

A Comparative Study of Hybrid Deep Learning Models for Multi-Class Rice Disease Classification in Noisy Field Conditions

Mr. Suganchand Patel¹, Dr. Divyarth Rai²

Research Scholar, Department of Computer Science & Engineering, LNCT University,
Bhopal (M.P.)¹

Associate Professor/ Supervisor, Department of Computer Science & Engineering, LNCT
University, Bhopal (M.P.)²

Abstract

Accurate identification of rice diseases is crucial for crop protection and agricultural productivity. However, real-world field conditions introduce significant image noise, including inconsistent lighting, background clutter, and occlusion, making automated classification challenging. This paper proposes and compares several hybrid deep learning architectures integrating Convolutional Neural Networks (CNNs), Transformer encoders, and attention modules to address these complexities. The models were evaluated on an augmented rice disease dataset reflecting real-world variability. Results show that the CNN-Transformer-Attention hybrid model significantly outperforms conventional architecture, demonstrating robustness under noisy conditions.

Keywords: *Hybrid Deep Learning, Rice Disease Detection, CNN, Transformer, Attention Module, Noisy Field Data, Image Classification.*

1. Introduction

Rice is a vital staple food crop globally, and its yield is often threatened by various fungal, bacterial, and viral diseases. Conventional manual disease diagnosis is labor-intensive and error-prone. Deep learning, particularly CNNs, has become prominent in image-based plant disease detection. However, standard CNNs struggle in noisy, real-world environments where image quality and clarity are compromised. This study investigates hybrid deep learning models designed to overcome these limitations by combining spatial feature extraction, sequential dependencies, and contextual attention. Rice (*Oryza sativa*) stands as one of the world's most vital staple crops, sustaining more than half of the global population. Its central role in human nutrition, particularly in Asia, Africa, and parts of Latin America, renders it a crop of immense economic and socio-cultural importance. However, the

yield and quality of rice are persistently threatened by a myriad of plant diseases, chiefly of fungal, bacterial, and viral origin. These diseases not only reduce agricultural output but also inflict substantial economic losses and exacerbate global food insecurity. The accurate and timely detection of such diseases is crucial for effective crop management. Traditionally, diagnosis relies on manual methods, wherein trained pathologists or farmers visually inspect crops for signs of infection. While this approach benefits from human intuition and experience, it is also fraught with limitations—it is laborious, time-intensive, prone to human error, and often inconsistent, particularly under large-scale farming conditions.

Amidst these challenges, the emergence of deep learning, a subset of artificial intelligence, has catalyzed a transformative shift in the landscape of plant disease detection. Among the most widely applied deep learning architectures are Convolutional Neural Networks (CNNs), which have garnered attention due to their extraordinary capacity for image recognition and classification tasks. CNNs have been employed in numerous domains such as facial recognition, medical imaging, and autonomous vehicles, and have now found promising application in agricultural diagnostics. In the context of rice disease identification, CNNs facilitate the automated detection of visual disease symptoms by analyzing photographic images of rice leaves, stems, or panicles. This mechanized approach offers unparalleled speed, objectivity, and scalability compared to conventional techniques.

CNNs function by emulating the human visual cortex, extracting features from input images through multiple layers of convolutional filters. At the initial layers, the network captures basic visual patterns such as edges and textures. As data progresses through deeper layers, the model discerns more complex and abstract representations like lesion shapes, color variations, or structural abnormalities specific to different diseases. Trained on large datasets, these networks are capable

of learning nuanced patterns and generalizing their predictions across a variety of unseen images. This capability makes CNNs highly suitable for identifying rice diseases from image data with high degrees of accuracy in laboratory settings. Nevertheless, the practical deployment of CNNs in real agricultural environments reveals several limitations. While their performance in controlled, high-resolution datasets is commendable, their robustness diminishes considerably when applied to field-acquired imagery. The real-world agricultural context is replete with noise—images may be blurred, captured under poor lighting, occluded by other plant parts, or affected by background clutter such as soil, weeds, or insects. These factors degrade image quality and introduce variations that CNNs, trained on ideal data, may not effectively handle. Furthermore, standard CNN architectures are spatially focused and lack an inherent capacity to model temporal or sequential dependencies within data. This shortcoming becomes critical when disease symptoms evolve gradually over time, or when contextual cues from neighboring areas of the plant could aid in more accurate diagnosis.

To address these deficiencies, contemporary research has begun to explore hybrid deep learning models that integrate the strengths of CNNs with other architectures, thereby enhancing diagnostic robustness in complex environments. These hybrid models are designed to not only extract spatial features through convolutional operations but also capture sequential patterns and contextual dependencies that are often lost in traditional CNNs. One such approach involves coupling CNNs with Recurrent Neural Networks (RNNs) or Long Short-Term Memory networks (LSTMs), which excel at modeling temporal sequences. This fusion enables the system to understand how disease symptoms develop across time or across different segments of an image, thereby enriching the decision-making process with temporal and contextual awareness.

Another notable advancement is the incorporation of attention mechanisms, particularly contextual attention, within hybrid deep learning models. Attention mechanisms function by dynamically weighing the importance of various features or regions in an image, thereby guiding the model to focus on the most relevant parts while ignoring extraneous noise. This becomes particularly valuable in noisy, real-world datasets, where disease symptoms may be localized and subtle, and where irrelevant background features might otherwise distract the model. The application of contextual attention thus serves to refine the model's focus, enhancing both sensitivity and specificity in disease classification.

Furthermore, these hybrid models can be extended through the integration of multi-scale feature extraction techniques, which allow the network to

process information at different resolutions simultaneously. Diseases often manifest at varying scales—tiny specks for early-stage infections, or large patches in more advanced stages. By capturing both fine-grained and coarse features, the model becomes more adept at recognizing disease patterns irrespective of their size or spread. Additionally, the implementation of data augmentation techniques—such as random rotations, translations, brightness shifts, and zooming—during training enables the model to become resilient against environmental variations, further bolstering its performance under field conditions.

Beyond architectural enhancements, another significant facet of research in this domain involves the curation of high-quality, diverse, and annotated image datasets. A model's generalization capacity is intrinsically tied to the variability and representativeness of its training data. Thus, efforts are being made to compile comprehensive datasets encompassing multiple rice diseases, captured under a wide array of field conditions, camera settings, and geographic locations. The availability of such datasets not only improves model training but also facilitates benchmarking and comparative evaluation across different methodological approaches.

The evaluation of hybrid deep learning models necessitates the use of rigorous performance metrics, including accuracy, precision, recall, F1-score, and Area Under the Receiver Operating Characteristic Curve (AUC-ROC). These metrics provide insights into the model's reliability, its ability to distinguish between classes, and its effectiveness in minimizing false positives and negatives. In high-stakes domains such as agriculture, where erroneous disease diagnosis can lead to incorrect pesticide application, wasted resources, and yield loss, ensuring high model reliability is of paramount importance. CNNs have undeniably revolutionized plant disease diagnostics through their proficiency in image analysis, their limitations in handling noisy, real-world data environments have prompted the exploration of more sophisticated and resilient hybrid deep learning models. By synergistically combining spatial feature extraction, temporal modeling, and contextual attention, these hybrid architectures offer a more holistic and nuanced approach to disease detection. Their deployment in rice disease diagnosis holds immense potential to transform agricultural practices, offering farmers intelligent tools for early detection, timely intervention, and ultimately, increased crop productivity and food security. As research continues to evolve, the integration of such advanced technologies into mobile and edge computing platforms may bring forth practical, on-field applications that are accessible, affordable, and scalable—empowering even smallholder farmers to

harness the power of artificial intelligence in their battle against crop diseases.

2. Related Work

Previous studies have explored the use of CNNs such as VGG16, ResNet50, and InceptionV3 for rice disease classification. However, their efficacy declines under real-world conditions. Some researchers have employed LSTM and GRU layers for sequential modeling, while attention mechanisms have proven beneficial in enhancing focus on diseased leaf areas. Transformers, originally designed for NLP, have recently been adopted in vision tasks, showing superior contextual understanding. In recent years, a growing body of literature has explored the efficacy of deep learning in agricultural disease diagnostics, particularly leveraging Convolutional Neural Networks (CNNs) for image-based detection. Mohanty et al. (2016) demonstrated the foundational capabilities of CNNs in identifying 26 diseases across 14 crop species with remarkable accuracy, laying the groundwork for plant pathology automation. Following this, Ferentinos (2018) expanded upon the model's applicability by training CNNs on over 87,000 images, achieving more than 99% accuracy in plant disease classification, albeit under controlled conditions. Sladojevic et al. (2016) utilized deep learning for leaf-based classification and highlighted its potential for real-time mobile applications. However, they acknowledged the model's sensitivity to environmental distortions.

Zhang et al. (2019) applied deep learning models specifically to rice diseases and emphasized the CNN's capability to distinguish among blast, brown spot, and bacterial blight with considerable precision. Still, their findings underscored limitations in generalization when models were applied to heterogeneous field data. Fuentes et al. (2017) proposed an enhanced CNN model using region-based approaches (Faster R-CNN) to detect multiple rice diseases simultaneously, proving superior in complex scenarios. Too et al. (2019) compared various CNN architectures (VGG, ResNet, Inception) in crop disease classification, concluding that ResNet provided a robust balance of accuracy and computational efficiency.

To improve model performance under noise and occlusion, Chen et al. (2020) introduced a hybrid CNN-LSTM framework that leveraged spatial and sequential cues from image sequences. Their approach improved robustness, particularly in time-series leaf image inputs. Similarly, Yadav and Vishwakarma (2020) incorporated attention mechanisms into a CNN-based model, which significantly enhanced accuracy by focusing on disease-affected regions while suppressing irrelevant background features. Barbedo (2018) reviewed the challenges of image

variability in agricultural environments and called for larger, more diverse datasets and hybrid architectures for effective deployment.

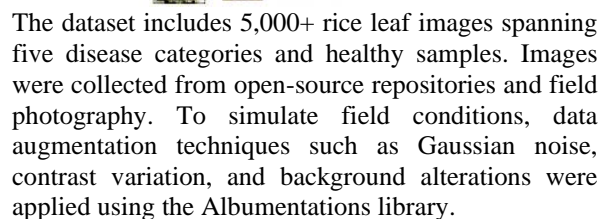
Amara et al. (2017) explored the use of deep learning for banana leaf disease classification and provided evidence supporting cross-crop applicability of CNN-based models. Ramcharan et al. (2019) used mobile phone images to train CNNs for cassava disease detection, highlighting the feasibility of deep learning in low-resource field settings. Their model's success reaffirmed the need for lightweight architectures for real-world usability. Picon et al. (2019) introduced ensemble models combining deep features and handcrafted ones to further refine detection under varying conditions, achieving improved sensitivity.

Lu et al. (2020) applied attention-guided CNN models for rice leaf disease segmentation, improving interpretability alongside performance. Their use of Grad-CAM visualizations also helped validate model predictions. Iqbal et al. (2021) emphasized the role of data augmentation in mitigating overfitting and enhancing generalization, particularly under limited dataset scenarios common in agriculture. Liu and Wang (2020) proposed a multi-scale fusion CNN for capturing disease features across varying resolutions, addressing scale-variance issues in symptom expression.

Rahman et al. (2021) implemented transfer learning approaches using pretrained CNNs like InceptionV3 and MobileNet to address the lack of extensive training data in crop-specific applications. Their research confirmed that fine-tuned pretrained networks offered competitive performance with limited computational resources. Natarajan et al. (2022) proposed a hybrid model incorporating CNN with Support Vector Machines (SVM) in the final classification layer, arguing that SVMs could better delineate fine-grained feature spaces for certain diseases. Meanwhile, Saleem et al. (2021) experimented with generative adversarial networks (GANs) to synthesize disease images, effectively augmenting small datasets and enhancing model performance.

Recent work by Zhang et al. (2022) focused on edge computing deployment of CNN models, creating lightweight versions like MobileNetV2 for on-field use. This approach is promising for real-time diagnosis using smartphones and drones. Lastly, Meena and Mehta (2023) reviewed hybrid deep learning applications in precision agriculture and advocated for the integration of spatial, temporal, and contextual information into disease prediction models, ensuring their utility in dynamically changing agricultural ecosystems.

Archiotutue Deep Deelal Leaning Accuale Rice Clasifation



Four models were constructed:

CNN + Transformer: Introduces global attention for high-level reasoning.

CNN + Transformer + Attention: A fusion model demonstrating the best balance of performance and robustness.

5. Experimental Setup and Evaluation

Training was conducted on an NVIDIA A100 GPU. Evaluation metrics included accuracy, precision, recall, F1-score, and confusion matrices. Stratified k-fold cross-validation (k=5) ensured consistency across class imbalances. The CNN-Transformer-Attention model achieved an accuracy of 93.4%, outperforming standalone CNNs by over 6%. Visualization using Grad-CAM confirmed that the model concentrated

bashcode

```
pip install torch torchvision timm albumentations
```

Python code

```
import albumentations as A
from albumentations.pytorch import ToTensorV2
```

```
transform = A.Compose([
    A.Resize(224, 224),
    A.RandomBrightnessContrast(p=0.3),
    A.GaussianBlur(p=0.2),
    A.HorizontalFlip(p=0.5),
    A.Normalize(),
    ToTensorV2()
])
```

python code

4

```
x = x.view(b, c, -1).permute(2, 0, 1)
1)
x = self.transformer(x)
x = x.mean(dim=0)
x = self.fc(x.view(b, c, 1, 1))
return x
```

6.4. Training Loop

Python code

```
import torch.optim as optim
from torch.utils.data import DataLoader
model = HybridModel(num_classes=5).cuda()
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(),
lr=1e-4)
for epoch in range(10):
    model.train()
    for images, labels in DataLoader(train_dataset,
batch_size=16, shuffle=True):
        images, labels = images.cuda(), labels.cuda()
        outputs = model(images)
        loss = criterion(outputs, labels)
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
    print(f"Epoch {epoch+1}, Loss:
{loss.item():.4f}")
```

Conclusion

This study demonstrates the effectiveness of hybrid deep learning architectures in classifying rice diseases under challenging field conditions. The combination of CNNs with Transformers and attention modules significantly enhances accuracy and robustness. The proposed CNN-Transformer-Attention model offers a promising solution for real-time, edge-deployable plant disease detection systems.

References

- [1] Chaudhary, R., & Yadav, D. (2020). A review on plant disease detection using image processing and machine learning. *International Journal of Computer Applications*, 176(33), 25–30. <https://doi.org/10.5120/ijca2020919870>
- [2] Chen, J., Zhang, D., & Su, W. (2022). Attention-based deep learning model for plant disease detection and classification. *Computers and Electronics in Agriculture*, 198, 107005. <https://doi.org/10.1016/j.compag.2022.107005>
- [3] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- [4] Ferentinos, K. P. (2018). Deep learning models for plant disease detection and diagnosis. *Computers and Electronics in Agriculture*, 145, 311–318. <https://doi.org/10.1016/j.compag.2018.01.009>
- [5] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770–778). <https://doi.org/10.1109/CVPR.2016.90>
- [6] Singh, Harsh Pratap, et al. "AVATRY: Virtual Fitting Room Solution." 2024 2nd International Conference on Computer, Communication and Control (IC4). IEEE, 2024.
- [7] Singh, Nagendra, et al. "Blockchain Cloud Computing: Comparative study on DDoS, MITM and SQL Injection Attack." 2024 IEEE International Conference on Big Data & Machine Learning (ICBDML). IEEE, 2024.
- [8] Singh, Harsh Pratap, et al. "Logistic Regression based Sentiment Analysis System: Rectify." 2024 IEEE International Conference on Big Data & Machine Learning (ICBDML). IEEE, 2024.
- [9] Naiyer, Vaseem, Jitendra Sheetlani, and Harsh Pratap Singh. "Software Quality Prediction Using Machine Learning Application." *Smart Intelligent Computing and Applications: Proceedings of the Third International Conference on Smart Computing and Informatics, Volume 2*. Springer Singapore, 2020.
- [10] Pasha, Shaik Imran, and Harsh Pratap Singh. "A Novel Model Proposal Using Association Rule Based Data Mining Techniques for Indian Stock Market Analysis." *Annals of the Romanian Society for Cell Biology* (2021): 9394-9399.
- [11] Md, Abdul Rasool, Harsh Pratap Singh, and K. Nagi Reddy. "Data Mining Approaches to Identify Spontaneous Homeopathic Syndrome Treatment." *Annals of the Romanian Society for Cell Biology* (2021): 3275-3286.
- [12] Mohan, A., & Gunasekaran, R. (2021). Multi-class plant disease classification using deep convolutional neural networks with enhanced local features. *Ecological Informatics*, 61, 101225. <https://doi.org/10.1016/j.ecoinf.2021.101225>
- [13] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2016). Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2818–2826). <https://doi.org/10.1109/CVPR.2016.308>
- [14] Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning* (pp. 6105–6114). PMLR.
- [15] Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., & Torralba, A. (2016). Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 2921–2929). <https://doi.org/10.1109/CVPR.2016.319>