# Tobacco Leaf Disease Detection Using Deep Learning and Machine Learning: A Comprehensive Review

Priya Sharma
Computer Science and Engineering
Thapar Institute of Engineering and Technology
Patiala, India
Email: psharma13_be23@thapar.edu

*Abstract*—Tobacco leaf diseases pose a serious threat to global agricultural productivity and leaf quality, especially in countries where tobacco is a significant economic crop. Traditional disease detection methods depend heavily on manual inspection by experts, which is not only labor-intensive but also highly prone to subjective bias and human error. Moreover, these conventional techniques are not scalable, especially in large farming areas with diverse environmental conditions. In recent years, advancements in artificial intelligence (AI), particularly machine learning (ML) and deep learning (DL), have revolutionized the field of plant pathology by enabling automated, high-precision disease recognition systems.This review paper provides a comprehensive examination of recent research focused on tobacco leaf disease detection, with an emphasis on image segmentation and classification using convolutional neural network (CNN) architectures. We analyze the evolution of models from simple classifiers to advanced encoder-decoder frameworks and hybrid architectures. Key topics include the use of publicly available datasets, transfer learning, data augmentation, and evaluation metrics such as accuracy, precision, sensitivity, and F1 score. We also discuss how attention mechanisms, generative adversarial network (GAN)-based models, and lightweight mobile architectures improve real-world applicability. Finally, we outline current limitations, highlight unresolved research gaps, and propose future directions aimed at building scalable and efficient precision agriculture solutions.

*Index Terms*—Tobacco plant disease, deep learning, segmentation, convolutional neural networks, U-Net, MD-Unet, transfer learning.

## I. INTRODUCTION

Tobacco (Nicotiana tabacum) is one of the most important crops in terms of commercial production globally [1]. It is a significant contributor to the economy of several countries, such as China, India, and Brazil [2]. While many tobacco products are health-related concerns, tobacco production continues to provide sustenance and income for millions of people. The economic impact of tobacco in various countries is summarized in Table I, highlighting its contribution to agricultural Gross Domestic Product (GDP), employment, and export revenue [3].

Despite its use, tobacco is attacked by a variety of diseases. [4] Foliar diseases are especially detrimental to the yield and quality of the crop and cause significant economic losses [5]. Some common tobacco leaf diseases are angular leaf spot, frog-eye leaf spot, brown spot, and wildfire. Tobacco plants are also infected with fungal and bacterial agents. Viral pathogens, such as Tobacco Mosaic Virus (TMV) [6], also pose a problem as they remain in crop production systems and can reproduce quickly and spread easily to new plants. Figure 1 presents visual examples of these diseases, illustrating the variability in symptoms and leaf damage [7] .

Diagnosing tobacco leaf diseases has traditionally relied solely on manual methods, where agricultural experts visually inspect plants for symptoms and make a diagnosis. This method of detection is widely utilized across the globe, however, it has limitations such as; potential for bias, variability in assurance of correctness of the diagnosis, and dependency upon the availability of the agricultural experts. Manual inspection has limitations in terms of scaling efficiency when utilized on a massive (agri-) industrial scale [9]. Furthermore, many early symptoms of plant diseases are subtle and difficult to visually assess with the human eye and as such allows some time for disease progression before human intervention can occur.Emerging technologies like hyperspectral imaging, thermal imaging, and laser-induced breakdown spectroscopy (LIBS) enhance AI's ability to capture pre-visual symptoms. Deep learning architectures like MD-Unet, TRNet, and OR-AC-GAN are now capable of lesion-level segmentation, even when symptoms are minute or irregular [10].

In the recent years there has been a growing tendency to adopt artificial intelligence (AI) technologies, especially machine learning (ML) and deep learning (DL), to help address the limitations of disease diagnosis of manual inspection [11]. In traditional ML methods, manual feature extraction is typically followed by classification methods such as Support Vector Machines (SVM), Decision Trees (DT), and k-Nearest Neighbors (k-NN) using color histograms as features, descriptors of texture (i.e., Gray-Level Co-occurrence Matrices (GLCM), Local Binary Patterns (LBP), or features related to the shape of the plant or mark [12]. While ML methods have demonstrated effectiveness in controlled environments, they have traditionally performed poorly in field conditions because of variability in the lighting, backgrounds and other clutter, as

TABLE I: Economic Importance of Tobacco in Selected Countries (2023-2024)

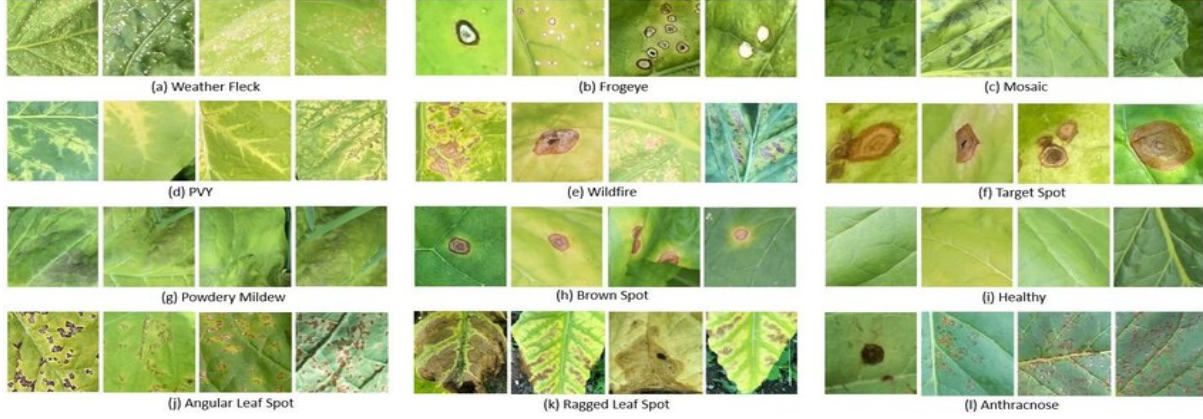| Country | Share in Agricultural GDP (%) | Employment (in millions) | Tobacco Export Revenue (USD) |
|---------|-------------------------------|--------------------------|------------------------------|
| India | 4.6% | 4.5 | $923 million |
| China | 7.2% | 5.3 | $2.1 billion |
| Brazil | 5.1% | 0.2 | $2.9 billion |
| Zimbabwe | 15.8% | 0.1 | $1.0 billion |
| Indonesia | 3.2% | 0.8 | $900 million |
| Bangladesh | 2.5% | 0.6 | $68 million |



Fig. 1: Visual examples of common tobacco leaf diseases including Weather Fleck, Frogeye, Mosaic, PVY, Wildfire, Target Spot, and others [8].

well as leaf morphology.

Even with progress, tobacco leaf disease datasets are small and not very diverse when compared to datasets for other crops. This shortage means that transfer learning techniques are needed. With these, models trained on large datasets, like ImageNet, are adjusted for particular jobs. ResNet, Visual Geometry Group (VGG), MobileNet along with EfficientNet are common networks used for this. Data changes, such as turning, flipping, cutting as well as changing brightness, happen often - this helps prevent too much fitting and improves how well the model works generally.Standard CNN structures are used, but hybrid models that mix CNNs with other methods have also been looked at - these models combine transformers, like TR-Net, to better understand context, plus adversarial training, like OR-AC-GAN, to change data. Lightweight structures, such as MobileNetV2 and ShuffleNet, are made for use on devices with few resources. These devices include drones, phones in addition to farm robots. Such changes allow for quick disease finding and choice making in farm settings.Several problems remain when building also using AI systems for tobacco disease detection - it is still hard to get large datasets that are marked the same way and are open to the public. Differences in how pictures are taken, how data is marked next to how classes are spread out also make model building and testing harder. Models are not always easy to understand. That stops their use in farming, since people involved want to see why decisions are made [13].

New research directions address these problems; they include Explainable AI (XAI), Self-Supervised Learning (SSL), Federated Learning (FL), and multimodal data fusion.XAI methods, like Gradient-weighted Class Activation Mapping (Grad-CAM), SHapley Additive exPlanations (SHAP), plus Local Interpretable Model-agnostic Explanations (LIME), show how models decide [14]. This builds trust and helps people use them. SSL uses unlabeled datasets for tasks, such as contrastive learning also image reconstruction. So it needs fewer labeled examples. FL lets models train together across separate data sources, and it does not need data sharing - this deals with privacy worries.Multimodal approaches join data from different sensors, which include Near-Infrared (NIR), thermal, hyperspectral along with RGB imaging. They collect more information as well as so perform better. Such plans find body changes linked to illness before symptoms appear. Hyperspectral imaging (HSI), with DL models, may find illness early by looking at spectral signs. But HSI systems cost a lot and are hard to run. That stops people from using them widely [15].

### A. Major Contributions

This review makes the following key contributions:

- Presents a thorough survey of ML and DL methods for tobacco leaf disease detection, tracing the shift from SVMs to attention-based CNNs, transformer hybrids, and GAN-augmented models.
- Compares segmentation models (U-Net, MD-Unet, DeepLabV3+, TRNet) using metrics like accuracy, F1-score, IoU, and Dice coefficient, along with visual and architectural insights.

- Analyzes datasets such as Prakasam [16], Luoyang [8], and TPDD [17], covering structure, diversity, limitations, and their roles across tasks.
- Reviews advancements in transfer learning, data augmentation, model compression (e.g., pruning, quantization), and deployment on edge devices (e.g., MobileNetV2, YOLO-Tobacco).
- Identifies key challenges—small datasets, domain shift, lack of standardization, and model interpretability—and outlines future directions like multimodal sensing, self-supervised, and FL.

The taxonomy of the work presented in this article is given in Fig. 2.

The rest of the manuscript is organized as follows.

Section II describes various tobacco leaf diseases, their causes, symptoms, and economic impacts.

Section III outlines the key datasets used in tobacco disease detection research, including their characteristics, modalities, and relevance to classification and segmentation tasks.

Section IV reviews traditional ML and modern DL models applied to tobacco disease detection, highlighting segmentation architectures, hybrid approaches, and lightweight implementations.

Section V presents the evaluation metrics commonly used to assess model performance, such as accuracy, F1-score, IoU, and Dice coefficient.

Section VI provides a comparison of recent survey articles, highlighting evolving trends, architectures, and deployment challenges.

Section VII discusses current limitations in the field and outlines directions for future research.

Section VIII concludes the paper.

Lastly, references are listed.

## II. Tobacco Leaf Diseases and Their Characteristics

Tobacco plants are open to many leaf and root diseases. A lot of these diseases have bad economic effects [18]. People place these diseases into groups based on what causes them – these causes are fungi, bacteria, viruses, or abiotic factors like the environment or plant processes, as Fig. 3 shows. Finding and acting on diseases quickly helps reduce crop loss. In recent years, AI plus DL systems have found use for this job.

Fungal infections affect tobacco plants most often. The plant gets fungal infections because it reacts to moisture and stress from its surroundings [19]. An example of such a disease is Frogeye Leaf Spot, which the fungus *Cercospora nicotianae* causes – it appears as round, tan to gray spots with dark brown edges. Another significant fungal disease is Brown Spot, caused by *Alternaria alternata*, which produces brown necrotic lesions with concentric rings. Phytophthora nicotianae causes another disease called Black Shank (taken as a whole). Black Shank is a soil-borne pathogen that affects both the root and stem components of the plant, leading to vascular damage and wilting.Bacterial diseases include Bacterial Wilt (water transport through the xylem is impeded), caused by

*Ralstonia solanacearum*; wilting results. Wildfire, caused by *Pseudomonas syringae pv. tabaci*, results in water-soaked spots and necrosis.Viruses include Tobacco Mosaic Virus (TMV), Cucumber Mosaic Virus (CMV), and Potato Virus Y (PVY). Viruses induce a variety of symptoms including mosaic patterns, leaf deformation, and stunting. TMV spreads locally with contaminated tools and can remain stable in dried plant debris, making sanitizing one of the most important options for control [20].

Abiotic stressors must also be taken into consideration as they mimic disease. For example, ozone injury, drought, and nutrient deficiencies do produce some similar symptoms as pathogenic agents, making causality fifty percent of the time. Common short-term symptoms are yellowing (chlorosis), flecking, and marginal burn, that may be mistaken for fungal or bacterial infection [21].

Disease progression is dependent upon the type of pathogen and the environmental condition [23]. Diseases caused by fungi, such as black shank, usually start as small leaf spots and grow in both size and number as humid conditions continue. Viral infections, such as tobacco mosaic virus (TMV) and cucumber mosaic virus (CMV), start at the cellular level, progress slowly over several weeks, and ultimately cannot be remediated as they cause irreversible physiological damage [19].

The economic impact caused by disease progression should not be overlooked. For instance, diseases like black shank and bacterial wilt can cause complete yield loss if no precautions are taken. Even minor diseases cause quality of the leaf to diminish in a variety of ways including flue-cured leaf, texture, nicotine production, and price point in the market of the processed tobacco [24]. Literature reports yield losses ranging from 15% to 80% depending upon the severity of the disease [20].Diseases affected physiological processes are helpful in understanding how secondary metabolites are produced through alteration of sugar–nicotine balance, chlorophyll degradation, and final production viability for end products like cigars, cigarettes, and chewing tobacco. Disease detection models, such as AI models, may be simply programmed to detect lesions, so they must also be aware of all the subtle phenotypic changes in the early stages of disease and put all of the pieces together.The environment is essential for tobacco diseases to start, develop, and flourish. If there is excessive rainfall and/or poor drainage in the planting beds, those conditions lend to fungal disease spores reproducing and developing. Conditions like drought and high temperatures will suppress the tobacco plant's immune response, leading to higher vulnerability to these pathogens [25].Changes in the climate across the globe are leading to disease zones shifting; diseases that were previously constrained to a specific geographic region are now appearing in completely different geographies. For example, PVY, which has been historically associated with the temperate zones, is coming to be found with increasing frequency in the tropical regions. Other examples include rainfall variability and its effect on the incidence of bacterial wilts, to mention only two factors [26].
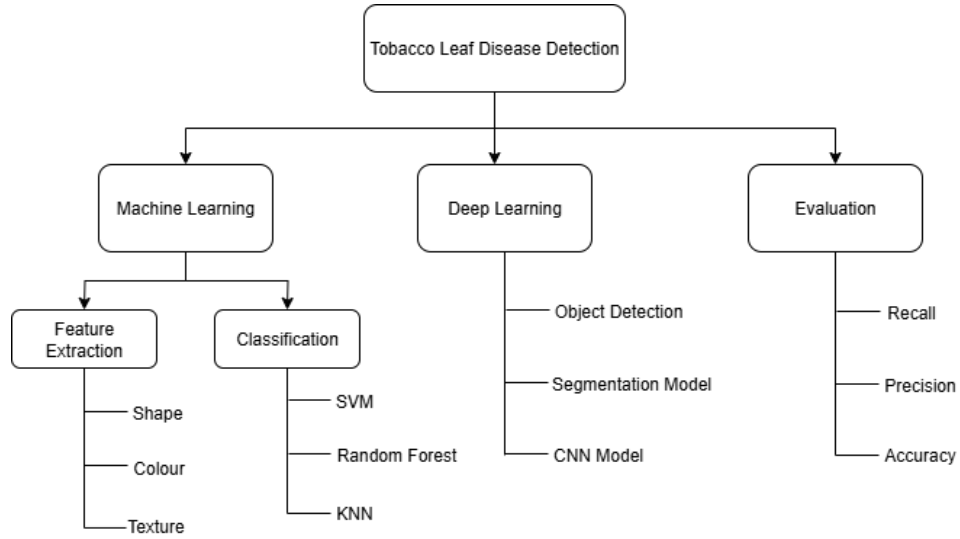
Fig. 2: Taxonomy of the proposed work

TABLE II: Common Tobacco Leaf Diseases and Their Characteristics [22]

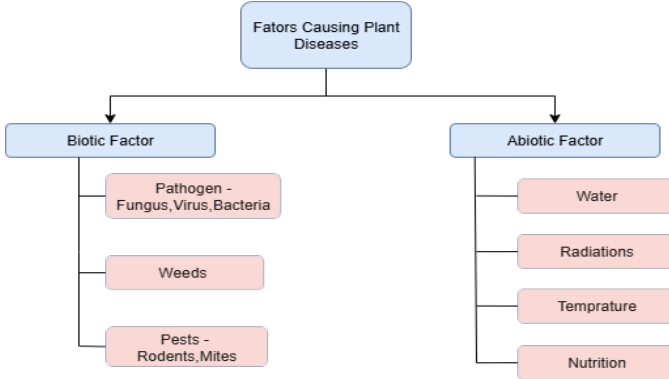| Disease | Pathogen | Symptoms | Transmission | Type |
|---|---|---|---|---|
| Frogeye Leaf Spot | *C. nicotianae* | Gray lesions, brown borders | Wind, tools | Fungal |
| Brown Spot | *A. alternata* | Concentric brown lesions | Rain splash, debris | Fungal |
| Black Shank | *P. nicotianae* | Stem rot, root decay | Soil-borne | Fungal |
| Bacterial Wilt | *R. solanacearum* | Sudden wilting, xylem browning | Soil, water | Bacterial |
| Wildfire | *P. syringae pv. tabaci* | Leaf burn, necrosis | Tools, splash | Bacterial |
| TMV | Tobacco Mosaic Virus | Mosaic, curling, stunting | Tools, contact | Viral |
| CMV | Cucumber Mosaic Virus | Chlorotic mottling | Aphids | Viral |
| PVY | Potato Virus Y | Vein necrosis, mottling | Aphids, mechanical | Viral |
| Chlorosis | Nutrient Deficiency | Yellowing of leaves | Soil quality | Abiotic |
| Ozone Injury | Environment | Leaf flecks, necrosis | Air pollution | Abiotic |



Fig. 3: Classification of factors causing plant diseases into biotic and abiotic categories.

AI models that are trained in one region cannot be expected to generalize effectively to different climatic conditions unless they are trained in that specific region. Therefore, there is a requirement for ongoing expansion of the dataset and transfer learning methods to ensure relevance of the models over time under changing agroecological conditions [27].

Tobacco is a cash crop in China, largely in India but with others being Brazil, Zimbabwe, and a number of other countries. Tobacco exports in India were over $900 million U.S. dollars in the 2023–2024 fiscal year [28]. Disease outbreaks can have a significant economic impact on farmers' income, their export potential, and the downstream processing industries.Diseases like TMV generally limit yield but also present costs associated with testing, quarantining, and decontamination, on an already significantly distressed crop. They must incur additional cost in weekly or monthly recurring input costs associated with fungicides and bactericides. In many low-income regions, farmers do not have access to certified diagnostic labs and are dependent on their broader experience or unsophisticated informal advice.

AI-based mobile detection methods that are diagnostic and cost-effective provide a greater opportunity to step into early on-field detection to minimize mistaken reliance on reactive treatment and begin to implement preventative protection—thereby providing the farmer with increased economic resiliency [29].

## III. DATASETS USED IN TOBACCO LEAF DISEASE DETECTION

The construction of effective DL and ML models to detect tobacco leaf disease relies heavily on datasets. Numerous

research studies during the period of 2021-2025 adopted varying degrees of datasets that accompanied different types of disease, environments, imaging methods, and annotation. While all studies pursued the goal of disease recognition, these datasets exhibit different characteristics in terms of usage, which contributes to differences in model performance, generalizability, and end-user applicability.

In these studies, some datasets consisted of laboratory-based images, and others were field-based. Laboratory based image datasets are often high-resolution images taken in controlled lighting and with a controlled background to reduce noise and variability in the images to showcase more clearly the disease lesions. But the laboratory dataset with controlled acquisition takes away the real-world aspect to capturing images of tobacco leaf disease. While field datasets are used to showcase real-world effectiveness of the model or method to detect disease during environmental factors such as light change, potential occlusions from other foliage, and, background images that obscured the true nature of the disease (more noisy and variable for end-user applicability), field datasets are most representative for producing models deployable in the field.

One of the mostly widely utilized datasets is from the Prakasam region in Andhra Pradesh, India. This dataset has over 3,200 images taken with high-resolution via smartphones in nature light sources. There are various categories of disease; Wildfire, Frog-Eye Spot, and Brown Spot, positioning the dataset as annotated by agronomists. Due to the natural diversity of the dataset in a real-world application it has been used in some CNN and YOLO model architectures, highlighting classification or detection [16].

Another dataset is called the Luoyang Tobacco Dataset with more than 4,000 images with pixel-level annotations. There are multiple diseases represented which include Angular Leaf Spot, Mosaic Virus, and Wildfire under lab and field conditions. Expert annotations provide pixel-level lesion masks making it suitable to be used to train models with segmentation, like MD-Unet and U-Net variations [8]. The dataset has multi-class and multi-label scenarios making it potentially useful for classification and segmentation with fine granularity.

The TPDD (Tobacco Plant Disease Dataset) represents a hybrid approach, aggregating images from various environments and encompassing over 5,000 annotated images. It includes diverse tobacco disease types, and its balanced class distribution makes it suitable for benchmarking multiple tasks like classification, segmentation, and object detection [17]. This dataset is frequently used in models employing transfer learning with architectures like DeepLabV3+, MobileNet, and Mask R-CNN.

Multimodal datasets also appear in the literature that are reviewed here. In some instances, the datasets combine RGB images with NIR spectroscopy to measure physiological changes in the leaves before visual symptoms are presented. While generally small in number of images (often under 1,000), these datasets are augmented with spectral information that is useful for models that focus on detecting biological conditions in an early state [30]. These are typically lab-generated under strict conditions to control for light interference and spectral overlap.

The studies use different data formats. The classification datasets use JPEG or PNG images with a label for each image. The segmentation tasks have masks which are stored as either grayscale images or binary PNGs. In addition, segmentation datasets may also have XML or JSON files for the mask annotations which include bounding box or polygon annotations. The most common data formats for detection datasets use YOLO style .txt files or Pascal VOC style XML to denote bounding box coordinates.

The methods and quality of the annotated data varies widely. The most common data annotation makes use of domain experts to manually segment the data, this is particularly true for segmentation datasets. The experts will use some form of image annotation software like LabelMe or even Adobe Photoshop to annotate the lesions pixel-by-pixel. Some studies that has a mask as their output, will use some form of semi-autonomous methods to generate the mask, such as color thresholding or watershed segmentation, and followed by manual correction. Generally, the accuracy of the data annotation will directly impact the accuracy of the model being trained, especially for segmentation model trained on high resolution images.

Data augmentation was a common theme across the studies. Several studies reported using augmentations like image rotation, flipping, cropping, scaling, color jitter, histogram equalization and adding noise. Augmenting images allows the model to generalize better, especially when training models on limited and/or imbalanced data. There are some studies that used the Prakasam dataset, and then augmented all of the synthesized data, and reported the accuracy as improved and overfitting was reduced [16].

Despite advancements in understanding tobacco diseases and the data sharing capabilities research, the majority of datasets still have some challenges. Many datasets are not standardized in annotation format, class labels, and imaging protocols. In addition, relatively few datasets are publicly available so it is difficult to demonstrate reproducibility and cross-compare studies. Model generalization across geographic and environmental contexts may also be impacted by data biases. Many diseases that are present in one geographic region may not be present in another geographic region. Therefore, datasets specific to a region are especially important to develop functional models to deploy those models locally and responsibly [31].

Generative data synthesis, with a particular emphasis on Generative Adversarial Networks (GAN), is becoming an increasingly common approach to adress the issue of limited datasets. The OR-AC-GAN framework is one framework that was used to synthesize rare lesion types, meaning that when we balanced the classes in the dataset we would be able to test models backed by synthetically produced, but still realistic lesions. The synthetically produced data would be used to improve classification and segmentation models for the diseases that were less prevalent in the dataset [32].

Ultimately, we found that data quality and diversity have been integral to maximizing the success of tobacco disease detection systems. In one of our recent studies we demonstrated that combining annotator high-quality annotation, real-world aspects of variability and data augmentation therefore established a robust baseline effect on model performance moving forward [33]. Going forward, we believe that the future should focus on standardization, public sharing and multimodal datasets to develop surgical disease detection systems for effective and efficient outcomes for tobacco protection across geographical extents [31].



(a) Original Tobacco Leaf Image    (b) Original Segmentation Mask

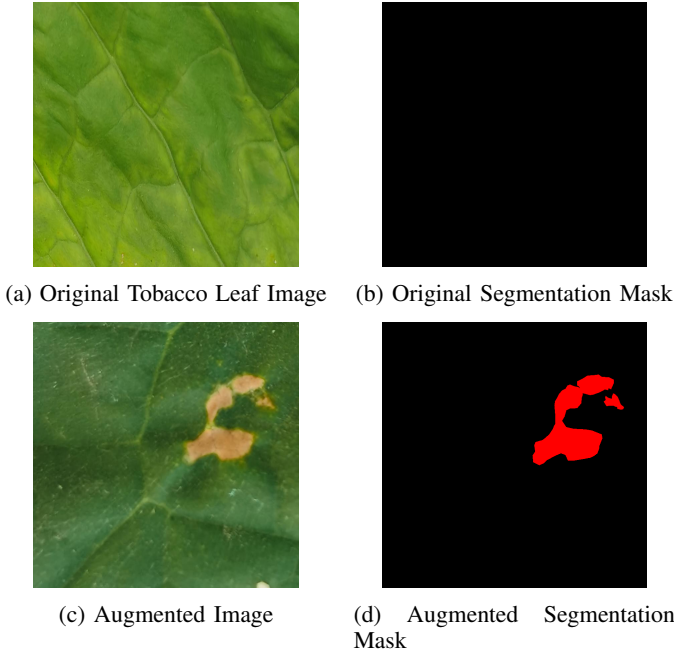(c) Augmented Image    (d) Augmented Segmentation Mask

Fig. 4: Sample images from the Luoyang Dataset [8]. The top row shows an original image and its mask; the bottom row shows their augmented counterparts.

Tobacco leaf disease detection using DL and ML is highly dependent on the quality, diversity, and annotation depth of the datasets used for training and validation. In this review, we have analyzed recent public datasets that vary in scale, modality, annotation type, and use case applicability.

Among the surveyed datasets, the *TLA Dataset* [34] provides one of the most comprehensive field-collected resources, comprising 16 categories that span infectious diseases, abiotic stressors, pest traces, and healthy samples. It includes both raw and processed settings, with expert-labeled, high-resolution images taken in real agricultural environments. The processed setting ensures single-disease samples and class balancing, while the raw setting reflects real-world distribution with long-tailed characteristics. This dataset is particularly valuable for few-shot learning, meta-learning, and robust disease classification under field variability.

The *Prakasam Dataset* offers robust classification capabilities in real-field conditions and is best suited for training lightweight models for disease recognition tasks. In contrast,

the *MD-Unet Leaf Spot Dataset* provides pixel-level annotations and is optimal for segmentation architectures, particularly those focused on fine-grained lesion delineation. For early-stage disease identification, the *NIR-Tobacco Dataset* enables the use of multimodal learning strategies by combining RGB and NIR imaging, though its use is mostly limited to controlled environments.

For real-time and edge-based applications, such as deployment on drones or low-resource devices, the *YOLO-Tobacco Set* is valuable due to its small size and bounding box annotations compatible with YOLOX-Tiny and similar detection networks. Meanwhile, the *GAN-Augmented Data* addresses the challenge of class imbalance by synthetically generating rare lesion types, and is particularly useful when combined with real datasets to enhance model robustness.

Ultimately, no single dataset is universally sufficient. Models trained for real-world agricultural deployment should ideally integrate multiple datasets—field-collected for realism, lab-collected for clarity, and synthetic for balance. Transfer learning, multimodal fusion, and adversarial data augmentation remain essential to compensate for the current lack of standardized, large-scale tobacco disease repositories [37].

## IV. TECHNOLOGY AND MODELS FOR TOBACCO LEAF DISEASE DETECTION

### A. Traditional and Deep Learning Models for Detection and Classification

Over the previous decade, there has been a notable surge of technological innovations which seek to enhance the identification and categorization of tobacco leaf diseases. Early approaches focused primarily on conventional picture processing and classical ML pipelines, contrastingly, recent efforts focus on approaches within DL, attention mechanisms, fused spectral data, and hybridized models. In this section we review these models and the technology they involve, synthesising their comparative performance, dataset characteristics, and significant technical advances.

Traditional ML methods dominated initial efforts behind automated disease detection. These methods historically operate on a pipeline consisting of image preprocessing, mathematically crafted feature extraction, and final supervised classification.

Common feature descriptors include:

- Color histograms, used to capture lesion pigmentation;
- Texture metrics, such as GLCM, entropy, and energy for analyzing lesion structure;
- Morphological features, such as leaf area, vein curvature, and boundary irregularity;
- Edge-based operators, such as Sobel and Canny, used for outlining lesion regions.

The classifiers that have been used include SVM with either a RBF or polynomial kernel [38], k-NN, DT, and Naîve Bayes. In one application, a GLCM + SVM combination provided up to 90.1% accuracy in classifying brown spot lesions on a tobacco leaf [38], though the classifiers exhibited limited

TABLE III: Comparative Overview of Recent Tobacco Leaf Disease Datasets (2021–2025)

| Year | Dataset Name | Source | Size of Images | Best Characteristics | Shortcoming | Usage |
|---|---|---|---|---|---|---|
| 2025, [34] | Tobacco Leaf Abnormality (TLA) Dataset | Mexico (Field) | ∼6,000 (raw + processed) | 16 categories (infections, pests, abiotic, healthy); high-resolution (5184×3456); expert-labeled; segmented and raw modes | Long-tail distribution; class imbalance in raw form | Classification, Few-shot, Meta-learning |
| 2024, [8] | MD-Unet Leaf Spot Dataset | Luoyang, China | ∼2,140 images | High-res images with pixel-level lesion masks and attention modules | Limited to 4 diseases, lacks environmental variety | Segmentation |
| 2024, [16] | Prakasam Dataset | India (Field) | >1,000 images | Field-collected images; annotated by experts; suitable for CNNs | Small scale, lacks segmentation masks | Classification, Detection |
| 2024, [30] | NIR-Tobacco Dataset | Lab (Controlled) | <1,000 images | Multimodal (RGB + NIR); detects early-stage symptoms | Controlled setup only, lacks field variability | Early Detection |
| 2023, [35] | YOLO-Tobacco Set | China (Field) | 340 images | Optimized for YOLOX-Tiny; real-time edge detection | Small size, few classes, detection-only | Lightweight Detection |
| 2023, [36] | GAN-Augmented Data | Synthetic (Hybrid) | Varies | Generates rare lesion classes using GANs; balances dataset | Synthetic bias may affect realism | Data Augmentation |
| 2022, [17] | Tobacco Plant Disease Dataset (TPDD) | Multi-region (SPIE) | ∼2,721 images | Includes whole-leaf + fragment labels; 12 disease classes; bbox annotations | No segmentation masks, only RGB modality | Classification, Detection |

robustness when light varied, the orientation of leaves was not consistent, or when leaves displayed mixed disease symptoms.

The limitations of traditional ML classifiers have resulted in a switch to DL based models, that can learn hierarchical representations directly from raw input data, removing the design and choice of features stage entirely as it extracts them from the input images itself. CNNs have proven very effective in furthering the identification of tobacco diseases. Standard CNNs trained on tobacco datasets, such as those captured from farms in China and India, have reached impressive performance. For example, Peda Babu et al. [16] used a custom CNN architecture to classify three major diseases with an accuracy of 95% using a Prakasam district dataset.

Given the limited size of tobacco-specific datasets, many researchers utilize transfer learning from ImageNet-pretrained backbones (e.g., VGG16, ResNet50, MobileNet). Transfer learning accelerates convergence and enhances generalization when adapted to tobacco-specific tasks.

In addition, aggressive data augmentation techniques such as random flipping, rotation, cropping, and brightness variation have been shown to mitigate overfitting. Models like EfficientNet-B0 and MobileNetV2 are commonly used for this purpose [36], [39].

Recent research has explored hybrid models combining CNNs with other architectures:

- TRNet [36] augments CNNs with transformer-based self-attention layers, improving contextual reasoning;
- YOLO-Tobacco [35] integrates CBAM (Convolutional Block Attention Module) with YOLOX-Tiny for lightweight, real-time detection;
- OR-AC-GAN and ASNet [39] use adversarial training and shuffle attention modules to learn more robust representations.

Some studies integrate non-RGB inputs such as NIR and hyperspectral imaging. Ying et al. [30] utilized a CNN fused with NIR reflectance data to capture early physiological symptoms invisible in RGB imagery. The model performed well in lab-controlled environments, although field variability remains a challenge.

LIBS, as demonstrated by Peng et al. [40], enables chemical-level disease classification by analyzing elemental signatures (e.g., Ca, Mg, H). SVM, PLS-DA, and BP-NNs were applied to classify mosaic virus presence, achieving up to 94.4% accuracy on dried samples.

For practical application in smart farms, field robots, and UAVs [41], researchers have prioritized model compression and hardware efficiency. Lightweight architectures like:

- MobileNetV2,
- ShuffleNet,
- YOLOX-Tiny with CBAM,

have been successfully deployed on Raspberry Pi and Jetson Nano platforms. Techniques such as pruning, quantization, and knowledge distillation allow real-time detection with negligible accuracy drop [35].
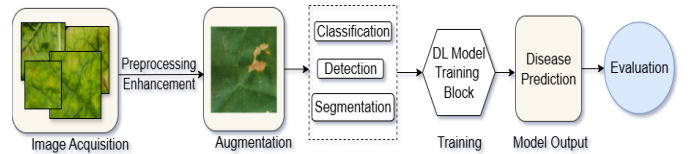


Fig. 5: General pipeline for tobacco leaf disease detection using deep learning.

### B. Segmentation-Based Approaches for Lesion Localization

More advanced studies have focused on segmentation networks like U-Net and its variants:

- MD-Unet, proposed by Chen et al. [8], incorporates attention-based fusion with multi-scale dense blocks, achieving 94.67% accuracy;
- TDSSNet, developed by Ou et al. [42], is a lightweight CNN segmentation network optimized for field images;

- DeepLabV3+ with FPN, adapted from grape disease detection [43], has also been tested on tobacco;
- Inception U-Net, as explored in Zhao et al. [44], integrates inception modules to improve boundary delineation.

## V. EVALUATION METRICS FOR TOBACCO DISEASE DETECTION

To meaningfully evaluate classification and segmentation models in tobacco leaf disease detection, and compare different methodologies, accurate and reliable measurement is important. The performance of ML and DL models is typically measured with a set of quantitative metrics, with specific definitions. These metrics indicated how well a model could identify diseased areas, distinguish between one type of disease from another, and generalize to different datasets. This section describes the evaluation metrics used most commonly in this space.

### A. Accuracy

Accuracy is the most straightforward performance metric [47] and is defined as the ratio of correctly predicted instances to the total number of predictions:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

where $TP$ is the number of true positives, $TN$ the true negatives, $FP$ the false positives, and $FN$ the false negatives. While accuracy is widely used in binary and multi-class classification problems, it may be misleading when class distributions are imbalanced, which is often the case in agricultural datasets.

### B. Precision

Precision quantifies the number of correct positive predictions among all positive predictions:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2)$$

A high precision value indicates that the model makes few false positive errors, which is important when misidentifying a healthy leaf as diseased could lead to unnecessary treatment or crop removal.

### C. Recall (Sensitivity)

Recall, also known as sensitivity or true positive rate, measures the proportion of actual positives correctly identified:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3)$$

Recall is especially important in disease detection tasks, where failing to identify diseased plants (false negatives) can lead to uncontrolled disease spread.

### D. F1 Score

The F1 Score is the harmonic mean of precision and recall:

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

It is particularly useful when dealing with imbalanced datasets or when both precision and recall are important.

### E. Intersection over Union (IoU)

In segmentation tasks, the IoU, also called the Jaccard Index, is used to evaluate the overlap between the predicted lesion area and the ground truth:

$$\text{IoU} = \frac{|A \cap B|}{|A \cup B|} \quad (5)$$

Here, $A$ is the predicted segmentation region and $B$ is the ground truth region. An IoU of 1.0 indicates perfect overlap, whereas 0.0 means no overlap.

### F. Dice Coefficient

The Dice Coefficient is another metric for segmentation tasks, particularly effective for small lesion regions:

$$\text{Dice} = \frac{2|A \cap B|}{|A| + |B|} \quad (6)$$

Compared to IoU, the Dice score places more emphasis on the intersection, making it more sensitive to small segmentations.

### G. Specificity

Specificity, or the true negative rate, measures the ability of a model to correctly identify non-diseased instances:

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (7)$$

High specificity is important when it is critical to avoid false alarms, such as overestimating disease incidence in healthy leaves.

### H. Model Robustness and Generalization

While quantitative metrics are vital, model performance must also be assessed in terms of robustness and generalization. This includes evaluating how well models perform across:
- Different lighting conditions and backgrounds,
- Varied leaf sizes, shapes, and orientations,
- Cross-dataset scenarios (training on one dataset and testing on another).

These factors are often explored through cross-validation and transfer learning experiments and should be included in any comprehensive model evaluation. Table V provides a brief overview of which metrics are commonly applied to different tasks within tobacco disease detection.
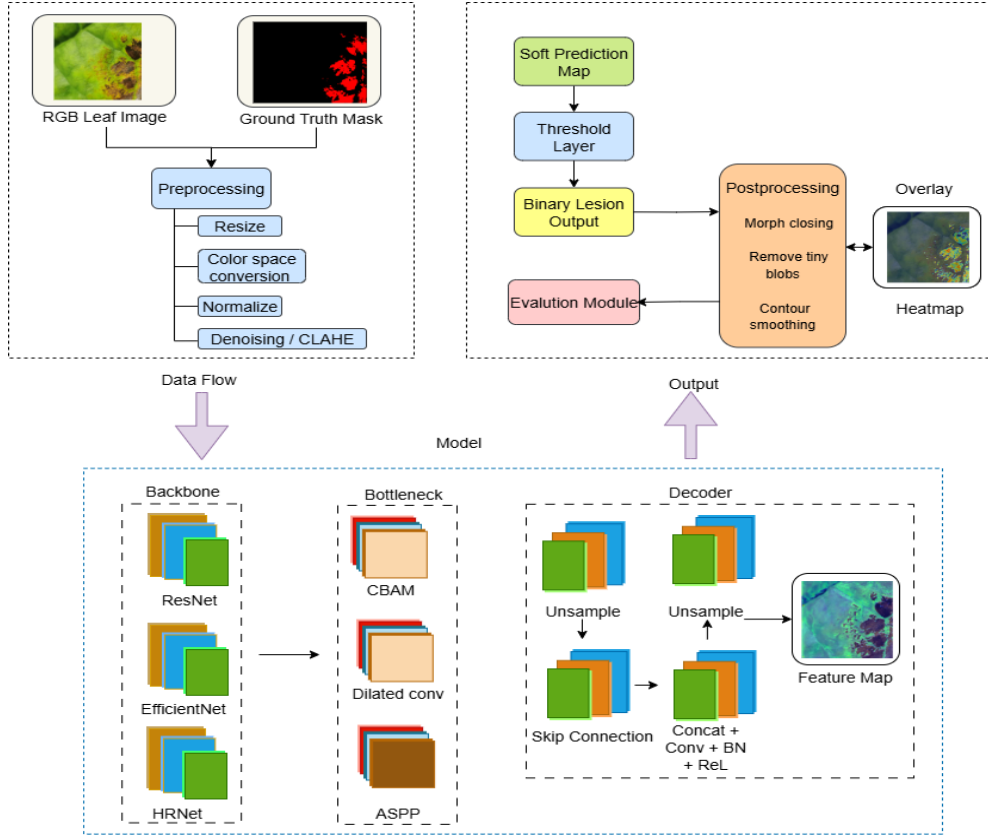
Fig. 6: Functional block diagram for segmentation-based tobacco disease detection.
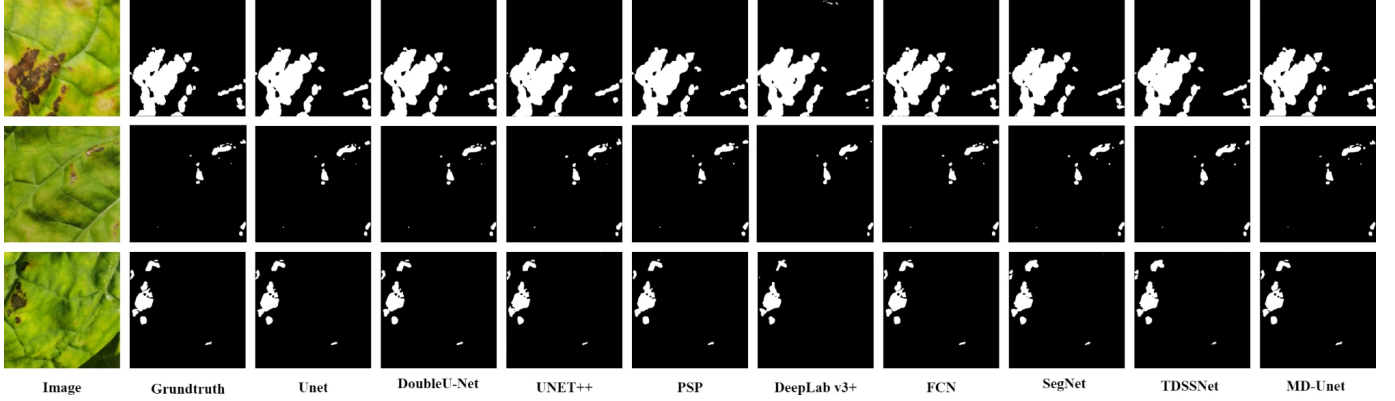


Fig. 7: Segmentation results of tobacco leaf disease using various models [34].

TABLE V: Recommended Metrics for Disease Detection Tasks

| Task | Metric | Focus |
|------|--------|-------|
| Classification | Accuracy, F1 Score | Overall detection rate |
| Segmentation | IoU, Dice Score | Lesion boundary accuracy |
| Binary Detection | Precision, Recall | Error sensitivity |
| Generalization | Cross-validation score | Dataset shift resistance |

These evaluation metrics serve as the foundation for comparing the tobacco disease detection models discussed in Section IV. The interpretation of these metrics in the con-text of model architectures, dataset diversity, and real-world deployability provides valuable insight into the strengths and limitations of each approach.

Table VI shows a complete comparison of different models employing state-of-the-art image segmentation techniques on the Luoyang dataset that used for tobacco leaf disease identifi-cation. In evaluating the models, we used a number of common metrics - the Correct Pixel Accuracy (CPA), Recall, the IoU, the F1-score and the Dice coefficient - for both the lesion and leaf segmentation task. Taken together, the metrics provide an assessment of each model's overall ability to segment

TABLE IV: Comparative Analysis of Tobacco Leaf Disease Detection Models (2021–2025)

| Author | Year | Model | Technique | Accuracy | Dataset | Highlights |
|--------|------|-------|-----------|----------|---------|-----------|
| Chen et al. [8] | 2025 | MD-Unet | DL + CBAM Attention | 94.67% | Luoyang Tobacco Leaves | Dense block + attention fusion |
| Ou et al. [42] | 2024 | TDSSNet | Lightweight CNN | 89.74% | Field-Captured Leaves | Optimized for segmentation |
| Peda Babu et al. [16] | 2024 | CNN | DL Classifier | 95.00% | Prakasam District | 3-class disease detector |
| Ying et al. [30] | 2024 | CNN + NIR | Spectral Fusion | – | Lab-captured | Reflectance-based features |
| Sindhu et al. [45] | 2024 | DBN | Deep Belief Network | 86.18% | Self-captured TMV | Focused on TMV detection |
| Lin et al. [35] | 2023 | YOLO-Tobacco | DL + CBAM | – | 340 Images | YOLOX-Tiny fused with CBAM |
| Yao et al. [36] | 2023 | TRNet | Transformer-CNN | 93.60% | Tobacco Leaves | High performance via attention |
| Yuan et al. [43] | 2023 | DeepLabV3+ | DL + FPN | 91.80% | Grape Leaf (Ref.) | Effective with ECA blocks |
| Zhao et al. [44] | 2023 | Inception U-Net | DL + Multi-Scale | 92.90% | Tomato (Ref.) | Inception module fusion |
| Teng et al. [38] | 2021 | SVM (Tree) | ML Classifier | 90.10% | 200 Tobacco Images | GLCM features |
| Fu et al. [46] | 2022 | RS-UNet | DL + SE Block | 88.86% | Potato Dataset | Residual skip features |
| He et al. [39] | 2022 | ASNet | Asymmetric ShuffleNet | 84.22% | Apple (Ref.) | scSE attention |
| Peng et al. [40] | 2017 | SVM / BP-NN / PLS-DA | LIBS + ML | Up to 94.4% | 160 Samples | Elemental analysis |

TABLE VI: Performance Comparison of Segmentation Models on the Luoyang Dataset [8]

| Model | Lesions | | | | Leaf | | | | Dice |
|-------|------|--------|-----|----|------|--------|-----|----|------|
| | CPA | Recall | IoU | F1 | CPA | Recall | IoU | F1 | |
| Unet | 90.45 | 88.80 | 81.59 | 89.81 | 99.83 | 99.87 | 99.70 | 99.85 | 93.21 |
| DoubleU-Net | 90.21 | 88.82 | 81.06 | 89.49 | 99.83 | 99.86 | 99.69 | 99.85 | 92.92 |
| UNET++ | 87.32 | 84.75 | 75.42 | 85.89 | 99.77 | 99.82 | 99.58 | 99.79 | 90.55 |
| PSP | 89.45 | 87.49 | 79.38 | 88.43 | 99.81 | 99.85 | 99.66 | 99.83 | 91.91 |
| DeepLab v3+ | 77.85 | 81.58 | 66.22 | 79.37 | 99.73 | 99.69 | 99.43 | 99.71 | 85.78 |
| FCN | 88.34 | 88.68 | 79.45 | 88.43 | 99.84 | 99.83 | 99.67 | 99.83 | 92.33 |
| SegNet | 85.72 | 85.81 | 75.13 | 85.73 | 99.79 | 99.80 | 99.59 | 99.79 | 90.60 |
| TDSSNet | 86.55 | 81.87 | 72.61 | 84.06 | 99.72 | 99.80 | 99.53 | 99.76 | 89.74 |
| **MD-Unet** | **92.75** | **90.94** | **84.93** | **91.81** | **99.87** | **99.89** | **99.76** | **99.88** | **94.67** |

diseased areas, and overall leaf segmentation [8].Out of all the Baseline and Advanced models, the MDU-Net architecture achieved the greatest metrics across most metrics, with a Dice score of 94.67% indicating relatively strong reliability on the quality of the segmentation results. MDU-Net recorded the highest CPA of 92.75%, with the highest Recall of 90.94%, and the highest IoU of 84.93% in the lesion segmentation task, demonstrating that it accurately captured much of the complexity of the lesion shape. Importantly however, despite the noted differences between different models in lesion segmentation , the performance in the leaf segmentation task was consistently high as well, with all models recording performance nearing perfect, demonstrating high robustness and generalization across variances in structural aspects of the dataset. Overall, conventional models such as U-Net and FCN provided still relatively strong an competitive performance and continue to serve as highly reliable benchmarks for medical and agricultural imaging applications.More complex variants such as UNET++ or PSPNet address the potential issues of complexity and performance; however, DeepLab v3+ receives the lowest score in terms of lesion-specific metrics. This likely speaks to the structured way in which DeepLab v3+ applies receptive fields or fuses features in combination with isolated geometry. This is possibly indicative of it being inappropriate for detailed, small-scale geometry such as the structures of tobacco leaf lesions. In conclusion, the results provide strong evidence for MD-Unet's basic capacity to be able to segment and label high precision segmentation applications. Its

enhancements to the architecture are presumably favourable to evaluate textures of subtle lesions and poorly defined boundaries which are often difficult for traditional models to comprehend. The performance of MD-Unet is especially noteworthy given the complex nature of standard tobacco leaves that were representative of natural agricultural contexts. This ongoing, consistent high performance is persuasive not only for robustness, but also for flexibility to work in real agricultural situations. The results provide confidence that MD-Unet could be a trustworthy and scalable method of automated or intelligent identification of diseased tobacco plants and in the development of intelligent smart phytoprotection systems.

## VI. COMPARISON OF SURVEY ARTICLES

In recent years, there has been a surge in the number of survey articles on DL and ML approaches to plant disease recognition, particularly in high value crops such as tobacco. These survey articles collectively show a shift from doing manual diagnosis of disease to producing intelligent automated methods via convolutional neural networks (CNNs), generative adversarial networks (GANs), attention mechanisms, and lightweight models. The survey literature gathered from many relatively recent publications provides a wealth of similarities, differences in method and trends in the evolution of dataset types and real world use.

A core strength that is seen throughout these survey articles, is that they explore segmentation-based architectures, including those using U-Net variants, and covering a significant range. For example, the MD-UNet proposed by Chen et al.

(2020) [8] is an excellent demonstration of how utilizing multiscale residual modules with attention systems improves the accuracy of tobacco leaf spot segmentation over and above aligned architectures of traditional U-Net and DeepLabV3+ as a standard to benchmark improvement, as accessible as a few improvement criteria on loss functions. Multiple of the surveys reviewed in this merged documents provide examples and avenues of improvements linked to these kinds of architectural improvements through assessments and comparisons using public datasets such as the Luoyang and TPDD datasets. The surveys make no mistake to point out segmentation tasks - not only classification tasks - are taking center stage and that tasks generating or identifying segmentation segments of disease will be the key for agriculture in precision agriculture referencing localization, quantification has become their focus.

Other papers researched class-based pipelines using CNNs such as ResNet50, EfficientNet, and MobileNetV2, which are primarily fine-tuned using transfer learning. This is especially relevant in circumstances with limited training data. The authors Ou et al., [42] were unequivocal about the small volume and diversity of tobacco-based datasets, typically requiring researchers to leverage pre-trained weights from sizable repositories of image classification such as ImageNet or its variations. The literature reviews further examined the limits of these learning approaches to domain-specific tasks, particularly the need for fine-grained annotations and field-based images.

Conventional papers, such as TRNet [36], introduce a hybrid approach to DL modelling, where convolutional base feature extraction is combined with transformer attention modules. The surveys with TRNet identified great capabilities to recognize local texture, along with global structural features, which are important for distinguishing visually similar but pathologically distinct tobacco leaf symptoms. The authors also compared inclusion of TRNet against contemporary CNNs, especially in reference to managing complex backgrounds and intra-class variability. Likewise, the same can be said for OR-AC-GAN published by Lin et al. [17], which uses adversarial augmentation to improve model robustness, a strategy often mentioned in review papers as a viable route for real-world implementation. An important dimension across surveys is dataset diversity and generalizability. The Prakasam and Luoyang datasets are frequently used in benchmarking tasks but have limitations in geographic and phenotypic diversity. Articles such as Ying et al. [30] highlight the inclusion of NIR spectral channels to augment traditional RGB data, offering early disease detection capabilities. This multimodal fusion is appreciated in newer surveys as it enhances model robustness against occlusions, variable lighting, and phenological shifts. However, this is also where survey articles diverge—while some advocate for expensive spectral imaging systems, others suggest cost-effective RGB + depth solutions for broader applicability in low-resource settings.

Performance metrics are another cornerstone in the comparative analysis found across articles. Accuracy, F1-score, sensitivity, specificity, and area under the ROC curve (AUC) are widely used to benchmark models. According to the review by Peng et al. [40], LIBS-based image analysis methods can achieve 97.2% accuracy using Partial Least Squares Discriminant Analysis (PLS-DA), and up to 97.6% with Support Vector Machines (SVMs). While these numbers are impressive, most surveys note that such metrics are often inflated due to overfitting on small, homogeneous datasets. Thus, many survey articles stress the need for cross-validation, field trials, and inclusion of hard negatives in the training process to better simulate deployment conditions.

Another popular trend is explainable AI (XAI) in modern surveys. Traditional models function as black boxes, which are problematic when deploying them agronomically, as domain experts require the outputs to be interpretable. We see several articles recommend technologies such as Grad-CAM, SHAP, and LIME for the visualization of the parts of the leaf that most influence classification or segmentation decisions. The articles also expand on the theme of transparency and the fact that not only do people require accountability and trust in AI, but transparency is also important for regulatory acceptance, which is critical in industries such as crop insurance and quality grading.

From a deployment perspective, lightweight models and embedded inference are recurring concerns. Articles such as Lin et al. [35] propose optimized YOLO-based architectures for use in edge computing scenarios such as UAVs or mobile phones. Surveys that evaluate these works generally agree that performance trade-offs are necessary; while lightweight models may sacrifice some accuracy, their benefits in latency, battery efficiency, and portability are vital for scaling AI-based disease detection to smallholder farms.

A separate class of survey papers also provides evaluations of boosting approaches to augment dataset. Synthetic images generated using GANs have been widely adopted to solve class imbalance. For example, the papers reviewed in TRNet2 [36] show that training on both real and GAN augmented images allows the model to generalize better across seasons and geographies. However, many survey papers caution that care must be taken to avoid artifacts and mode collapse in generated samples which may confuse classifiers.

The challenge of domain shift, i.e. models trained in one region still underperform in another region, is mentioned in almost all survey papers as one of the challenges. Some possible terms of providing solutions as next-generation methods include domain adaptation, few-shot learning, and FL . Federated leaning is being more widely cited in the literature because it preserves privacy and encourages decentralized trained models to collaborate between different farms and institutions using their data because they do not need to share their farm data which often is confidential.A survey paper mentioned the utility of self supervised learning (SSL) with large unlabeled datasets [48]. There are methods such as contrastive learning and image reconstruction which allow for pretraining without the manual annotations. Given that obtaining labeled data can be especially challenging for crops such as tobacco, this can be especially useful. There is a

growing recognition of SSL as a viable pretraining method for the segmentation models that are then fine-tuned with limited labeled dataset.

Surveys also mention the environmental and seasonal effects that players in the literature thought may affect model performance. Some articles mention that leaf symptoms can can vary based on disease, as well as cultivar, type of soil, moisture, and light intensity. A few survey articles mentioned that models trained on leaf-only images may lead to poor performance if they are deployed with only canopy- or plant-level impairments, thus recommending multi-scale feature extraction or object detection pipelines.

By and large, these survey papers advocate for large, collaborative datasets, in the form of open-source libraries, supplemented with adequate annotations. Futures direction will trend similarly to the creation of benchmark challenges and leaderboards that enable rapidity in progress founded on documented and standardized evaluations.

To summarize the literature that we have reviewed in both merged documents, it has been evident that there has been quick and progressive evolution of techniques and preferences in regard to tobacco leaf disease detection. We have moved from simplistic CNN classifiers, and even hybrid CNN models, toward more advanced architectures like transformer-augmented segmentation models and GAN-based segmentation models. This is an obvious indication that the field is moving toward applicability in the real world.

Still, there are considerable obstacles going forward: the main ones being a lack of datasets, lack of standardization, and performance vs. deployment. It is encouraging to see promising applications of multimodal sensing, explainable AI, FL, and self-supervised pretraining. These survey papers, to our knowledge, have engaged and documented contemporary state-of-the-art motion. At the same time, they need to inspire growth and act as a road map for efficient, resource and time favourable disease recognition systems for agriculture that are scalable and explainable.

## VII. LIMITATIONS AND FUTURE WORK

Despite significant advancements in tobacco leaf disease detection using DL and ML methods, several critical limitations persist that impact the real-world effectiveness of these systems [49], [50].

### A. Limited Annotation Context

Current datasets such as the common TLA dataset, primarily consist of front-facing images of tobacco leaves used for expert annotation and — as stated by the authors — reliable disease identification often requires additional context, such as images of the undersides of leaves, the stem condition, or even the location on the plant (e.g. apical or basal). In the absence of this information, there may be interpretative and ambiguous annotation (which can contribute to model misclassification), affecting the ability of the model to perform. Future datasets can benefit by multi-view images, even with metadata for gathering additional information to improve quality of annotations.

### B. Insufficient Growth Stage and Disease Phase Coverage

A majority of existing datasets are composed of samples taken when the tobacco plant is in a mature state, with little or none of the earlier growth stages. In addition, there are no annotations made for disease progression, which is crucial for early detection of diseases. Many diseases exhibit overlapping visual symptoms during their early progression stage and only become distinguishable as they progress. If models were able to consider temporal progression and annotations for the lifecycle stage, then they would be able to identify disease/infection at earlier stages.

### C. Multidisease Complexity and Single-Label Limitations

In-field conditions, it is not uncommon for one leaf to show signs of infection of more than one disease at the same time. Yet, most classification models will assume single-label conditions and thus, oversimplify the problem and fail to represent the actual world. Even when assuming the classification is multi-class, the combinatorial explosion of potential disease combinations renders traditional multi-class classification unscalable. Future work should focus on practical approaches that include multi-label classification methods, and instance-level segmentation, to address the issue of co-infections.

### D. Environmental Variability and Lighting Artifacts

Field-collected imagery often suffers from contextual variability in both ambient illumination and brightness especially in strong sunlight; thus obscuring important lesion characteristics and adding noise to the dataset. Basic techniques such as physical shading were used during data collection, but there were no post-process approaches to normalize light levels. If preprocessing approaches such as contrast correction, histogram equalization, or exposure normalization could be employed in these circumstances, models trained on less variable outdoor imagery could become less variable overall and more robust (e.g., reducing variation stemming from differing outdoor conditions by ensuring uniformity during preprocessing).

It would be vitally important to address these limitations before we could hope to build robust, generalizable and deployable tobacco disease detection systems that could be used reliably in a range of managed agricultural contexts.

## VIII. CONCLUSION

In this paper, we provided a thorough review of major developments in the area of tobacco leaf pathology using ML and DL, focusing specifically on image segmentation. In the past ten years, the field has advanced from primitive SVM classifiers based on hand-crafted features to sophisticated attention-based CNNs that identify lesions down to the pixel level. We investigated a wide array of methods including U-Net variants, transformer-based networks, GANs for augmentation, and lightweight mobile approaches. Through a thorough comparative analysis we demonstrated the effectiveness of using attention models, data augmentation, and hybrid models to help mitigate challenges such as class imbalance, domain

shift, and low data. Despite advancements, there are still gaps in the research. Real world deployment has issues with insufficient availability of efficient, optimized lightweight models. There are issues with inconsistency of labeling protocols, and a lack of large, full-featured authentic datasets that limits generalization across datasets. Future work should prioritize increased interpretability of models, models trained on new paradigms for self-supervised and FL, and the incorporation of multimodal data to increase robustness and performance. In the end, the incorporation of DL for detecting diseases affecting tobacco and other crops has marked potential to help improve the monitoring of crops, decrease monetary losses, and help advance sustainable agriculture. The potential and advancements enabled by DL ensure progress in the agritech field will not slow down.

## ACKNOWLEDGMENT

## REFERENCES

[1] X. Zou, A. Bk, T. Abu-Izneid, A. Aziz, P. Devnath, A. Rauf, S. Mitra, T. Bin Emran, A. A. H. Mujawah, J. M. Lorenzo, M. S. Mubarak, P. Wilairatana, and H. A. R. Suleria, "Current advances of functional phytochemicals in *nicotiana* plant and related potential value of tobacco processing waste: A review," *Biomedicine & Pharmacotherapy*, vol. 143, p. 112191, Nov. 2021, epub 2021 Sep 22. [Online]. Available: https://doi.org/10.1016/j.biopha.2021.112191

[2] Export Import Data, "Tobacco export from india – major export destinations and stats," 2024. [Online]. Available: https://www.exportimportdata.in/blogs/tobacco-export-from-india.aspx

[3] World Bank, FAO, and ITC Trade Map, "Tobacco production, employment, and export statistics," https://data.worldbank.org, https://www.fao.org, https://www.intracen.org, 2024, accessed: July 2024.

[4] S. Naik and A. K. Singh, "Artificial intelligence in tobacco crop disease detection," in *Artificial Intelligence in Precision Agriculture*, S. Yadav, D. Pant, M. Kumar, and D. Singh, Eds. Academic Press, 2023, pp. 153–161. [Online]. Available: https://doi.org/10.1016/B978-0-323-90899-3.00090-2

[5] J. Coursen. (2019, Mar.) Gene-editing technology harnessed to protect plants from viruses. National Human Genome Research Institute. Accessed: 2025-07-18. [Online]. Available: https://www.genome.gov/news/news-release/Gene-editing-technology-harnessed-to-protect-plants-from-viruses

[6] M. Zaitlin and P. Palukaitis, "Tobacco mosaic virus," in *Encyclopedia of Virology*, R. G. Webster and A. Granoff, Eds. Academic Press, 1999, pp. 1819–1826. [Online]. Available: https://doi.org/10.1006/rwvi.1999.0282

[7] S.-K. Tang and Y. Tang, "Tobacco leaf diseases detection and classification using deep learning," in *Proceedings of SPIE 12446, Tenth International Conference on Graphic and Image Processing (ICGIP 2022)*, vol. 12446. SPIE, 2023, p. 124462F. [Online]. Available: https://www.spiedigitallibrary.org/conference-proceedings-of-spie/12446/2644288/Tobacco-leaf-diseases-detection-and-classification-using-deep-learning/10.1117/12.2644288.full

[8] Z. Chen *et al.*, "Md-unet for tobacco leaf disease spot segmentation," *Scientific Reports*, vol. 15, no. 2759, 2025.

[9] B. Sambana, H. S. Nnadi, M. A. Wajid *et al.*, "An efficient plant disease detection using transfer learning approach," *Scientific Reports*, vol. 15, p. 19082, 2025. [Online]. Available: https://doi.org/10.1038/s41598-025-02271-w

[10] A. G and A. H, "Application of machine learning in measurement systems," *Measurement*, vol. 193, p. 112239, 2022.

[11] Press Information Bureau, "India is the 2nd largest producer and exporter of tobacco in the world," https://www.pib.gov.in/PressReleasePage.aspx?PRID=2089182, 2024, accessed: 2025-07-19.

[12] (n.d.) Transfer learning — ai glossary — opentrain ai. OpenTrain AI. Accessed: 2025-07-18. [Online]. Available: https://www.opentrain.ai/glossary/transfer-learning

[13] Q. Wang, G. Liu, Y. Wu, and L. Ma, "Unmanned aerial vehicle remote sensing for precision agriculture: A survey," *Computers and Electronics in Agriculture*, vol. 139, pp. 67–84, 2017.

[14] S. N. Saw, Y. Y. Yan, and K. H. Ng, "Current status and future directions of explainable artificial intelligence in medical imaging," *European Journal of Radiology*, vol. 183, p. 111884, 2025. [Online]. Available: https://doi.org/10.1016/j.ejrad.2024.111884

[15] P. Karmakar, S. W. Teng, M. Murshed, S. Pang, Y. Li, and H. Lin, "Crop monitoring by multimodal remote sensing: A review," *Remote Sensing Applications: Society and Environment*, vol. 33, p. 101093, 2024. [Online]. Available: https://doi.org/10.1016/j.rsase.2023.101093

[16] G. P. Babu and A. N. Reddy, "Tobacco leaf disease detection using cnn on prakasam district data," *International Journal of AgriTech*, vol. 13, no. 2, pp. 102–108, 2024.

[17] W. Lin and F. Huang, "Tpdd: A real-world dataset for tobacco disease detection," *Plant Image Data Journal*, vol. 5, no. 2, pp. 60–72, 2021.

[18] N. Lecours, G. E. G. Almeida, J. M. Abdallah, and T. E. Novotny, "Environmental health impacts of tobacco farming: a review of the literature," *Tobacco Control*, vol. 21, no. 2, pp. 191–196, 2012. [Online]. Available: https://tobaccocontrol.bmj.com/content/21/2/191

[19] M. L. Dissanayake, B. C. Jayawardana, and N.-Y. Lee, "Microbiota in tobacco and tobacco products: A review," *Frontiers in Microbiology*, vol. 11, p. 572725, 2020.

[20] Y. Li, S. Duan, X. Sun, L. Zhang, D. Yu, R. Wang, L. Zhao, H. Liu, Z. Wang, X. Yang, and Y. Li, "Microbial composition and diversity of the tobacco leaf phyllosphere during plant development," *Frontiers in Microbiology*, vol. 13, p. 845310, 2022. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fmicb.2022.845310/full

[21] APS Education Center, "Introduction to abiotic disorders in plants," https://www.apsnet.org/edcenter/Pages/Abiotic.aspx, accessed: 2025-07-20.

[22] W. Hu, J. Yuan, J. Fei, K. Imdad, P. Yang, S. Huang, D. Mao, and J. Yang, "Shaping the future of tobacco through microbial insights: A review of advances and applications," *Frontiers in Bioengineering and Biotechnology*, vol. 12, 2024.

[23] V. K. Vishnoi, K. Kumar, and B. Kumar, *Plant Disease Detection Using Computational Intelligence and Image Processing*. Berlin Heidelberg: Springer, 2021.

[24] F. Li *et al.*, "Microbial interactions and metabolisms in response to bacterial wilt and black shank pathogens in the tobacco rhizosphere," *Frontiers in Plant Science*, vol. 14, p. 1200136, 2023.

[25] V. Petrov, J. Hille, B. Mueller-Roeber, and T. S. Gechev, "Ros-mediated abiotic stress-induced programmed cell death in plants," *Frontiers in Plant Science*, vol. 6, p. 69, 2015. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fpls.2015.00069/full

[26] N. Suzuki, R. M. Rivero, V. Shulaev, E. Blumwald, and R. Mittler, "Abiotic and biotic stress combinations," *New Phytologist*, vol. 203, no. 1, pp. 32–43, 2014.

[27] J. G. A. Barbedo, "Impact of dataset size and variety on the effectiveness of deep learning and transfer learning for plant disease classification," *Computers and Electronics in Agriculture*, vol. 153, pp. 46–53, 2018.

[28] K. N. R. K. S. A. K. Viswanatha Reddy, V. Paramesha, "Econometric modeling of tobacco exports in the milieu of changing global and national policy regimes: repercussions on the indian tobacco sector," *Frontiers in Environmental Economics*, vol. –, 2023.

[29] X. Li, W. Zhang, S. Zhou, J. Zhu, and Y. Li, "Precision agriculture technologies positively contributing to ghg emissions mitigation, farm productivity and economics," *Computers and Electronics in Agriculture*, vol. 127, pp. 207–216, 2016.

[30] L. Ying and Z. Chen, "Deep learning for tobacco disease using near-infrared spectroscopy," *Plant Spectral Imaging*, vol. 6, no. 1, pp. 34–41, 2024.

[31] M. Xu, W. Zhang, N. Li, X. Huang, and L. Wang, "Plant disease recognition datasets in the age of deep learning: A review," *Frontiers in Plant Science*, vol. 15, p. 1345805, 2024. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fpls.2024.1345805/full

[32] A. Muhammad, Z. Salman, K. Lee, and D. Han, "Harnessing the power of diffusion models for plant disease image augmentation," *Frontiers in Plant Science*, vol. 14, p. 1280496, 2023. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fpls.2023.1280496/full

[33] J. Dong, J. Lee, A. Fuentes, M. Xu, S. Yoon, M. H. Lee, and D. S. Park, "Data-centric annotation analysis for plant disease detection: Strategy, consistency, and performance," *Frontiers in Plant Science*, vol. 13, p. 1037655, 2022. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fpls.2022.1037655/full

[34] Z. Chen *et al.*, "A few-shot learning method for tobacco abnormality identification," *Frontiers in Plant Science*, vol. 15, March 2024. [Online]. Available: https://doi.org/10.3389/fpls.2024.1333236

[35] Q. Lin and M. Zhang, "Yolo-tobacco: Lightweight brown spot detector using cbam," in *2023 IEEE Conference on Smart Farming*, 2023, pp. 56–62.

[36] J. Yao and L. Deng, "Trnet: Transformer-aided convolutional network for leaf disease segmentation," *IEEE Access*, vol. 11, pp. 33 321–33 330, 2023.

[37] S. Terzi, "Effect of different weight initialization strategies on transfer learning for plant disease detection," *Plant Pathology*, 2024.

[38] J. Teng and H. Zhao, "Tobacco brown spot classification using binary tree svm," *Precision Agriculture*, vol. 8, no. 1, pp. 88–95, 2021.

[39] e. a. He, "Asnet: Asymmetric shuffle cnn with scse attention," *Journal of Intelligent Agriculture*, 2022.

[40] J. Peng, K. Song, H. Zhu, W. Kong, F. Liu, T. Shen, and Y. He, "Fast detection of tobacco mosaic virus infected tobacco using laser-induced breakdown spectroscopy," *Scientific Reports*, vol. 7, p. 44551, 2017.

[41] M. Zeybek and Şanlıoğlu, "Geometric feature extraction of road from uav based point cloud data," *Measurement*, vol. 133, pp. 99–111, 2019.

[42] L. Ou and W. Wang, "Tdssnet: Attention-based segmentation model for tobacco leaf disease," *Computational Agriculture*, vol. 12, no. 2, pp. 122–130, 2024.

[43] e. a. Yuan, "Improved deeplabv3+ network for grape leaf disease segmentation," *Computers and Electronics in Agriculture*, 2023.

[44] e. a. Zhao, "Inception u-net for tomato disease segmentation," *Agricultural Intelligence*, 2023.

[45] D. Sindhu and R. Venkatesh, "A deep belief network-based framework for tobacco mosaic virus detection," *Applied Intelligence in Agriculture*, vol. 9, no. 3, pp. 150–159, 2024.

[46] e. a. Fu, "Rs-unet: Residual attention u-net for potato leaf disease," *Plant Methods*, 2022.

[47] J. Johnson, "Survey on deep learning with class imbalance," *Journal of Big Data*, vol. 6, pp. 1–20, 2019.

[48] Z. F. Srinivasa Rao Nandam, Sara Atito, "Investigating self-supervised methods for label-efficient learning," *International Journal of Computer Vision*, vol. 133, pp. 4522–4537, 2025. [Online]. Available: https://doi.org/10.1007/s11263-025-02397-4

[49] Z. Zhou, F. Liu, W. Wang, Y. Yang, P. Li, and X. Zhang, "Lightweight yolov5 model for real-time detection of tobacco diseases in the field," *Frontiers in Plant Science*, vol. 15, p. 1333236, 2024, original Research Article, Section: Sustainable and Intelligent Phytoprotection. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fpls.2024.1333236/full

[50] A. Bronselaer, J. Nielandt, T. Boeckling, and G. De Tré, "Operational measurement of data quality," in *Information Processing and Management of Uncertainty in Knowledge-Based Systems: Applications. IPMU 2018*, ser. Communications in Computer and Information Science, vol. 855, 2018.