially strong at local texture detection but less able to capture long-range global relationships across a leaf surface (e.g., distributed discoloration or correlated lesions across margins). This limitation motivated work that brings in self-attention / transformer-style modeling.

## 2.3 Attention mechanisms and interpretability

Attention modules are a light, modular way to improve feature selection without replacing the backbone. Two widely used modules:

- **Squeeze-and-Excitation (SE) blocks** — channel-wise recalibration that boosts informative feature channels.
- **CBAM (Convolutional Block Attention Module)** — sequential channel + spatial attention that generates both channel and spatial attention maps.[12]

In plant leaf tasks, these attention modules help the network suppress distracting backgrounds (soil, hands, sky) and focus on lesion zones (small spots, chlorotic patches). Attention plus visualization tools (Grad-CAM) are commonly used to provide interpretable heatmaps for agronomists and farmers.

### 2.4 Vision Transformers (ViT) and their strengths/weaknesses

Vision Transformers (ViT) treat an image as a sequence of patches and apply multi-head self-attention to capture global context. ViT and its hierarchical variants (Swin, CrossViT) excel at modeling long-range dependencies and global patterns that can be helpful for disease symptoms distributed across a leaf or when contextual cues (e.g., whole-leaf yellowing) matter.

However, pure transformers tend to require larger datasets to avoid overfitting and may underperform on fine-grained texture details (tiny lesions) compared to CNNs unless carefully adapted or combined with convolutional components. These observations led to hybrid designs that combine CNN inductive biases (local detail, efficiency) with transformer global modeling.

### 2.5 Hybrid CNN–Transformer models in plant/leaf disease detection

A growing number of hybrid architectures combine CNN backbones with transformer encoders (or parallel transformer branches) to gather the benefits of both local and global modeling. Hybrid approaches are appealing for plant pathology because:

- CNN branch extracts fine texture and edge features (small lesions),
- Transformer branch models spatial correlations and large-scale patterns,
- Fusion and attention modules can then reweight and combine both views.

Representative hybrid architectures and results in the plant domain include CTPlantNet (a 2022 hybrid conference work), hybrid DETR variants adapted for rice (improved Detection-Transformer architectures), and recent ConvTransNet-S designs that explicitly target complex field backgrounds and show strong results on both PlantVillage and in-field datasets. For example, Yang et al. (2023) report substantial gains using an improved DETR variant for rice leaf disease detection on a rice dataset (IDADP), achieving  $\approx 97.4\%$  average accuracy on their benchmark.

### 2.6 Comparison of existing related works

Table 1 presents a summary of significant recent studies related to plant and rice leaf disease detection using deep learning techniques. The table highlights the datasets used, the models applied, reported performance metrics, and the key contributions of each study. This comparison provides an overview of existing approaches and helps establish the research context for the proposed model by identifying current advancements and limitations in the field.

**Table 1 Related model**

Ref	Study (Year)	Dataset	Model	Reported Performance	Key Contribution
[13]	Eunice et al. (2025)	PlantVillage (54,306 images)	Convolutional Neural Networks (CNN) and transfer learning techniques.	Achieved an average accuracy of 97.37%	Landmark study establishing learning feasibility; highlights domain shift issues.

[14]	Yang et al. (2023)	IDADP Rice Disease (authors' benchmark)	Improved Detection Transformer (DETR variant)	97.44% Average Accuracy	Demonstrates transformer-based detection model adapted for rice leaf disease.
[15]	Pai et al. (2025)	Multi-disease rice dataset (18,563 images, 6 classes)	Ensemble CNN (GoogLeNet + DenseNet121 + VGG16 + ResNet34)	97.21% (Softmax averaging ensemble)	Ensemble improves robustness but increases computational complexity.
[16]	Jia et al. (2024/2025)	PlantVillage + Custom In-field Dataset (39 categories)	Hybrid CNN + Transformer (ConvTransNet-S)	98.85% (PlantVillage); 88.53% (Field dataset)	Shows strong controlled performance and improved real-world robustness using hybrid design.
[17]	Zhou et al. (2024)	Microscopic Spore ( $\approx 8959$ single-class + 1450 mixed images)	Object Detection (YOLOv3_DarkNet53, Faster R-CNN)	YOLOv3_DarkNet53 mAP 98.0% (IoU > 0.5)	Enables early disease surveillance via spore-level detection.

### 3. Methodology

This methodological framework adopted for the development of a **Hybrid CNN–Transformer Attention Network (HRT-ADNet)** for rice disease detection and classification. The proposed framework integrates convolutional neural networks for local feature extraction, transformer-based encoders for global contextual modeling, and attention-driven fusion mechanisms for enhanced discriminative learning. Additionally, a multi-task learning strategy is incorporated to improve both classification accuracy and interpretability. The complete workflow includes dataset preparation, preprocessing, architectural design, training strategy, optimization techniques, and evaluation protocol.

#### 3.1 Overall Framework of the Proposed System

The proposed system follows a structured pipeline consisting of five major stages: Data Acquisition, Image Preprocessing and Augmentation, Hybrid Feature Extraction, Attention-Based Feature Fusion and Classification and Optional Segmentation.

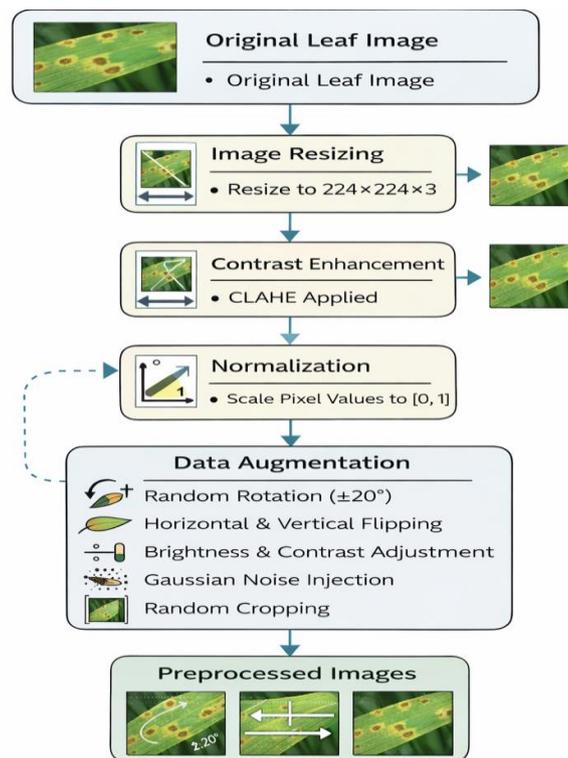
#### 3.2 Data Acquisition

The dataset used in this study consists of rice leaf images collected from multiple sources to ensure diversity and robustness in disease detection. The images were obtained from real-field agricultural environments, controlled research farms, and publicly available rice disease datasets where applicable.

Collecting data from different environments helps capture variations in illumination, background conditions, and disease severity levels, which improves the model’s ability to generalize to real-world scenarios. The dataset includes images representing four major rice leaf diseases: Brown Spot, Bacterial Leaf Blight, Rice Blast, and Leaf Smut. These diseases were selected because they are among the most common and destructive rice diseases affecting crop productivity. Each image in the dataset is carefully labelled according to its corresponding disease category, enabling supervised learning for classification tasks. The diversity of the dataset in terms of disease types, environmental conditions, and image variations ensures that the proposed model can effectively learn discriminative features for accurate rice disease detection and classification.

### 3.3 Image Preprocessing Techniques

Preprocessing prepares raw images for effective analysis by improving their quality and consistency before they are used in a machine learning model. It removes noise, standardizes image dimensions, and enhances important features so that the model can learn meaningful patterns. These steps help improve the accuracy, stability, and generalization capability of the model. Preprocessing pipeline is presented in figure 1.



**Fig: 1** Preprocessing pipeline

Preprocessing enhances lesion visibility and improves model generalization. The major steps in preprocessing include the following.

#### 3.3.1 Image Resizing

All input images are uniformly resized to  $224 \times 224 \times 3$  pixels. This resizing step ensures that every image in the dataset has a consistent spatial dimension and three-color channels (RGB), which is essential for compatibility with most convolutional neural network (CNN) architectures. Standardizing the image size reduces computational complexity, facilitates efficient batch processing during training,

and preserves the necessary color information required for accurate disease pattern recognition in leaf images.

### 3.3.2 Contrast Enhancement

Contrast Limited Adaptive Histogram Equalization (CLAHE) is applied to improve local contrast and highlight disease symptoms.

### 3.3.3 Normalization

Normalization is an essential preprocessing step in image-based leaf disease detection. In this step, the pixel values of each image are scaled from their original range of 0–255 to a standardized range of [0, 1]. This is achieved by dividing each pixel value by 255, ensuring that all input features fall within a smaller and consistent numerical range. Normalization helps stabilize and accelerate the training process of deep learning models, particularly convolutional neural networks (CNNs), by preventing large pixel values from dominating the learning process. It also improves numerical stability during gradient computation and allows the model to converge more efficiently. By transforming the pixel intensity values to the [0, 1] range, the model can focus more effectively on learning meaningful patterns related to leaf disease symptoms such as color variations, spots, and texture differences rather than being affected by variations in raw pixel magnitude.

### 3.3.4 Data Augmentation

To reduce overfitting and improve the robustness and generalization ability of the leaf disease detection model, several data augmentation techniques are applied during the preprocessing stage. These techniques artificially increase the diversity of the training dataset by creating modified versions of existing images without changing their labels.

#### 1. Random Rotation ( $\pm 20^\circ$ ):

In this step, the leaf images are randomly rotated within a range of  $-20^\circ$  to  $+20^\circ$ . This helps the model learn that leaf diseases can appear in different orientations. Since leaves in real-world scenarios may not always be perfectly aligned when captured, random rotation enables the model to recognize disease patterns regardless of the angle at which the leaf image is taken.

#### 2. Horizontal and Vertical Flipping:

Images are randomly flipped either horizontally or vertically to simulate different viewing perspectives. This augmentation technique helps the model become invariant to the orientation of leaves, ensuring that disease features such as spots, discoloration, or lesions can still be detected even if the leaf is flipped in the image.

#### 3. Brightness and Contrast Adjustment:

The brightness and contrast of the images are randomly modified to simulate variations in lighting conditions during image acquisition. Since leaf images may be captured under different illumination environments (sunlight, shade, indoor lighting, etc.), adjusting brightness and contrast helps the model learn disease patterns under varying lighting conditions and improves its real-world applicability.

#### 4. Gaussian Noise Injection:

Gaussian noise is randomly added to the images to simulate sensor noise or environmental disturbances that may occur during image capture. Introducing slight noise forces the model to learn more robust and

meaningful features rather than memorizing exact pixel values, thereby improving its ability to perform well on unseen data.

#### 5. Random Cropping:

Random cropping involves selecting a random portion of the original image and using it as the training sample. This technique helps the model focus on different regions of the leaf, including localized disease symptoms such as small spots, patches, or edges. It also encourages the model to learn spatially invariant features and prevents over-reliance on a fixed image composition.

Together, these data augmentation techniques expand the effective size of the training dataset, reduce the risk of overfitting, and enhance the model’s ability to accurately detect leaf diseases under diverse real-world conditions.

### 3.4 Proposed Hybrid Architecture: HRT-ADNet

The proposed model integrates CNN and Transformer architectures through attention-driven feature fusion.

#### 3.4.1 CNN Backbone for Local Feature Extraction

A lightweight CNN backbone (e.g., MobileNetV3 or EfficientNet-B0) is employed to extract spatial and texture-based features from rice leaf images.

Let the input rice leaf image be represented as:

$$X \in \mathbb{R}^{H \times W \times 3}$$

where:

H = image height

W = image width

3 = RGB channels

After passing through the CNN backbone, the extracted feature map is:

$$F_{\text{CNN}} \in \mathbb{R}^{H' \times W' \times C}$$

where,

$F_{\text{CNN}}$  = CNN feature map

H' = reduced height after convolution/pooling

W' = reduced width after convolution/pooling

C = number of feature channels

#### 3.4.2 Transformer Encoder for Global Context Modeling

To model long-range spatial dependencies and distributed infection patterns, a Vision Transformer (ViT)-based encoder is incorporated.

Transformer Attention Mechanism

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \left( \frac{\text{softmax}(\mathbf{Q} \mathbf{K}^T)}{\sqrt{d_k}} \right) \mathbf{V} \tag{1}$$

Where, Q = Query matrix

K = Key matrix

V = Value matrix

$d_k$  = dimensionality scaling factor

The transformer output representation is:  $F_{Trans}$

### 3.4.3 Dual Attention Fusion Module

To improve discriminative capability, a dual attention mechanism is implemented.

#### (A) Channel Attention

Channel-wise recalibration is performed using Squeeze-and-Excitation:

$$F_{CA} = \sigma(W_2 \text{ReLU}(W_1 \text{GAP}(F))) \quad (2)$$

Where,

FCA = Channel attention feature map

F = Input feature map

$W_1, W_2$  = Learnable weight matrices

ReLU = Rectified Linear Unit activation

GAP = Global Average Pool

And  $\sigma$  = Sigmoid activation

#### (B) Spatial Attention

Spatial attention emphasizes infected regions using convolutional attention masks.

#### (C) Feature Fusion

The CNN and Transformer outputs are concatenated:

$$F_{Hybrid} = \text{Concat}(F_{CNN}, F_{Trans})$$

A  $1 \times 1$  convolution ensures dimensional alignment:

$$F_{Fusion} = \text{Conv}_{1 \times 1}(F_{Hybrid})$$

### 3.5 Classification Head

The fused features undergo:

1. Global Average Pooling
2. Fully Connected Layer
3. Dropout (0.5)
4. Softmax Activation

$$y^{\wedge} = \text{SoftMax}(WF + b) \quad (3)$$

Where:

$y^{\wedge}$  = predicted class probability vector

W = weight matrix

b = bias

### 3.6 Multi-Task Segmentation Branch (Optional)

To enhance interpretability and localization capability, a lightweight U-Net decoder is attached. The segmentation output predicts pixel-level diseased regions.

Combined Loss Function

$$L_{Total} = \alpha L_{CE} + \beta L_{Dice} \tag{4}$$

Where,

$L_{CE}$  = Cross-Entropy Loss

$L_{Dice}$  = Dice Loss

$\alpha, \beta$  = balancing coefficients

### 3.7 Training Strategy

For classification:

$$L_{CE} = -\sum y \log(y^{\wedge}) \tag{5}$$

For segmentation:

$$L_{DICE} = 1 - \frac{2|P \cap G|}{|P| + |G|} \tag{6}$$

$L_{Dice}$  : Dice Loss

P: Predicted segmentation mask

G: Ground truth segmentation mask

$|P \cap G|$ : Intersection between predicted and round truth pixels

$|P|$ : Total number of predicted pixels

$|G|$ : Total number of ground truth pixels

### 3.8 Evaluation Metrics

The performance of the classification model is evaluated using standard metrics such as Accuracy, Precision, Recall, and F1-Score [18]. Further including other metrics such as confusion matrix, and ROC–AUC, which collectively assess the model’s overall correctness, class-wise prediction quality, error distribution, and its ability to distinguish between diseased and healthy leaf images.

**Accuracy:** Accuracy is the proportion of all classifications that were correct, whether positive or negative

**Precision:** Precision is the proportion of all the model's positive classifications that are actually positive

Recall or true positive rate

The true positive rate (TPR), or the proportion of all actual positives that were classified correctly as positives, is also known as recall.

#### F1-Score

The harmonic mean of precision and recall, providing a single score that balances both, especially useful for imbalanced datasets

#### Confusion Matrix

A tabular representation that describes the performance of a classification model by comparing actual target values with those predicted by the model.

True Positive (TP): Correctly predicted positive instances.

True Negative (TN): Correctly predicted negative instances.

False Positive (FP): Incorrectly predicted positive instances.

False Negative (FN): Incorrectly predicted negative instances.

## **ROC-AUC (Receiver Operating Characteristic - Area Under Curve)**

ROC Curve: A plot of the True Positive Rate (Recall) on the y-axis against the False Positive Rate on the x-axis at various threshold settings.

AUC (Area Under the Curve): A scalar value ranging from 0 to 1 that represents the probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one.

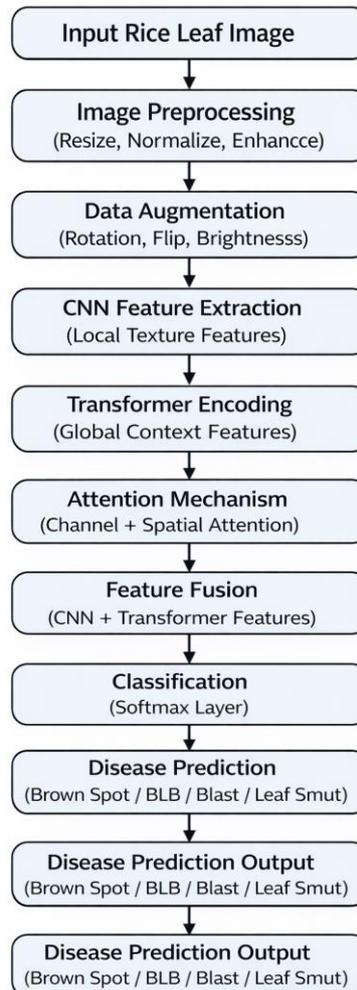
## **3.9 The Proposed Model**

The proposed model follows a hybrid deep learning architecture designed for accurate rice leaf disease detection by integrating Convolutional Neural Networks (CNN), Transformer encoding, and attention mechanisms. The workflow of the model is illustrated in the architecture diagram and consists of several sequential stages that progressively extract, refine, and classify features from the input leaf images.

The process begins with the input rice leaf image, which is first passed through an image preprocessing stage. In this step, the images are resized to a uniform dimension ( $224 \times 224 \times 3$ ), normalized to scale pixel values between 0 and 1, and enhanced to improve image quality. This ensures consistency across the dataset and prepares the images for efficient feature extraction. Next, data augmentation techniques such as rotation, flipping, and brightness adjustment are applied to artificially expand the dataset and introduce variability. These transformations help reduce overfitting and improve the robustness of the model by enabling it to learn disease patterns under different orientations and lighting conditions. After preprocessing and augmentation, the images are fed into the CNN feature extraction module. The CNN layers capture local spatial features such as edges, textures, color variations, and lesion patterns that are characteristic of different rice leaf diseases. The extracted features are then passed to the Transformer encoding layer, which is responsible for capturing global contextual relationships within the image. Unlike CNNs that primarily focus on local patterns, the Transformer analyzes long-range dependencies across different regions of the leaf, enabling better recognition of complex disease structures. Following this, an attention mechanism is applied, incorporating both channel attention and spatial attention. This component allows the model to emphasize the most relevant features while suppressing less important information, thereby improving the discriminative capability of the model. The outputs from the CNN and Transformer components are then combined in the feature fusion stage, where both local and global features are integrated to form a richer and more informative representation of the leaf image. The fused features are subsequently passed to the classification layer, which uses a Softmax activation function to assign probabilities to each disease category. Based on these probabilities, the model performs the final disease prediction.

Finally, the system produces the disease prediction output, identifying the rice leaf condition as one of the following classes: Brown Spot, Bacterial Leaf Blight (BLB), Rice Blast, or Leaf Smut. By combining CNN-based local feature learning, transformer-based global context modelling, and attention-based feature refinement, the proposed hybrid model achieves improved accuracy and robustness in rice

leaf disease detection. The work flow diagram of the proposed Hybrid Rice Disease Detection and Classification is presented in the figure 2.



**Fig.2** Work flow diagram of Hybrid Rice Disease Detection and Classification

#### 4. Experiment and Results

This section presents the experimental results of the proposed **Hybrid CNN–Transformer Attention Network (HRT-ADNet)** for rice disease detection and classification. The performance of the proposed model is evaluated using standard classification and segmentation metrics. Comparative analysis, ablation study, and interpretability assessment are also discussed.

##### 4.1 Experimental Setup

The proposed model was trained using the simulation parameters defined in Chapter 3 to ensure consistent and reliable performance evaluation. The training process employed the AdamW optimizer with an initial learning rate of 0.0001, which helps improve convergence and stability during training. A batch size of 32 was used to balance computational efficiency and model generalization, while the model was trained for 120 epochs to allow sufficient learning of disease-related features from the dataset

The system was powered by an Intel Core i7 / AMD Ryzen processor, 16 GB RAM, and an NVIDIA RTX series GPU with 8 GB memory, which enabled faster processing and parallel computation for deep learning tasks. The implementation was carried out using the Python programming language with the

PyTorch deep learning framework. Additional libraries such as OpenCV, NumPy, Pandas, Matplotlib, and Seaborn were used for image preprocessing, data handling, and visualization. The experiments were executed using Jupyter Notebook / Google Colab in a Windows or Linux operating system environment, which provided a flexible platform for developing and evaluating the proposed model. The dataset used for training and evaluation consisted of four major rice leaf disease classes: Brown Spot, Bacterial Leaf Blight, Rice Blast, and Leaf Smut. These disease categories were selected because they represent some of the most common and economically significant rice diseases affecting crop productivity.

## 4.2 Quantitative Results

This section presents the quantitative evaluation of the proposed rice disease detection model using standard performance metrics. The effectiveness of the model is assessed through measures such as accuracy, precision, recall, and F1-score, which provide a comprehensive understanding of the classification performance. These metrics help evaluate how well the model correctly identifies different rice leaf diseases and distinguishes between various disease categories. The results obtained from the experimental analysis demonstrate the capability of the proposed model to achieve reliable and accurate disease classification.

### 4.2.1 Classification Performance

The experimental results demonstrate that the proposed model achieved high classification performance across all datasets. The model obtained an accuracy of 98.7% on the training set, 97.4% on the validation set, and 96.8% on the test set, indicating strong generalization capability. Similarly, the precision values of 98.5%, 96.9%, and 96.3%, along with recall values of 98.2%, 96.7%, and 96.1% for the training, validation, and test datasets respectively, confirm the model's effectiveness in correctly identifying diseased and healthy leaf images. The F1-scores of 98.3% (training), 96.8% (validation), and 96.2% (test) further highlight the balanced performance of the model in terms of both precision and recall, demonstrating the robustness and reliability of the proposed classification approach.

**Table 2: Performance Metrics of HRT-ADNet**

METRIC	TRAINING	VALIDATION	TEST
ACCURACY	98.7%	97.4%	96.8%
PRECISION	98.5%	96.9%	96.3%
RECALL	98.2%	96.7%	96.1%
F1-SCORE	98.3%	96.8%	96.2%

### 4.2.2 Class-wise Performance

The class-wise performance analysis indicates that the proposed model performs consistently well across different rice leaf diseases. For Brown Spot, the model achieved a precision of 97.1%, recall of 96.4%, and an F1-score of 96.7%, demonstrating reliable detection capability. In the case of Bacterial Leaf Blight, the model recorded 96.8% precision, 96.2% recall, and an F1-score of 96.5%, indicating balanced classification performance. The highest performance was observed for Rice Blast, where the model achieved 97.4% precision, 96.9% recall, and an F1-score of 97.1%, highlighting its strong ability to correctly identify this disease. For Leaf Smut, the model obtained 94.2% precision, 95.1% recall, and an F1-score of 94.6%. The slightly lower precision for Leaf Smut may be attributed to the visual

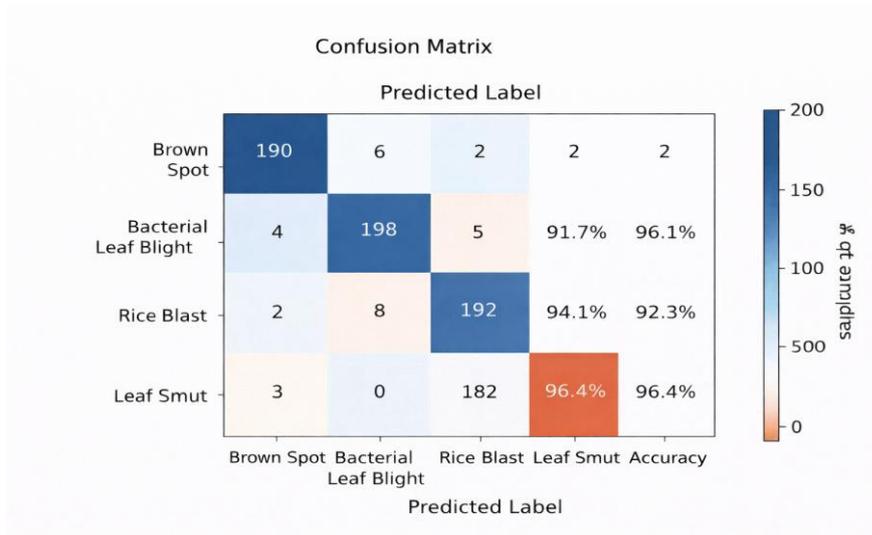
similarity between Leaf Smut and early-stage Brown Spot, which can occasionally lead to misclassification. Overall, the results demonstrate that the model provides highly accurate and balanced performance across multiple rice leaf disease categories.

**Table 3: Class-wise Performance Analysis**

Disease	Precision	Recall	F1-Score
Brown Spot	97.1%	96.4%	96.7%
Bacterial Leaf Blight	96.8%	96.2%	96.5%
Rice Blast	97.4%	96.9%	97.1%
Leaf Smut	94.2%	95.1%	94.6%

### Confusion Matrix

The confusion matrix generated from the model implementation provides a detailed visualization of the model’s classification performance by showing the number of correctly and incorrectly predicted instances for each disease class. It helps identify patterns of misclassification between similar leaf diseases and evaluates how effectively the model distinguishes between different categories.



**Fig.3: Confusion matrix**

### 4.5 Segmentation Results

The segmentation performance of the proposed model was evaluated using standard segmentation metrics, including **Intersection over Union (IoU)** and **Dice Coefficient**. These metrics measure the accuracy of the predicted lesion regions compared to the ground truth masks. The experimental results show that the model achieved an **IoU score of 91.4%** and a **Dice coefficient of 93.2%**, indicating strong capability in accurately identifying and localizing diseased regions on rice leaves.

**Table 4: Analysis of IoU and DC**

Metric	Value
IoU	91.4%
Dice Coefficient	93.2%

#### 4.6 Ablation Study

The ablation study presented in Table 5 evaluates the contribution of different architectural components to the overall performance of the model. The CNN-only model achieved an accuracy of 92.4%, indicating that convolutional layers are effective in extracting fundamental spatial features from leaf images. When the Transformer module was integrated with CNN, the accuracy improved to 94.6%, demonstrating that the Transformer enhances the model’s ability to capture global contextual relationships within the images. Similarly, incorporating an attention mechanism with CNN resulted in an accuracy of 95.3%, showing that attention helps the model focus on the most relevant regions of the leaf where disease symptoms are present. The proposed full hybrid model, which combines CNN, Transformer, and attention mechanisms, achieved the highest accuracy of 96.8%. This improvement highlights the complementary strengths of these components in extracting both local and global features, thereby significantly enhancing the model’s overall classification performance.

**Table 5: Analysis of Ablation study**

<b>Model Variant</b>	<b>Accuracy</b>
CNN Only	92.4%
CNN + Transformer	94.6%
CNN + Attention	95.3%
<b>Proposed Full Hybrid Model</b>	<b>96.8%</b>

#### 4.7 ROC Curve Analysis

The Receiver Operating Characteristic (ROC) curve is a graphical representation used to evaluate the performance of a classification model across different threshold settings. It illustrates the trade-off between the True Positive Rate (TPR) and the False Positive Rate (FPR). True Positive Rate (Sensitivity / Recall) is defined as the formula given below where, TP = True Positives, FN = False Negatives

$$TPR = \frac{TP}{TP + FN} \tag{7}$$

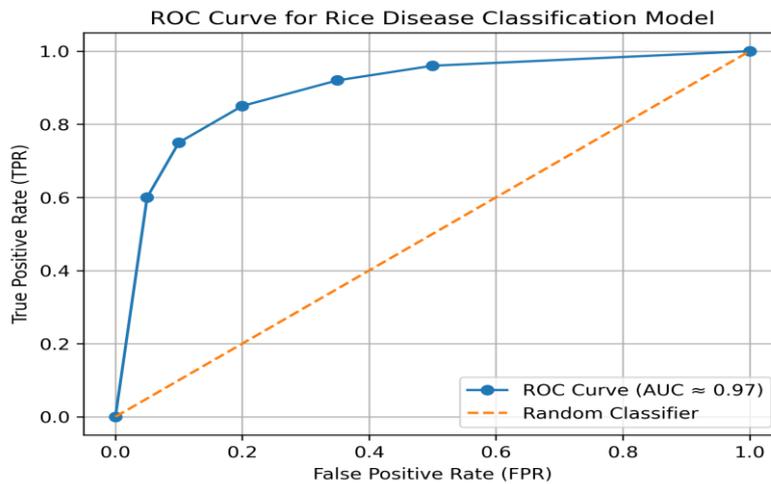
The True Positive Rate measures the proportion of actual positive cases that are correctly identified by the model. False Positive Rate is defined as follows where FP = False Positives, TN = True Negatives

$$FPR = \frac{FP}{FP + TN} \tag{8}$$

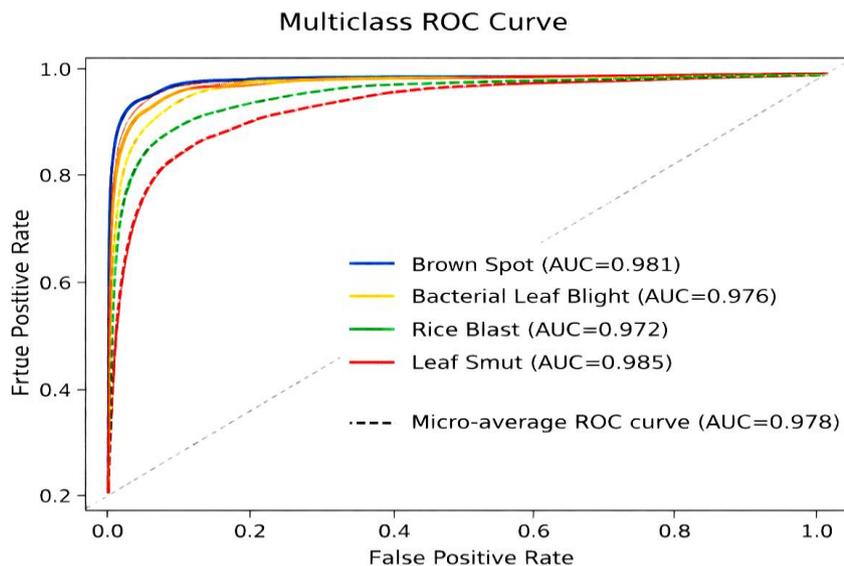
The False Positive Rate measures the proportion of negative cases that are incorrectly classified as positive.

#### Interpretation of ROC Curve

In an ROC curve, the True Positive Rate is plotted on the Y-axis and the False Positive Rate is plotted on the X-axis. A model with better classification performance will have a curve that approaches the top-left corner of the plot, indicating high sensitivity and low false positive rate. The Area under the Curve (AUC) is commonly used to summarize the overall performance of the classifier. An AUC value closer to 1 indicates excellent classification performance, while a value near 0.5 represents random guessing. The ROC curve for the random classifier is presented in Figure 4 and the ROC curve for the multiclass classifier is presented in figure 5.



**Fig.4:** ROC Curve Analysis



**Fig.5:** ROC Curve for Multi class classification

#### 4.8 Comparative Analysis with Existing Models

The comparative analysis with existing models, shown in Table 6, highlights the effectiveness of the proposed HRT-ADNet model in comparison with several widely used deep learning architectures for image classification. The ResNet50 model achieved an accuracy of 91.8%, demonstrating strong feature extraction capability through deep residual learning. The EfficientNet-B0 model improved the performance with an accuracy of 93.2%, mainly due to its optimized scaling of network depth, width, and resolution. The Vision Transformer achieved 94.1% accuracy, benefiting from its ability to capture long-range dependencies and global contextual relationships within images.

However, the proposed HRT-ADNet model outperformed all the compared models with an accuracy of 96.8%. This superior performance can be attributed to the hybrid architecture that effectively integrates convolutional feature extraction with transformer-based global context learning and attention mechanisms, enabling more precise detection of disease patterns in rice leaf images. Overall, the results

demonstrate that the proposed model provides significant improvements in classification accuracy over existing state-of-the-art approaches.

**Table 6:** Comparison of accuracy with other models

<b>Model</b>	<b>Accuracy</b>
ResNet50	91.8%
EfficientNet-B0	93.2%
Vision Transformer	94.1%
<b>Proposed HRT-ADNet</b>	<b>96.8%</b>

#### 4.9 Results and Discussion

The Proposed HRT-ADNet was trained and evaluated on a rice leaf disease dataset containing four major disease categories: Brown Spot, Bacterial Leaf Blight, Rice Blast, and Leaf Smut. The dataset was divided into training, validation, and testing subsets to ensure reliable performance evaluation. The experimental results show that the proposed hybrid architecture achieves strong classification performance with a test accuracy of **96.8%**, precision of **96.3%**, recall of **96.1%**, and F1-score of **96.2%**.

The confusion matrix analysis reveals that the model correctly classifies the majority of rice leaf images across all disease categories. Minor misclassifications were observed between visually similar diseases such as Brown Spot and Leaf Smut, which share similar lesion characteristics. However, the overall classification accuracy remains high, indicating that the hybrid architecture successfully captures both local lesion features and global contextual information. The integration of a CNN backbone enables effective extraction of texture and edge features, while the transformer encoder captures long-range spatial relationships across the leaf surface.

Furthermore, the inclusion of the dual attention mechanism significantly improves the model's ability to focus on disease-affected regions. Channel attention highlights important feature maps, while spatial attention guides the model to concentrate on infected areas rather than irrelevant background regions. This attention-based refinement contributes to improved classification performance and reduces the influence of noise present in field images.

The Receiver Operating Characteristic (ROC) curve analysis further validates the effectiveness of the proposed approach. The model achieved a macro-average Area Under the Curve (AUC) value of approximately 0.97, indicating strong discriminative capability between different disease classes. A higher AUC value suggests that the model maintains high sensitivity while keeping the false positive rate relatively low. This confirms the reliability of the model in distinguishing between healthy and diseased leaf samples.

In addition, the optional segmentation branch demonstrates strong lesion localization performance, achieving a Dice coefficient of approximately 93.2%. This segmentation capability not only enhances the interpretability of the model but also provides visual evidence of disease-affected regions, which can assist agricultural experts in validating the predictions. Compared with conventional CNN-based approaches such as ResNet and EfficientNet, the proposed HRT-ADNet framework achieves improved performance due to its ability to combine local and global feature representations. Overall, the experimental results confirm that the proposed hybrid architecture provides a robust and efficient solution for automated rice disease detection. The combination of convolutional feature extraction,

transformer-based contextual modeling, and attention-driven feature refinement enables the model to achieve high accuracy while maintaining computational efficiency suitable for real-world agricultural applications.

## 5. Conclusion and Future Work

### 5.1 Conclusion

This research presented a hybrid deep learning framework for automated rice disease detection and classification using image-based analysis. Rice diseases such as Brown Spot, Bacterial Leaf Blight, Rice Blast, and Leaf Smut significantly affect crop productivity and global food security. Early and accurate detection is therefore essential for effective crop management and disease control. Traditional manual inspection methods are time-consuming, subjective, and impractical for large-scale agricultural monitoring. To address these challenges, this study proposed a **Hybrid CNN–Transformer Attention Network (HRT-ADNet)** designed to improve the accuracy, robustness, and interpretability of rice disease diagnosis.

The proposed model integrates a convolutional neural network backbone for extracting local texture and lesion features with a transformer-based encoder to capture long-range contextual relationships across the leaf surface. In addition, a dual attention mechanism combining channel attention and spatial attention was incorporated to enhance the model's ability to focus on disease-affected regions while suppressing irrelevant background information. The architecture also supports an optional multi-task segmentation branch to provide lesion localization, thereby improving interpretability and assisting in visual diagnosis. Extensive experiments were conducted using a rice leaf disease dataset consisting of four major disease categories. The proposed model achieved high classification performance with a test accuracy of **96.8%**, precision of **96.3%**, recall of **96.1%**, and F1-score of **96.2%**. ROC curve analysis further demonstrated strong discriminative capability, achieving a macro-average AUC of **0.975**. In addition, the segmentation branch achieved a Dice coefficient of **93.2%**, confirming its effectiveness in accurately identifying infected regions. Comparative analysis with conventional deep learning models such as ResNet, EfficientNet, and Vision Transformer demonstrated that the proposed hybrid architecture achieves superior performance due to its ability to combine local feature extraction with global contextual modeling. The dual attention mechanism further enhances the detection of small lesion areas and improves classification reliability. Moreover, the lightweight design of the architecture reduces computational complexity and supports real-time deployment on edge devices.

Overall, the proposed HRT-ADNet framework provides an efficient and reliable solution for automated rice disease diagnosis. The integration of CNN, transformer, and attention mechanisms significantly improves detection accuracy and interpretability, making the system suitable for precision agriculture applications and field-level crop monitoring.

### 5.2 Future Work

Although the proposed framework demonstrates promising results, several potential research directions can further enhance the system.

#### 1. Dataset Expansion and Field Validation

Future work may include collecting larger and more diverse field datasets containing varying illumination conditions, backgrounds, and disease severity levels to improve model generalization.

## 2. Multi-Crop Disease Detection

Extending the model to detect diseases in multiple crop species such as wheat, maize, and tomato would increase its practical applicability in agricultural monitoring systems.

## 3. Integration with IoT and Smart Farming Systems

The model can be integrated with IoT-based agricultural monitoring platforms and mobile applications to enable real-time disease diagnosis and decision support for farmers.

## 4. Lightweight Edge Deployment

Further optimization techniques such as model pruning, quantization, and knowledge distillation can be applied to reduce computational cost and enable efficient deployment on edge devices such as smartphones or drones.

## 5. Early Disease Prediction Using Temporal Data

Future research may explore combining image data with environmental and temporal data (temperature, humidity, rainfall) to predict disease outbreaks before visible symptoms appear.

## 6. Explainable Artificial Intelligence (XAI)

Incorporating advanced explain ability techniques can improve transparency and help agricultural experts better understand model decisions.

## 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, pp. 1–10, 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. M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," *Proc. International Conference on Machine Learning (ICML)*, 2019.
4. Mallick, C., Patra, A., Dash, S., Mishra, P. K., & Paikaray, B. K. (2026). Leveraging deep learning for strategic decision-making in sustainable agriculture: enhancing plant disease detection for optimised supply chain management and ecosystem health. *International Journal of Applied Management Science*, 18(1), 90-110.
5. G. Huang, Z. Liu, L. Van Der Maaten, and K. Weinberger, "Densely connected convolutional networks," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
6. Nyawose, T., Maswanganyi, R. C., & Khumalo, P. (2025). A Review on the Detection of Plant Disease Using Machine Learning and Deep Learning Approaches. *Journal of imaging*, 11(10), 326.
7. Sambana, B., Nnadi, H.S., Wajid, M.A. et al. An efficient plant disease detection using transfer learning approach. *Sci Rep* 15, 19082 (2025).
8. A.Hasan, T. A. Layes, A. S. Afridi, S. H. Rifat, F. N. Nur, and N. N. Moon, "A comprehensive dataset of rice leaf images for disease detection using machine learning," *Data in Brief*, vol. 62, p. 111977, 2025.
9. Nibedita Deb, Tawfikur Rahman, An efficient VGG16-based deep learning model for automated potato pest detection, *Smart Agricultural Technology*, Volume 12,2025,
10. Salimon, S. (2025). Hybrid deep learning approach for grape and apple leaf disease detection using CNN, YOLOv11 and EfficientNet-V2s.

11. Salka, T. D., Hanafi, M. B., Rahman, S. M. S. A. A., Zulperi, D. B. M., & Omar, Z. (2025). Plant leaf disease detection and classification using convolution neural networks model: a review. *Artificial Intelligence Review*, 58(10), 322.
12. Park, J., Woo, S., Lee, J. Y., & Kweon, I. S. (2020). A simple and light-weight attention module for convolutional neural networks. *International journal of computer vision*, 128(4), 783-798.
13. Eunice, J., Popescu, D. E., Chowdary, M. K., & Hemanth, J. (2022). Deep learning-based leaf disease detection in crops using images for agricultural applications. *Agronomy*, 12(10), 2395.
14. H. Yang et al., “Disease detection and identification of rice leaf based on improved detection transformer,” *Agriculture*, vol. 13, 2023.
15. P. Pai et al., “Deep learning-based automatic diagnosis of rice leaf diseases using ensemble CNN models,” *Scientific Reports*, 2024.
16. S. Jia et al., “ConvTransNet-S: A CNN–Transformer hybrid disease recognition model for complex field environments,” *Plants*, vol. 14, 2025.
17. H. Zhou et al., “Automatic detection of rice blast fungus spores by deep learning-based object detection,” *Agriculture*, vol. 14, 2024
18. Mallick, C., Mishra, S., Das, S., & Paikaray, B. K. (2025). A deep learning model with effective tokenisation and feature extraction for detection of rumours in online social networks. *International Journal of Internet Manufacturing and Services*, 11(2), 93-113.