



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** V **Month of publication:** May 2026

DOI: <https://doi.org/10.22214/ijraset.2026.79729>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Multimodal Maple Plant Disease Detection Using EfficientNet and Transformer-Based Semantic Fusion

Pramod Mishra¹, Meetendra Singh Chahar², Sapna Balani³

¹Student, M.Tech, CSE Department, Eshan college of engineering, Farah, Mathura, Uttar Pradesh

²Associate Professor, CSE Department, Eshan college of engineering, Farah, Mathura, Uttar Pradesh

³Assistant Professor, CSE Department, Eshan college of engineering, Farah, Mathura, Uttar Pradesh

Abstract: This paper presents a novel multimodal deep learning framework for maple plant disease detection by integrating visual and semantic information. Traditional plant disease detection systems rely primarily on visual features extracted from leaf images, which often leads to misclassification in cases of visually similar disease symptoms. To address this limitation, the proposed approach combines EfficientNet-based convolutional neural networks for visual feature extraction with transformer-based language models, including BERT and FLAN-T5, for semantic feature encoding. A Multilayer Perceptron (MLP)-based fusion mechanism is employed to integrate visual and textual features, enabling effective cross-modal learning. The proposed model is evaluated on a balanced dataset of 2,000 maple leaf images and associated disease descriptions. Experimental results demonstrate that the multimodal framework achieves an accuracy of 94.8%, outperforming vision-only and text-only models by a significant margin. Ablation studies and comparative analysis confirm the effectiveness of multimodal fusion and transformer-based semantic encoding. The proposed framework provides a robust and scalable solution for intelligent plant disease detection and has potential applications in smart agriculture systems.

Keywords: Plant Disease Detection, Multimodal Learning, EfficientNet, BERT, FLAN-T5.

I. INTRODUCTION

Agriculture plays a vital role in sustaining global food security and economic stability, particularly in countries where a significant portion of the population depends on farming and forestry. Among various factors affecting agricultural productivity, plant diseases represent one of the most critical challenges, leading to substantial yield loss and economic damage. Early and accurate detection of plant diseases is therefore essential to ensure timely intervention and effective crop management. In forestry applications, such as maple plant cultivation, disease detection becomes even more complex due to environmental variability and the subtle nature of disease symptoms[1].

Traditionally, plant disease identification relies on manual inspection by agricultural experts, who analyze visual symptoms such as discoloration, lesions, and texture variations on leaves. While this approach can be effective, it is inherently limited by subjectivity, scalability issues, and dependence on expert knowledge. In large-scale agricultural settings, manual diagnosis becomes impractical, time-consuming, and prone to human error. These limitations have motivated the development of automated disease detection systems using machine learning and computer vision techniques.[2]

In recent years, deep learning has emerged as a powerful tool for plant disease detection. Convolutional Neural Networks (CNNs), in particular, have demonstrated remarkable success in image-based classification tasks due to their ability to automatically learn hierarchical feature representations. Models such as VGGNet, ResNet, and EfficientNet have been widely applied to leaf-based disease detection, achieving high classification accuracy. EfficientNet, with its compound scaling strategy, offers an optimal balance between accuracy and computational efficiency, making it particularly suitable for real-world applications.[3]

Despite these advancements, vision-only approaches suffer from inherent limitations. Many plant diseases exhibit similar visual characteristics, especially during early stages, making it difficult for CNN-based models to distinguish between them. For example, diseases such as anthracnose and leaf scorch may present comparable discoloration patterns, leading to misclassification. Furthermore, variations in lighting conditions, background noise, and leaf orientation can adversely affect model performance. These challenges highlight the need for incorporating additional contextual information beyond visual features.

Human experts, when diagnosing plant diseases, do not rely solely on visual observation. Instead, they consider a combination of visual symptoms, environmental conditions, disease progression patterns, and prior knowledge. This observation suggests that integrating semantic knowledge with visual data can significantly improve disease detection accuracy. Natural Language Processing (NLP) techniques, particularly transformer-based models such as BERT and FLAN-T5, have shown exceptional capability in capturing contextual and semantic relationships in textual data. These models can encode disease descriptions, symptom narratives, and expert knowledge into meaningful feature representations.[4]

The integration of visual and semantic information leads to the concept of multimodal learning, where multiple data modalities are combined to enhance model performance. Multimodal approaches have gained attention in various domains, including medical diagnosis, autonomous driving, and multimedia analysis. However, their application in plant disease detection, especially for forestry species such as maple plants, remains relatively underexplored. Existing studies primarily focus on unimodal approaches or employ shallow fusion techniques that do not fully exploit cross-modal relationships.

To address these limitations, this paper proposes a novel multimodal maple plant disease detection framework that integrates EfficientNet-based visual feature extraction with transformer-based semantic encoding. The proposed approach utilizes BERT and FLAN-T5 to capture disease-related textual information and employs a Multilayer Perceptron (MLP) to perform feature-level fusion. This design enables the model to learn complex interactions between visual symptoms and semantic knowledge, thereby improving classification performance.[4]

The effectiveness of the proposed framework is evaluated on a balanced dataset comprising 2,000 maple leaf images across five disease classes, along with corresponding textual descriptions. Experimental results demonstrate that the multimodal approach significantly outperforms both vision-only and text-only models, achieving an accuracy of 94.8%. The integration of semantic information proves particularly beneficial in resolving ambiguities associated with visually similar diseases.

The main contributions of this work are summarized as follows:

- A novel multimodal framework for maple plant disease detection
- Integration of EfficientNet with transformer-based language models
- A feature-level fusion strategy using an MLP network
- Comprehensive experimental validation demonstrating improved performance

The remainder of this paper is organized as follows. Section 2 reviews related work in plant disease detection and multimodal learning. Section 3 presents the proposed methodology, including visual and semantic feature extraction and fusion. Section 4 describes the experimental setup and results. Finally, Section 5 concludes the paper and outlines future research directions.

II. LITERATURE REVIEW

The problem of plant disease detection has attracted significant research attention in recent years, driven by advancements in machine learning, computer vision, and deep learning. This section reviews existing work in three major areas: (i) traditional image processing approaches, (ii) deep learning-based visual models, and (iii) emerging multimodal learning techniques.

A. Traditional Image Processing Approaches

Early research in plant disease detection primarily relied on traditional image processing and machine learning techniques. These methods involved manual feature extraction using color, texture, and shape descriptors, followed by classification using algorithms such as Support Vector Machines (SVM), k-Nearest Neighbors (k-NN), and decision trees.[5]

Color-based segmentation techniques were widely used to identify diseased regions on leaf surfaces by distinguishing healthy and infected areas. Texture analysis methods, such as Gray-Level Co-occurrence Matrix (GLCM), were employed to capture variations in leaf surface patterns. Although these approaches provided initial solutions for disease detection, they suffered from several limitations. The reliance on handcrafted features made them sensitive to variations in lighting, background noise, and image quality. Moreover, these methods lacked generalization capability and were not suitable for complex real-world scenarios.

B. Deep Learning-Based Plant Disease Detection

The introduction of deep learning, particularly Convolutional Neural Networks (CNNs), significantly improved the performance of plant disease detection systems. CNNs automatically learn hierarchical feature representations from raw images, eliminating the need for manual feature engineering.

Early CNN-based models such as AlexNet and VGGNet demonstrated promising results in leaf classification tasks. Subsequently, deeper architectures like ResNet addressed the problem of vanishing gradients and enabled the training of very deep networks. These models achieved higher accuracy and robustness compared to traditional approaches.[6]

More recently, EfficientNet has gained attention due to its compound scaling strategy, which balances network depth, width, and resolution. EfficientNet-based models have been successfully applied to plant disease detection, achieving high accuracy while maintaining computational efficiency. However, despite their strong performance, CNN-based approaches rely solely on visual information and often struggle in cases where diseases exhibit similar visual characteristics.

Several studies have attempted to improve CNN performance through data augmentation, transfer learning, and ensemble methods. While these techniques enhance accuracy, they do not address the fundamental limitation of relying exclusively on visual features.

C. Transformer-Based Models and Semantic Understanding

In parallel with advancements in computer vision, Natural Language Processing (NLP) has witnessed significant progress with the introduction of transformer-based models such as BERT and T5. These models use self-attention mechanisms to capture contextual relationships in textual data and have achieved state-of-the-art performance in various NLP tasks.

Transformer models have been applied in agricultural domains for tasks such as crop recommendation, disease description analysis, and knowledge extraction from agricultural texts. These models are capable of encoding complex semantic relationships, making them suitable for representing expert knowledge related to plant diseases.[7]

FLAN-T5, an instruction-tuned variant of T5, further enhances semantic understanding by learning from task-oriented instructions. This makes it particularly effective in capturing disease-related contextual information. However, the application of transformer models in plant disease detection remains limited, especially in integration with visual models.[8]

D. Multimodal Learning for Disease Detection

Multimodal learning involves the integration of multiple data modalities, such as images and text, to improve model performance. In recent years, multimodal approaches have shown success in domains such as medical diagnosis, video understanding, and autonomous systems.

In the context of plant disease detection, a few studies have explored the combination of visual and textual information. These approaches typically use simple fusion strategies, such as concatenation of features or late decision fusion. While such methods demonstrate improved performance compared to unimodal models, they often fail to capture deep interactions between modalities.

One of the key challenges in multimodal learning is designing an effective fusion mechanism that can learn meaningful relationships between heterogeneous features. Shallow fusion techniques are insufficient for capturing complex cross-modal dependencies, leading to suboptimal performance.[9], [10], [11], [12]

E. Research Gap and Motivation

From the literature review, several key gaps can be identified:

- Most existing approaches rely solely on visual information and ignore semantic knowledge
- Transformer-based language models are underutilized in plant disease detection
- Existing multimodal methods use shallow fusion techniques without deep integration
- Limited research has been conducted on forestry-specific disease detection, particularly for maple plants

These limitations highlight the need for a robust multimodal framework that can effectively integrate visual and semantic information.

F. Positioning of the Proposed Work

To address the identified gaps, this paper proposes a multimodal deep learning framework that combines EfficientNet-based visual feature extraction with transformer-based semantic encoding using BERT and FLAN-T5. Unlike existing approaches, the proposed method employs a feature-level fusion strategy using a Multilayer Perceptron (MLP), enabling the model to learn complex relationships between visual and textual features.

The proposed framework not only improves classification accuracy but also enhances robustness and generalization capability, particularly in scenarios involving visually similar diseases. This positions the work as a significant advancement over existing unimodal and shallow multimodal approaches.

III. PROPOSED METHODOLOGY

This section presents the proposed multimodal framework for maple plant disease detection. The methodology integrates visual and semantic information through a unified deep learning architecture. The overall framework consists of three primary components: (i) visual feature extraction using a convolutional neural network, (ii) semantic feature extraction using transformer-based language models, and (iii) feature-level fusion using a Multilayer Perceptron (MLP).

The objective of the proposed approach is to leverage complementary information from image and text modalities to improve classification accuracy, particularly in cases where visual symptoms alone are insufficient for reliable disease identification.

A. Overall Framework

The proposed system takes two inputs:

- A maple leaf image I
- A disease-related textual description T

The framework processes these inputs through separate feature extraction modules and then combines them using a fusion mechanism to produce the final disease classification output.

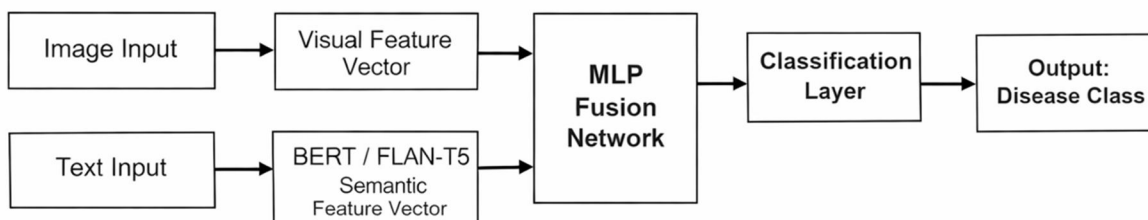


Figure 1: Proposed Multimodal Framework

B. Visual Feature Extraction

Visual features are extracted from maple leaf images using an EfficientNet-based Convolutional Neural Network (CNN). EfficientNet employs a compound scaling strategy that uniformly scales network depth, width, and resolution, resulting in improved performance and efficiency.

Transfer learning is employed by initializing EfficientNet with pretrained weights, followed by fine-tuning on the maple leaf dataset to adapt the model to domain-specific features.

C. Semantic Feature Extraction

Semantic features are extracted from disease-related textual descriptions using transformer-based language models, specifically BERT and FLAN-T5. These models utilize self-attention mechanisms to capture contextual relationships between words and generate meaningful embeddings.

Tokenization, padding, and attention masking are applied prior to embedding generation to ensure uniform input representation.

D. Loss Function

The categorical cross-entropy loss function is used since the maple plant disease diagnosis task can be formulated as a multi-class classification problem. A loss term is used to quantify the difference between the estimated classes probabilities and the actual classes in a multi-class deep learning task and it is popular in multi-class learning.

Let y denote the ground-truth label and \hat{y} denote the predicted probability distribution over disease classes. The loss function is defined as:

$$L = \sum_{i=1}^C y_i \log(\hat{y}_i)$$

where C represents the total number of disease classes.

This loss formulation encourages the model to assign high probability to the correct disease class while penalizing incorrect predictions.

E. Training Strategy

The model is trained in an end-to-end manner using the Adam optimizer. Key training steps include:

- Initialization with pretrained weights
- Fine-tuning of CNN and language model
- Batch-wise training with backpropagation
- Early stopping based on validation performance

Dropout regularization is applied to prevent overfitting and improve generalization.

F. Algorithmic Workflow

The overall workflow of the proposed system is summarized below:

Algorithm 1: Proposed Multimodal Maple Plant Disease Detection

Input:

- Maple leaf image I
- Disease-related textual description T

Output:

- Predicted disease class D

Steps:

1. Acquire input maple leaf image I
2. Resize and normalize image I
3. Pass I through EfficientNet-based CNN
4. Extract visual feature vector V
5. Acquire textual input T
6. Preprocess text (tokenization and cleaning)
7. Pass TTT through transformer-based language model (BERT / FLAN-T5)
8. Extract semantic feature vector S
9. Normalize V and S
10. Concatenate features to form fused vector

$F = [V || S]$

11. Pass F through MLP-based fusion network
12. Compute class probabilities using softmax
13. Assign disease class with highest probability as output D

End

IV. RESULTS AND DISCUSSION

This section presents the experimental results of the proposed multimodal maple plant disease detection framework and provides a detailed analysis of model performance. The evaluation is conducted using the dataset described earlier, consisting of 2,000 maple leaf images across five disease classes, along with corresponding textual descriptions.

All models were evaluated on a held-out testing dataset of 300 images, and performance was measured using accuracy, precision, recall, and F1-score.

A. Performance Comparison

The performance of the proposed multimodal model is compared with vision-only and text-only baseline models. The results are summarized in Table 2.

Table 2: Performance Comparison of Models

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)
EfficientNet (Vision-only)	89.2	88.5	87.9	88.2
BERT (Text-only)	81.6	80.9	80.2	80.5
FLAN-T5 (Text-only)	83.4	82.8	82.1	82.4
Proposed (EfficientNet + FLAN-T5)	94.8	94.2	93.9	94.0

The results clearly demonstrate that the proposed multimodal model significantly outperforms both vision-only and text-only approaches. The EfficientNet-based vision model achieves strong baseline performance, confirming its capability in extracting discriminative visual features. However, its performance is limited in scenarios involving visually similar diseases.

Text-only models based on BERT and FLAN-T5 achieve lower accuracy, indicating that semantic information alone is insufficient for precise disease classification. Among the two, FLAN-T5 performs better due to its enhanced contextual understanding.

The proposed multimodal model achieves the highest performance across all evaluation metrics. The integration of visual and semantic features results in an accuracy improvement of approximately 5–6% over the vision-only baseline, demonstrating the effectiveness of multimodal learning.

B. Confusion Matrix Analysis

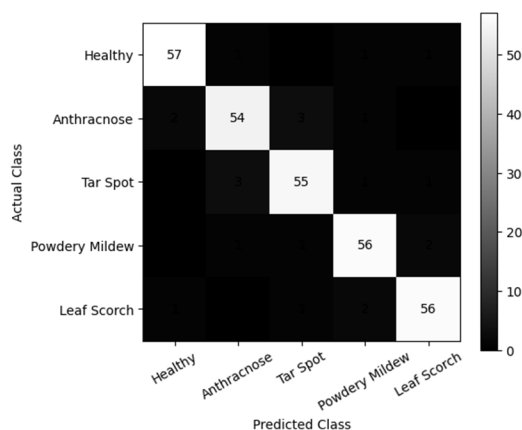


Figure 3: Confusion Matrix

The confusion matrix provides detailed insight into class-wise model performance. Most predictions are concentrated along the diagonal, indicating correct classification across all disease categories.

Minor misclassifications are observed between anthracnose and tar spot, as well as between powdery mildew and leaf scorch, due to similarities in visual symptoms. However, the frequency of such errors is significantly reduced compared to vision-only models, highlighting the advantage of incorporating semantic information.

C. Ablation Study

To evaluate the contribution of different components, an ablation study was conducted. The results are presented in Table 3.

Table 3: Ablation Study Results

Configuration	Accuracy (%)
Vision-only (EfficientNet)	89.2
Multimodal (No Text)	90.1
Multimodal (No Fusion)	92.3
Full Multimodal Model	94.8

The ablation study confirms that each component of the proposed framework contributes to overall performance. Removing semantic information results in a noticeable decrease in accuracy, indicating the importance of textual knowledge. Similarly, disabling feature fusion reduces performance, highlighting the necessity of effective multimodal integration.

The full multimodal model consistently achieves the highest accuracy, demonstrating that both semantic encoding and fusion mechanisms are critical for optimal performance.

D. Comparative Analysis with Existing Methods

The proposed model is compared with existing deep learning-based plant disease detection approaches.

Table 4: Comparison with Existing Methods

Method	Accuracy (%)
VGGNet-based CNN	85.6
ResNet-based CNN	87.8
EfficientNet	89.2
Basic Multimodal	91.3
Proposed Model	94.8

The proposed framework achieves superior performance compared to existing methods. The improvement is primarily attributed to:

- Effective visual feature extraction using EfficientNet
- Rich semantic representation using transformer models
- Deep feature-level fusion using MLP

This demonstrates that multimodal learning provides a significant advantage over traditional unimodal approaches.

E. Key Observations

The following key observations can be drawn from the experimental results:

- 1) Multimodal learning significantly improves classification accuracy
- 2) Semantic information enhances discrimination between visually similar diseases
- 3) FLAN-T5 provides better performance than BERT
- 4) Feature-level fusion is more effective than unimodal approaches

V. CONCLUSION AND FUTURE DIRECTIONS

A. Conclusion

This paper presented a novel multimodal deep learning framework for maple plant disease detection by integrating visual and semantic information. The proposed approach addresses the limitations of traditional vision-based methods, which often struggle to distinguish between diseases with similar visual characteristics.

The framework combines EfficientNet-based convolutional neural networks for visual feature extraction with transformer-based language models, including BERT and FLAN-T5, for semantic feature encoding. These heterogeneous features are effectively fused using a Multilayer Perceptron (MLP), enabling the model to learn complex cross-modal relationships.

Experimental evaluation on a balanced dataset of 2,000 maple leaf images demonstrated that the proposed multimodal model achieves superior performance compared to unimodal approaches. Specifically, the model achieved an accuracy of 94.8%, outperforming the vision-only EfficientNet model (89.2%) and text-only models (81.6%–83.4%). The improvement of approximately 5–6% highlights the effectiveness of integrating semantic knowledge with visual features.

The confusion matrix analysis showed strong diagonal dominance, indicating high classification accuracy across all disease classes. Minor misclassifications were primarily observed between visually similar diseases such as anthracnose and tar spot. The ablation study further confirmed that both semantic information and feature fusion play a critical role in improving performance.

Overall, the results demonstrate that multimodal learning provides a robust and effective solution for plant disease detection. The proposed framework enhances classification accuracy, reduces ambiguity, and improves generalization capability, making it suitable for real-world agricultural applications.

B. Future Directions

Although the proposed multimodal framework demonstrates strong performance in maple plant disease detection, several avenues exist for future research and improvement. One important direction is the expansion of the dataset to include a larger number of samples, diverse environmental conditions, and multiple plant species to enhance model generalization. In addition, optimizing the model for real-time deployment on mobile and edge devices can significantly improve its practical applicability in field conditions. Future work may also explore advanced multimodal fusion techniques, such as attention-based or graph-based approaches, to better

capture complex interactions between visual and semantic features. Furthermore, integrating the system with IoT-based smart agriculture platforms can enable continuous monitoring and automated disease detection. The incorporation of explainable artificial intelligence techniques can improve model interpretability and user trust by providing insights into prediction decisions. Another promising direction is the inclusion of temporal data to analyze disease progression over time, allowing early detection and predictive analysis. Finally, the proposed framework can be extended to other crops and forestry applications, making it a scalable and generalized solution for intelligent plant disease management systems.

REFERENCES

- [1] G. Gupta and S. Kumar Pal, "Applications of AI in precision agriculture," *Discov. Agric.*, 2025, doi: 10.1007/s44279-025-00220-9.
- [2] Olabimpe Banke Akintuyi, "AI in agriculture: A comparative review of developments in the USA and Africa," *Open Access Res. J. Sci. Technol.*, 2024, doi: 10.53022/oarjst.2024.10.2.0051.
- [3] G. Hampel and Z. Fabulya, "The Risks of AI in Agriculture," *Analecta Tech. Szeged.*, 2024, doi: 10.14232/analecta.2024.4.32-44.
- [4] P. Sharma, A. Sanghi, G. Agarwal, and R. Agarwal, "AI in Agriculture: the Future of Farming.," *Grenze Int. J. Eng. & Technol.*, 2025.
- [5] A. A. Mana, A. Allouhi, A. Hamrani, S. Rahman, I. el Jamaoui, and K. Jayachandran, "Sustainable AI-based production agriculture: Exploring AI applications and implications in agricultural practices," *Smart Agric. Technol.*, 2024, doi: 10.1016/j.atech.2024.100416.
- [6] U. A. Okengwu, L. N. Onyejebu, L. U. Oghenekaro, M. O. Musa, and A. O. Ugbari, "Environmental and ethical negative implications of AI in agriculture and proposed mitigation measures," *Sci. Africana*, 2023, doi: 10.4314/sa.v22i1.13.
- [7] T. Miller, G. Mikiciuk, I. Durlik, M. Mikiciuk, A. Łobodzińska, and M. Śnieg, "The IoT and AI in Agriculture: The Time Is Now—A Systematic Review of Smart Sensing Technologies," 2025. doi: 10.3390/s25123583.
- [8] Adebunmi Okechukwu Adewusi, Onyeka Franca Asuzu, Temidayo Olorunsogo, Temidayo Olorunsogo, Ejuma Adaga, and Donald Obinna Daraojimba, "AI in precision agriculture: A review of technologies for sustainable farming practices," *World J. Adv. Res. Rev.*, 2024, doi: 10.30574/wjarr.2024.21.1.0314.
- [9] M. Karnawat, S. K. Trivedi, D. Nagar, and R. Nagar, "Future of AI in Agriculture," *Biot. Res. Today*, 2020.
- [10] M. Javaid, A. Haleem, I. H. Khan, and R. Suman, "Understanding the potential applications of Artificial Intelligence in Agriculture Sector," *Adv. Agrochem*, 2023, doi: 10.1016/j.aac.2022.10.001.
- [11] A. Ahmad et al., "AI can empower agriculture for global food security: challenges and prospects in developing nations," 2024. doi: 10.3389/frai.2024.1328530.
- [12] Olabimpe Banke Akintuyi, "Adaptive AI in precision agriculture: A review: Investigating the use of self-learning algorithms in optimizing farm operations based on real-time data," *Open Access Res. J. Multidiscip. Stud.*, 2024, doi: 10.53022/oarjms.2024.7.2.0023.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)