

Journal Pre-proof

End-to-end plant disease detection using
transformers with collaborative hybrid
assignment trainings

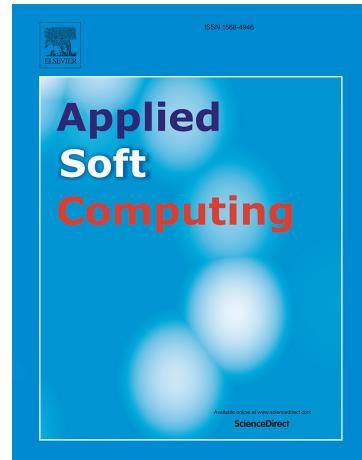
Yanfen Li, Muhammad Fayaz,
Sufyan Danish, Lilia Tightiz,
Hanxiang Wang, Tan N. Nguyen,
L.Minh Dang

PII: S1568-
4946(25)01450-4
DOI: [https://doi.org/10.1016/
j.asoc.2025.114137](https://doi.org/10.1016/j.asoc.2025.114137)
Reference: ASOC 114137

To appear in: *Applied Soft
Computing*

Received Date: 7 May 2025
Revised Date: 3 September 2025
Accepted Date: 27 October 2025

Please cite this article as: Li Y, Fayaz M, Danish S, Tightiz L,
Wang H, Nguyen TN, Dang LM, End-to-end plant disease detection
using transformers with collaborative hybrid assignment trainings,
Applied Soft Computing (2025),
doi: <https://doi.org/10.1016/j.asoc.2025.114137>.



This is a PDF of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability. This version will undergo additional copyediting, typesetting and review before it is published in its final form. As such, this version is no longer the Accepted Manuscript, but it is not yet the definitive Version of Record; we are providing this early version to give early visibility of the article. Please note that Elsevier's sharing policy for the Published Journal Article applies to this version, see: <https://www.elsevier.com/about/policies-and->

[standards/sharing#4-published-journal-article](#). Please also note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2025 Published by Elsevier B.V.

Journal Pre-proof

Highlights

- A huge fruit disease dataset of 6 different diseases that contains over 81,000 images.
- An efficient transformer-based fruit disease detection framework.
- Analysis of the disease region using the transformer's feature discriminability scores.
- The proposed model outperformed previous state-of-the-art object detection models.

End-to-End Plant Disease Detection Using Transformers with Collaborative Hybrid Assignment Trainings

Yanfen Li^a, Muhammad Fayaz^b, Sufyan Danish^b, Lilia Tightiz^c, Hanxiang Wang^a, Tan N. Nguyen^{d,*}, L. Minh Dang^{e,f,*}

^a*School of Computer Science, Qufu Normal University, Rizhao, 276826, China;*

^b*Department of Computer Science and Engineering, Sejong University, Seoul 05006, South Korea;*

^c*School of Computing, Gachon University, 1342 Seongnamdaero, Seongnam-si, Gyeonggi-do, 13120, Republic of Korea;*

^d*Department of Architectural Engineering, Sejong University, 209 Neungdong-ro, Gwangjin-gu, Seoul 05006, Republic of Korea;*

^e*The Institute of Research and Development, Duy Tan University, Da Nang, 550000, Viet Nam;*

^f*Faculty of Information Technology, Duy Tan University, Da Nang, 550000, Viet Nam*

Abstract

Plant diseases pose a significant threat to fruit production and quality if not detected and managed promptly. Precise and efficient recognition of these diseases is critical for ensuring plant health and maximizing fruit production. To tackle this issue, a range of image processing and deep learning techniques have been preferred for plant disease recognition due to their superior performance. This paper proposes an end-to-end transformer-based model that improves both the accuracy and detection rate of fruit diseases. The model is based on a state-of-the-art transformer model and trained using the Collaborative Hybrid Assignment (Co-DETR) scheme. Moreover, several targeted modifications to the original model are conducted to optimize its performance. These modifications enable the model to detect six types of plant diseases with a mean average precision (mAP) of 0.89 while maintaining efficient training times. The proposed model consistently outperforms state-of-the-art detection models. In addition, the model offers interpretability through the visualization of feature

*Corresponding authors

Email addresses: tnnguyen@sejong.ac.kr (Tan N. Nguyen), danglienminh@duytan.edu.vn (L. Minh Dang)

discriminability scores to ensure that the prediction process is interpretable and understandable. Finally, the model demonstrates robust performance under challenging environmental conditions, such as poor lighting and image blurring, which is essential for real-world applications in disease management and precision agriculture.

Keywords: image processing, transformer, deep learning, precision agriculture, fruit disease

1 1. Introduction

2 According to the Food and Agriculture Organization (FAO), global food
3 demand is projected to surge by 70% by 2050 as the world population surpasses
4 9.1 billion [1]. Fruits, as critical sources of essential nutrients, play a pivotal
5 role in ensuring food security and combating malnutrition [2]. However, it is
6 increasingly challenging to achieve and sustain high fruit yields due to factors
7 such as limited farmland, climate change, and the devastating impact of pests
8 and diseases [3]. Among these threats, fruit diseases, such as mango scab and
9 citrus thrips, often cause catastrophic yield losses and economic devastation
10 when left undetected.

11 Traditional fruit disease detection relied on manual inspection, a labor-
12 intensive and error-prone process with delays in identifying early-stage symp-
13 toms [4]. The lag between symptom appearance and detection often results
14 in significant losses. To address these challenges, automated detection systems
15 using machine learning (ML) and deep learning (DL) have emerged as transfor-
16 mative solutions for scalable, accurate, and efficient disease monitoring [5].

17 While early ML approaches utilized handcrafted features, such as color, tex-
18 ture, and shape, with classifiers like support vector machines (SVMs) [6] and
19 random forests (RFs) [7], their performance was constrained by domain-specific
20 feature engineering and environmental variability [3]. Recent DL advancements,
21 particularly convolutional neural networks (CNNs), have demonstrated superior
22 performance in disease classification [8], segmentation [9], and detection [10, 11].

23 However, CNN-based models often require manual hyperparameter tuning, like
24 anchors, proposals, and post-processing to reduce redundant predictions [12].

25 Transformers were initially developed for natural language processing (NLP).
26 Their self-attention mechanisms [13] enable global context modeling, which ad-
27 dresses CNN limitations in capturing long-range dependencies [12]. For in-
28 stance, Longformer [14] introduced sliding window attention to process long
29 documents efficiently. Reformer [15] reduced computational complexity using
30 locality-sensitive hashing for large-scale NLP tasks. Beyond NLP, the adapt-
31 ability of transformers was further enhanced by specialized variants customized
32 to domain-specific challenges. For example, in finance, transformer variants
33 have been trained to model temporal patterns and forecast price movements
34 [16, 17], while in remote sensing and fault detection, they have enabled pre-
35 cise anomaly identification in high-resolution imagery and industrial systems
36 [18]. In manufacturing, transformers powered quality inspection and predictive
37 maintenance [19]. In protein sequence modeling, Performer with kernel-based
38 attention was introduced to effectively model the scalable protein sequence [20].
39 However, their application to fruit disease detection remains underexplored,
40 with challenges in convergence, data scarcity, and subtle symptom recognition
41 in complex agricultural environments [21].

42 To bridge this gap, this study introduces FD-TR, a modified transformer-
43 based fruit disease detection model based on Co-DETR training scheme [22].
44 The key contributions of this study are:

- 45 • The proposed model was trained on a large-scale fruit disease dataset of
46 81,000 high-resolution images.
- 47 • Key modules (e.g., loss, optimizer) of Co-DETR were replaced, and hy-
48 perparameters were fine-tuned to address the unique challenges of fruit
49 disease detection.
- 50 • Introduction of a feature discriminability score visualization method to
51 enhance model interpretability for real-world deployment.

52 • The model demonstrated its robustness through systematic evaluation
53 across four benchmark datasets and an additional healthy fruit subset.

54 The outline of the manuscript is as follows. Section 3 provides a comprehensive
55 description of the fruit disease dataset used in this study. Section 4
56 discusses in detail each component of the proposed fruit disease detection frame-
57 work based on DINO with a Co-DETR training scheme. The results of various
58 experiments conducted to evaluate the model’s performance are reported in Sec-
59 tion 5. Section 6 discusses the main contributions and experimental results of
60 this study. Finally, Section 7 provides conclusions and outlines future research
61 directions.

62 2. Related Work

63 Table 1 provides an overview of recent fruit disease detection studies. It
64 highlights the diversity of models used, ranging from CNN models to hybrid
65 and transformer-based architectures, applied on various fruit types. While most
66 models achieved high accuracy on their respective datasets, the majority were
67 limited by small sample sizes, limited disease coverage, and lack of real-world
68 deployment validation. These limitations emphasized the need for a more gen-
69 eralized, scalable, and interpretable solution.

Table 1: Summary of recent fruit disease detection studies (2020–2025)

Author(s) & Year	Method	Dataset	Main findings	Limitations
Xie et al. (2020) [23]	Inception + SE-block	Grape leaf images (4449 images)	0.81 mAP at 15 FPS	Computationally intensive architecture
You et al. (2022) [24]	YOLO + Deep Metric Learning	Strawberry dataset (7230 images)	97.8% overall accuracy	Complex architecture; lab-based dataset
Syed et al. (2022) [25]	Two-stage CNN	Citrus leaf images (598 images)	94.37% accuracy	Limited generalizability
Huang et al. (2023) [26]	EfficientNet-Inception CNN + U-Net	Citrus dataset (800 images)	95.6% classification accuracy; 87.7% severity segmentation	Only 2 citrus diseases; small, lab-based dataset
Arifin et al. (2024) [27]	ResNet50 features + Logistic Regression	Citrus dataset (1814 images)	99.69% accuracy	Small, imbalanced dataset; no lesion localization
Sun et al. (2024) [11]	YOLOv5 + shuffle-channel blocks	Natural orchard images (4252 images)	0.93 mAP	Manual hyperparameter tuning and post-processing needed
Aksoy et al. (2025) [28]	ResNet152V2 (transfer learning)	Kaggle apple fruit disease images (502 images)	92% classification accuracy	Small dataset (4 classes)
Faye et al. (2025) [29]	ResNet50 for severity grading	SenMangoFruitDDS (862 images)	97.8% accuracy	Only on mango; limited background variability
He et al. (2025) [30]	Sparse Attention YOLOv11	Passion fruit dataset (10,000 annotated images)	90% F1-score	Only passion fruit; stem-focused labels; high computational cost

70 *2.1. Traditional Machine Learning Approaches*

71 Early efforts in fruit disease detection focused on ML models using hand-
72 crafted features. For instance, SVMs trained on color and texture features
73 achieved moderate success in classifying diseases on fruits [6]. RFs were em-
74 ployed to distinguish apple fruit diseases based on color and texture descriptors
75 [7]. However, these methods struggled with environmental variability and re-
76 quired extensive domain expertise for feature design [3].

77 *2.2. Deep Learning-Based Approaches*

78 The developments of CNNs revolutionized fruit disease detection. One-stage
79 detectors like You Only Look Once (YOLO) and Single Shot MultiBox Detector
80 (SSD), and two-stage frameworks such as Region-based CNN (R-CNN), were
81 progressively adopted for precise recognition of fruit diseases [3]. For example,
82 Sun et al. [11] introduced an innovative method for identifying fruit diseases in
83 natural orchard settings using a combination of binocular cameras and DL tech-
84 niques. They implemented a Unimatch stereo-matching algorithm to generate
85 depth maps that focused detection on leaves and proposed a lightweight disease
86 detection model based on YOLOv5-augmented with shuffle-channel blocks and
87 attention modules. The experimental results reveal that it outperformed the
88 YOLOv5-s architecture with 0.93 mean average precision (mAP). Syed et al.
89 [25] presented a two-stage CNN for citrus disease detection. Firstly, the model
90 employed a region proposal network to identify potential diseased areas on citrus
91 leaves. After that, it classified these regions into specific disease categories using
92 a classifier. The model demonstrated a high detection accuracy of 94.37% for
93 citrus black spot, citrus bacterial canker, and Huanglongbing. In another study,
94 Xie et al. [23] addressed real-time detection of common grape leaf diseases us-
95 ing a customized Faster R-CNN with Inception-v1, Inception-ResNet0-v2, and
96 SE-block. The model achieved a mAP of 0.81 at a real-time detection speed of
97 15.01 frames per second (FPS). Although these DL models enabled early and
98 accurate disease detection, they still required manually fine-tuned hyperparam-
99 eters like anchors and proposals during training and additional post-processing

¹⁰⁰ algorithms to reduce duplicate predictions [12].

¹⁰¹ *2.3. Transformer-Based Approaches*

¹⁰² Transformers have introduced paradigm shifts in object detection. Vision
¹⁰³ Transformers (ViTs) effectively processed entire images as sequences of patches,
¹⁰⁴ which enhanced global context modeling and motivated researchers to extend
¹⁰⁵ their use to more complex tasks such as object detection [31]. For example,
¹⁰⁶ Carion et al.[32] proposed detection transformer (DETR), an end-to-end object
¹⁰⁷ detector that directly predicted bounding boxes (BB) and classes via learned ob-
¹⁰⁸ ject queries. DETR did not require extensive manual tuning and was proved to
¹⁰⁹ handle varying object sizes and overlapping objects more effectively. Subsequent
¹¹⁰ extensions, such as Deformable DETR [33], DN-DETR [34], and DAB-DETR
¹¹¹ [35], aimed to improve DETR’s convergence and performance. While these ex-
¹¹² tensions showed better detection performances, they still performed poorer than
¹¹³ the CNN counterparts [12]. The recent introduction of a collaborative hybrid
¹¹⁴ assignments training scheme for DETR (Co-DETR) [22] addressed the issue of
¹¹⁵ sparse supervision in DETR models by utilizing multiple auxiliary heads with
¹¹⁶ one-to-many label assignments to enhance the learning of both the encoder and
¹¹⁷ decoder. Co-DETR improved the training efficiency and discriminative feature
¹¹⁸ learning of DETR-based detectors without adding any extra computational cost
¹¹⁹ or parameters during inference. The experiment results demonstrated a signifi-
¹²⁰ cant performance gain on various DETR variants. The integration of Co-DETR
¹²¹ into DINO-Deformable-DETR achieved 66.0% AP on the Common Objects in
¹²² Context (COCO) test-development set.

¹²³ **3. Materials**

¹²⁴ Table 2 highlights the evolution of benchmark datasets in plant disease re-
¹²⁵ search. Earlier datasets, such as PlantDoc [36] and PlantVillage [37], included
¹²⁶ diseases affecting both fruits and leaves on multiple species but did not specifi-
¹²⁷ cally focus on fruit diseases. In contrast, smaller self-collected datasets, such as

¹²⁸ the Pomegranate Fruit Diseases [38] and Citrus Diseases [39], primarily focus on
¹²⁹ diseases of single fruit types and contain fewer than 3,000 images, which limit
¹³⁰ their scalability and generalizability.

¹³¹ This research stands out by training the proposed model on a large fruit
¹³² disease identification dataset containing roughly 81,000 images that cover six
¹³³ different fruit disease types [40]. Provided by the National Information Society
¹³⁴ Agency of Korea (NIA)¹, this extensive dataset exceeds the scope and size of
¹³⁵ most existing datasets. The collection of data was made possible through the
¹³⁶ collaboration of Jeju Special Self-Governing Province², with additional support
¹³⁷ from Flexink³ and Bgrinfo⁴ for data acquisition, and GDS Consulting⁵ for data
¹³⁸ refinement and processing. The scale and diversity of this dataset significantly
¹³⁹ contribute to the strength and practical relevance of this study.

Table 2: Descriptions of several widely used plant disease datasets. **Note:** # stands for the number of something

Dataset	Year	Category	# species	# classes	# images
PlantVillage [37]	2015	Classification	14	38	54,305
PlantDoc [36]	2020	Classification	13	27	2,598
Citrus diseases [39]	2024	Classification	1	5	759
Pomegranate fruit diseases [38]	2024	Classification	1	5	5,099
Fruit disease dataset [40]	2024	Detection	8	6	81,000

¹⁴⁰ For details on the data collection process, including camera settings and
¹⁴¹ acquisition methods, please refer to [40]. Figure 1 presents representative images

¹https://www.nia.or.kr/site/nia_kor/main.do

²<https://www.jeju.go.kr/index.htm>

³<https://flexink.com/en/home/home-en/>

⁴<http://www.bgrinfo.co.kr/>

⁵<http://gdsconsulting.co.kr/>

142 from each class of the fruit disease dataset on eight different plant species,
 143 including banana, fig, lemon, mango, mandarin, olive, passion fruit, and pitaya.



Figure 1: Depiction of the six classes of fruit diseases from the dataset used in this study, with the affected regions highlighted by red BB.

144 Fruits displaying signs of disease, such as spots, lesions, or other visible
 145 deformities, are visually inspected in both natural environments like orchards
 146 and controlled settings such as research greenhouses. Annotations are made at
 147 the lesion or affected region level. Each symptom is evaluated using specific
 148 attributes, including texture, spread, and severity, to ensure accurate labeling.
 149 Annotation guidelines follow established diagnostic criteria specific to each dis-
 150 ease, as outlined below.

- 151 • Anthracnose (*Colletotrichum spp.*): Anthracnose affects a wide variety
 152 of plants, including pitaya, passion fruit, and olive [41]. Anthracnose
 153 typically presents small, sunken, dark brown to black lesions on the fruit's
 154 skin. These lesions may extend and finally lead to significant areas of rot.
 155 The disease can lead to premature fruit drop, leaf loss, and a significant

156 reduction in overall fruit yield.

- 157 • Bacterial fruit blotch (*Acidovorax citrulli*): a serious disease caused by
158 the bacterium *Acidovorax citrulli* [42]. The disease typically manifests
159 as dark, water-soaked lesions on the fruit's surface. These lesions often
160 start small but can rapidly expand to cover large portions of the fruit. As
161 the disease progresses, the affected areas may crack and release a sticky,
162 amber-colored bacterial exudate. The lesions can combine and lead to
163 large, irregular blotches that severely influence the fruit's appearance and
164 marketability. In severe cases, the entire fruit may become soft and rot.
- 165 • Broad mite (*Polyphagotarsonemus latus*): Broad mite is a tiny pest that
166 can cause significant damage to various plants. The mites can infest young
167 lemon fruits [43] and cause russetting or scars on the fruit surface. The af-
168 fected fruits may be deformed and dropped prematurely in extreme cases.
- 169 • Weevil (*Curculionoidea*): Weevils [44] are small beetles that can cause sig-
170 nificant damage to a variety of plants, including fig. Some weevil species
171 burrow into the fruit and cause internal damage that may not be imme-
172 diately visible from the outside. As a consequence, weevil infestation can
173 lead to premature fruit drop, and the affected fruits may become con-
174 tracted. The entry points created by weevils can also serve as gateways
175 for secondary infections by fungi or bacteria, which can further degrade
176 the fruit's quality.
- 177 • Thrips (*Thysanoptera*): Thrips feed by piercing the surface of plant tissues
178 and sucking out the contents of the cells [45], which leads to a range of
179 symptoms that can seriously affect the health and yield of the plants. The
180 most common symptom is surface scars, which affect the quality of the
181 fruits.
- 182 • Fungal infection: Fungal infections can significantly impact the quality,
183 marketability, and production of fruits such as bananas, lemon, mango,
184 and fig [46]. Each type of fruit can be affected by specific fungal pathogens,

which lead to distinct symptoms and potential economic losses. For example, black mildew forms a thin, black layer that can cover significant portions of the fruit's surface, such as lemon and mango. Although the fungus does not penetrate the fruit, it can lead to an unsightly appearance on the affected fruits. Powdery mildew can appear as a white to greyish powdery growth on the skin of figs. This fungal layer can lead to a rough fruit's surface and cause the fruit to crack in severe cases.

The annotation process focused on capturing both the visual characteristics of lesions and any related symptoms or traits that could improve the disease detection performance of the models. A dedicated team of 15 experts from MKG Engineering and Construction (MKGENC) were tasked with a five-month image annotation assignment. Each person annotated approximately 55 images per day to ensure that various disease symptoms were labeled precisely. An open-source annotation tool developed in Python was used to facilitate the entire annotation process [47]. Figure 2 provides an overview of the dataset by showing the number of images for each disease class. It includes a total of 81,000 labeled images, which were split into 80% for training, 10% for validation, and 15% for testing. Therefore, 64,800 images were used for training, while 8,100 images were designated for both validation and testing.

4. Methods

4.1. System Overview

Figure 3 illustrates the primary steps of the fruit disease detection framework, referred to as FD-TR. In this framework, “FD” represents fruit disease detection, while “TR” refers to the transformer-based model. The two core components of the framework are outlined as follows.

- Data pre-processing: Real-world data often presents significant variability due to factors such as inconsistent lighting (e.g., shade, overexposure, underexposure), blurriness (caused by camera motion or low-quality optics), diverse angles (e.g., oblique views, close-ups), and noise (introduced

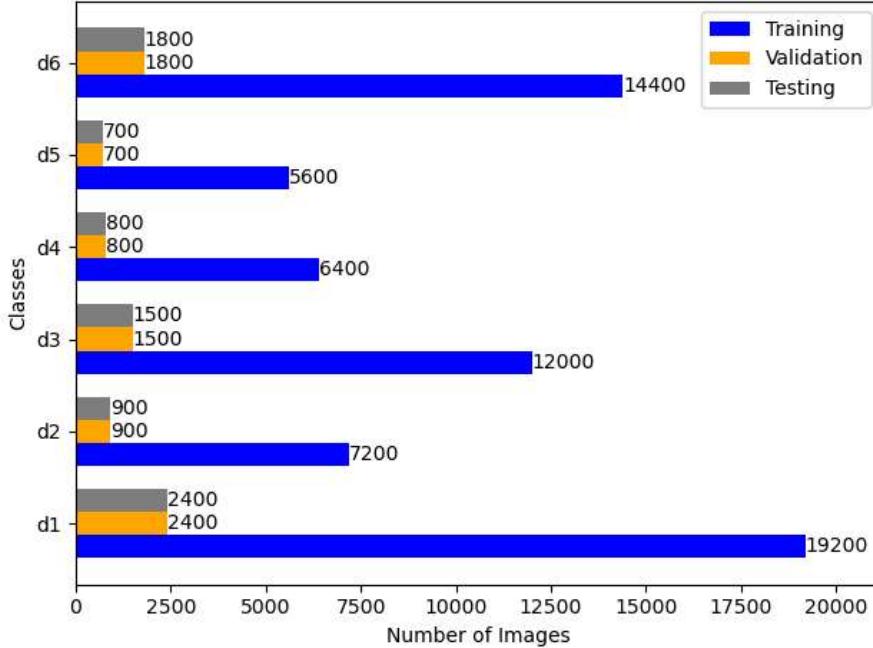


Figure 2: A horizontal bar chart revealing the distribution of images per each disease class from d1 to d6.

214 by sensor imperfections or compression artifacts). Therefore, data aug-
 215mentation is essential to improve the model’s robustness against these
 216 real-world challenges and its ability to generalize to unseen data [48]. The
 217 data augmentation technique involves artificially replicating these condi-
 218 tions within the dataset to effectively increase its size and diversity.

- 219
- 220 • Fruit disease detection: While existing object detection models like Mask-
 221 RCNN [49], YOLO [50], and SSD [51] achieved strong performance on
 222 benchmarks such as COCO [52] and Pascal VOC [53], they rely on man-
 223 ual hyperparameter tuning and multi-stage training. To address these
 224 limitations, we propose FD-TR, a transformer-based architecture with ef-
 ficient end-to-end training. FD-TR focuses on specific parts of the input

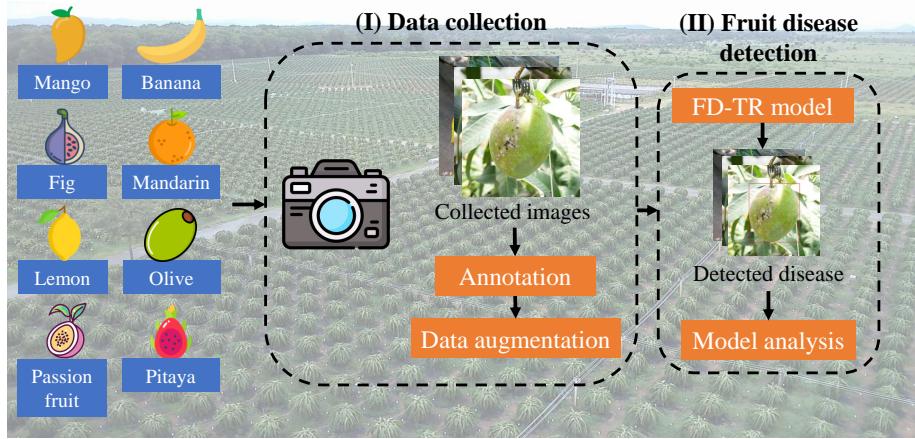


Figure 3: Description of the primary processes of the proposed fruit disease detection framework (FD-TR).

225 image most relevant for identifying diseases. Moreover, feature discrim-
 226 inability scores analysis provides insights into the model's decision-making
 227 process for practical applications [13].

228 *4.2. Data Augmentation*

229 This section outlines the image augmentation process applied to the fruit
 230 disease training dataset to improve the model's robustness and generalization
 231 by simulating various real-world conditions. These augmentation methods were
 232 performed on the original training set to better represent the variability encoun-
 233 tered in real-world agricultural settings. The augmentation techniques expanded
 234 the original training set of 64,800 images by five-fold to 324,000 images.

235 This process involved a series of transformations applied to the original im-
 236 ages, including random horizontal and vertical flips to replicate different ori-
 237 entations of fruits on trees, and rotations at angles of 90°, 180°, and 270° to
 238 enhance the model's invariant to fruit positioning. In addition, color jittering,
 239 where the brightness, contrast, saturation, and hue of input images were ran-
 240 domly adjusted within predefined ranges to mimic varying lighting conditions
 241 and potential color distortions caused by natural environments. To increase
 242 the model's robustness against the effects of camera noise and environmental

²⁴³ factors, Gaussian noise was introduced to the images. Furthermore, random
²⁴⁴ cropping and resizing were performed to expose the model to fruits at different
²⁴⁵ scales and viewpoints. Figure 4 provides a visual representation of the sampled
²⁴⁶ augmented images obtained through different augmentation techniques.

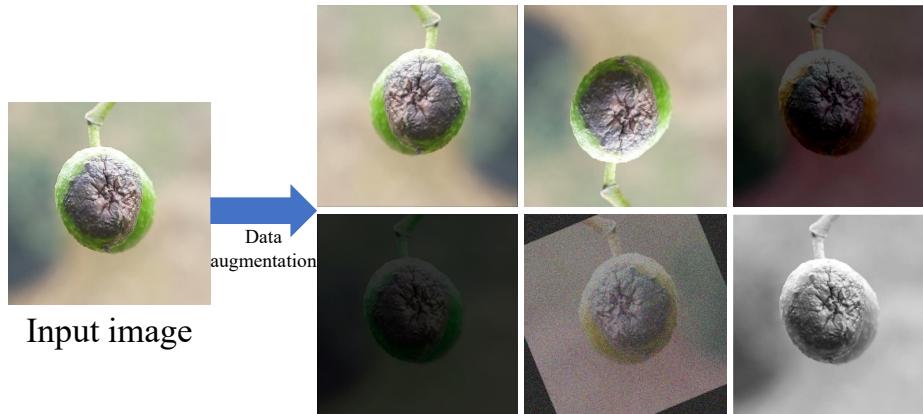


Figure 4: Output images of applying predefined data augmentation techniques on the original dataset.

²⁴⁷ *4.3. Co-DETR Framework*

²⁴⁸ Co-DETR introduces a novel collaborative hybrid assignments training scheme
²⁴⁹ designed to enhance the efficiency and effectiveness of DETR-based detectors.
²⁵⁰ This scheme relies on versatile label assignment strategies to significantly boost
²⁵¹ the encoder's learning capabilities in end-to-end detection frameworks [22]. Co-
²⁵² DETR also optimizes the encoder's learning process by training multiple parallel
²⁵³ auxiliary heads with one-to-many label assignments. In addition, Co-DETR im-
²⁵⁴ proves the overall detection performance by optimizing the attention learning
²⁵⁵ of the decoder through customized positive queries derived from the positive
²⁵⁶ coordinates identified by the auxiliary heads. Figure 5 illustrates the Co-DETR
²⁵⁷ model, which includes three primary modules: a backbone, a transformer en-
²⁵⁸ coder, and a decoder.

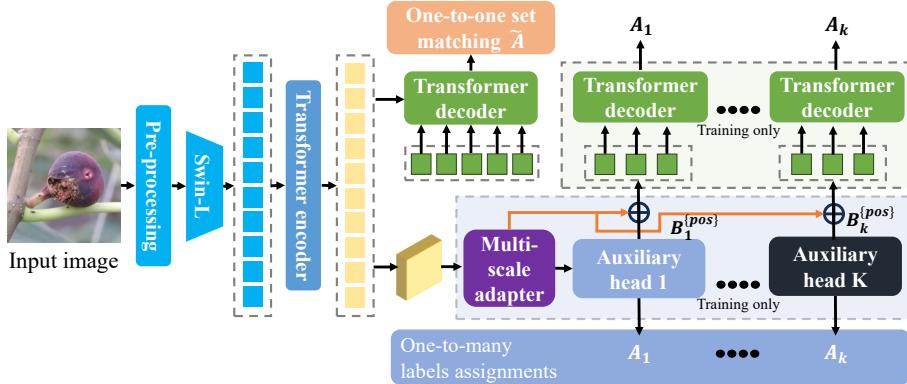


Figure 5: Illustration of the architecture of the Co-DETR approach.

259 According to the standard DETR protocol, the input image is fed into the
 260 backbone and encoder to extract latent features. Several predefined object
 261 queries subsequently interact with the decoder through cross-attention mecha-
 262 nisms. Co-DETR improves this process by integrating a collaborative hybrid
 263 assignment learning and a custom positive query generation module, which op-
 264 timize feature learning in the encoder and attention learning in the decoder.

265 4.3.1. Collaborative Hybrid Assignments Training

266 To address the insufficient supervision of encoder outputs caused by the lim-
 267 ited positive queries in the decoder of standard DETR architectures, Co-DETR
 268 integrates multiple label assignment strategies (e.g., Adaptive Training Sample
 269 Selection (ATSS), Faster R-CNN) with auxiliary supervision heads. These aux-
 270 illiary heads strengthen encoder supervision by refining discriminative learning.
 271 Specifically, after processing the latent features \mathcal{F} , the encoder transforms them
 272 into a feature pyramid $\mathcal{F}_1, \dots, \mathcal{F}_J$ via a multi-scale adapter, where J denotes the
 273 number of feature maps with downsampling stride of 2^{2+J} . Following the ViT-
 274 Det framework, Co-DETR constructs its feature pyramid using a single-scale
 275 encoder feature map, which is upsampled using bilinear interpolation.

276 For example, the feature pyramid is built by sequentially applying upsam-

pling (stride 2 with 3×3 convolution) or downsampling to the encoder’s single-scale feature. In multi-scale encoders, only the coarsest resolution features are downsampled to generate the feature pyramid. For each K collaborative heads, the predicted output \hat{P}_i is sequentially propagated through the feature pyramid $\mathcal{F}_1, \dots, \mathcal{F}_J$. Within the i -th head, module A_i computes supervised targets for positive and negative samples, $P_i^{\text{pos}}, B_i^{\text{pos}}, P_i^{\text{neg}}$, using the supervised target set G , as follows:

$$P_i^{\{\text{pos}\}}, B_i^{\{\text{pos}\}}, P_i^{\{\text{neg}\}} = A_i(\hat{P}_i, G) \quad (1)$$

where pos and neg represent the spatial coordinates classified as positive and negative by A_i . The index j corresponds to the feature index within the feature pyramid \mathcal{F}_j . B_i^{pos} denotes the spatial coordinates of the positive samples, while P_i^{pos} and P_i^{neg} refer to the supervised targets associated with these coordinates, including both category labels and BB regression offsets.

The encoder loss function can be defined as follows:

$$\mathcal{L}_i^{\text{enc}} = \mathcal{L}_i(\hat{P}_i^{\{\text{pos}\}}, P_i^{\{\text{pos}\}}) + \mathcal{L}_i(\hat{P}_i^{\{\text{neg}\}}, P_i^{\{\text{neg}\}}) \quad (2)$$

For negative samples, the regression loss is excluded from consideration. The objective of optimization for the K auxiliary heads is therefore defined as follows:

$$\mathcal{L}^{\text{enc}} = \sum_{i=1}^K \mathcal{L}_i^{\text{enc}} \quad (3)$$

4.3.2. Customized Positive Queries Generation

In the one-to-one matching paradigm, each ground-truth box is paired with a single specific query as its supervised target. However, when the number of positive queries is insufficient, this can lead to inefficient cross-attention learning within the transformer decoder. To address this issue, Co-DETR generates a diverse set of customized positive queries. Specifically, in the i -th auxiliary head,

²⁹⁹ the customized positive query $Q_i \in \mathbb{R}^{M_i \times C}$ (where M_i represents the number
³⁰⁰ of positive samples) is generated through the following process:

$$Q_i = \text{Linear}(\text{PE}(B_i^{\{\text{pos}\}})) + \text{Linear}(\mathbb{E}(\{\mathcal{F}_*\}, \{\text{pos}\})) \quad (4)$$

³⁰¹ Here, $\text{PE}(\cdot)$ represents positional encoding, which extracts the relevant fea-
³⁰² ture from $\mathbb{E}(\cdot)$ based on the spatial positive and negative coordinates (j, \mathcal{F}_j) .

³⁰³ Therefore, there are $K+1$ query groups involved in the one-to-one matching
³⁰⁴ process, including those with label assignments. The auxiliary label assignment
³⁰⁵ shares weights with the standard L decoder layers. In the auxiliary branches, all
³⁰⁶ queries are conditioned on the positive query, eliminating the need for redundant
³⁰⁷ matching. The loss for the l -th decoder layer in the i -th auxiliary head is
³⁰⁸ formalized as follows:

$$\mathcal{L}_{i,l}^{\text{dec}} = \tilde{\mathcal{L}}(\tilde{P}_{i,l}, P_i^{\{\text{pos}\}}) \quad (5)$$

³⁰⁹ where $\tilde{\mathcal{L}}^{\text{dec}}$ denotes the loss from the original one-to-one matching branch.

³¹⁰ Finally, the global objective function of Co-DETR is defined as:

$$\mathcal{L}^{\text{global}} = \sum_{l=1}^L (\tilde{\mathcal{L}}_l^{\text{dec}} + \lambda_1 \sum_{i=1}^K \mathcal{L}_{i,l}^{\text{dec}} + \lambda_2 \mathcal{L}^{\text{enc}}) \quad (6)$$

³¹¹ Here, λ_1 and λ_2 are the coefficients that balance the different losses.

³¹² 4.4. Model Customization

³¹³ Although Co-DETR can be applied to state-of-the-art transformer architec-
³¹⁴ tures such as DETR with Improved deNoising anchOr box (DINO) [54] and
³¹⁵ Deformable DETR [33] for fruit disease detection, the performance of these
³¹⁶ base models remains sensitive to critical factors like label assignment strategies,
³¹⁷ robustness to complex backgrounds, and adaptability under varying environ-
³¹⁸ mental conditions. To optimize transformer-based detection for fruit disease
³¹⁹ detection, several targeted adjustments were introduced to the original trans-
³²⁰ former models' architecture and optimization process. These modifications were
³²¹ implemented before applying the Co-DETR approach, as outlined in Table 4.

322 The modifications include integrating BatchFormerV2 to enhance feature
 323 representation through batch-based learning, adopting the LAMB optimizer,
 324 known for its efficiency in training large-scale models [55], and utilizing the
 325 Complete Intersection over Union (CIoU) loss function instead of the GIoU to
 326 improve localization accuracy. These modifications are expected to improve the
 327 baseline models' performance and generalization capabilities on the fruit disease
 328 domain.

Table 3: Comprehensive explanation of the model fine-tuning process

Model	Auxiliary	Loss	Optimizer
Deformable DETR	N/a	Hybrid (L1 + GIoU)	AdamW
DINO	N/a	Hybrid (L1 + GIoU)	AdamW
FD-TR (This study)	BatchFormerV2 [56]	Hybrid (L1 + CIoU)	LAMB

- 329 • BatchFormerV2 (BF): Proposed by Hou et al. [56], BF enhances trans-
 330 formers' capacity to model inter-sample relationships within mini batches.
 331 Unlike conventional transformer blocks that operate on pixel- or patch-
 332 level feature maps, BF processes feature structured by batch size. In
 333 FD-TR framework, BF implements a two-stream architecture where both
 334 branches share weights and merge into a unified transformer decoder. This
 335 design ensures efficiency and coherence during the training process as all
 336 shared blocks are consistently trained with the same weights. Moreover,
 337 the original transformer blocks retain their full functionality without BF,
 338 which minimizes any additional computing during inference. The appli-
 339 cation of BatchFormerV2 into various transformer models, such as DETR
 340 [32] and Deformable-DETR [33], consistently demonstrated a performance
 341 improvement of over 1.3 mAP on the benchmark MS COCO dataset.
- 342 • Complete Intersection over Union (CIoU): The Generalized IoU (GIoU)
 343 extends the standard IoU metric by measuring the overlap between the
 344 predicted and ground truth BB while considering areas outside their in-

345 tersection [57]. CIoU improves GIoU by introducing additional terms that
 346 account for localization precision and aspect ratio alignment. This refine-
 347 ment enables better convergence and improved detection accuracy com-
 348 pared to GIoU loss. Therefore, CIoU and L1 loss are utilized to calculate
 349 the box regression reconstruction loss for FD-TR model in this study.

$$\mathcal{L}_{CIoU} = 1 - IoU + \frac{d^2(p, p^{gt})}{c^2} + \alpha V \quad (7)$$

350 The variable c denotes the diagonal length of the smallest enclosing box
 351 that covers both the predicted and ground truth BB, while d represents
 352 the Euclidean distance between their center points. p and p^{gt} refer to the
 353 central points of the predicted and ground truth BB, respectively. The
 354 variable V measures the consistency of the aspect ratios, and α serves
 355 as a trade-off parameter that assigns less weight when the overlap is low
 356 and more weight when the overlap is high. The value of α is computed
 357 dynamically as:

$$v = \frac{4}{\pi^2} \left(\arctan \frac{w^{gt}}{h^{gt}} - \arctan \frac{w}{h} \right)^2, \quad \alpha = \frac{v}{(1 - IoU) + v}, \quad (8)$$

358 We compared a dynamically computed α with fixed values $\alpha \in \{0.25, 0.5, 1.0\}$
 359 on the validation set (Table 4). The dynamic α showed the highest peak
 360 validation mAP (0.77), but $\alpha = 0.5$ achieved a comparable validation
 361 mAP (0.75). To improve reproducibility and make cross-experiment com-
 362 parisons more straightforward, we therefore use $\alpha = 0.5$ in all subsequent
 363 experiments. Moreover, fixed α also reduces hyperparameter tuning. If
 364 the aim is to maximize single-run peak mAP, dynamic α remains an ap-
 365 propriate choice.

- 366 • LAMB optimizer: While AdamW is commonly considered the default opti-
 367 mizer for a variety of vision transformer-based models [12], [58] have iden-
 368 tified potential training instability, particularly when there is an increased

Table 4: Ablation study comparing the dynamic α and fixed candidates {0.25, 0.5, 1.0}

α values	Mean mAP
Dynamic	0.77
0.25	0.68
0.5	0.75
1	0.7

369 ratio between the L2-norm of weights and gradients. To mitigate this is-
 370 sue, this study adopts the Layer-wise Adaptive Large Batch Optimization
 371 (LAMB) optimizer as an alternative. LAMB combines the strengths of
 372 both the Adam and Layer-wise Adaptive Rate Scaling (LARS) optimizers
 373 [55]. In particular, the layer-wise adaptive technique from LAMB normal-
 374 izes each dimension based on the square root of the second moment, while
 375 also applying layer-wise normalization. This method has been proved to
 376 be effective for distributed training and has demonstrated effectiveness in
 377 transformer models on large-scale datasets.

$$\begin{aligned}
 m_t &= \beta_1 m_t^{(\text{prev})} + (1 - \beta_1) g_t \\
 v_t &= \beta_2 v_t^{(\text{prev})} + (1 - \beta_2) g_t^2 \\
 m_t &= \frac{m_t}{1 - (\beta_1)^t} \\
 v_t &= \frac{v_t}{1 - (\beta_2)^t} \\
 r_t &= \frac{m_t}{\sqrt{v_t} + \epsilon} \\
 x_{t+1}^{(i)} &= x_t^{(i)} - \eta_t \frac{\phi \left(\|x_t^{(i)}\| \right)}{\left\| r_t^{(i)} + \lambda x_t^{(i)} \right\|} \left(r_t^{(i)} + \lambda x_t^{(i)} \right)
 \end{aligned} \tag{9}$$

378 where the hyperparameters β_1 and β_2 regulate momentum and weight
 379 decay, respectively. m_t refers to the first moment estimate at time step t ,
 380 and v_t indicates the second moment estimate. The parameter λ manages
 381 the degree of layer-wise adaptiveness, while η_t represents the learning rate

382 vector at time t , and ϕ denotes the parameter vector at the same instance.
 383 A small constant ϵ is introduced to prevent division by zero. In addition,
 384 r_t represents the update ratio used in the LAMB optimizer.

385 *4.5. Feature Discriminability Scores Analysis*

386 After training, feature discriminability maps are generated by analyzing the
 387 multi-scale feature outputs from the DETR-based model [32]. They offer valua-
 388 ble insights into how the model distributes its focus on different regions of
 389 the input image. The feature discriminability scores are obtained by extract-
 390 ing multi-scale features from the model’s final layers. For each feature map,
 391 the L2-norm is computed across the channel dimension to quantify the activa-
 392 tion strength at each spatial location, consistent with established visualization
 393 practices for CNN activations [59]. The resulting feature discriminability scores
 394 are then normalized by their maximum values to ensure consistent intensity of
 395 different scales.

396 To visualize the feature discriminability scores, each normalized feature map
 397 is resized to match the dimensions of the input image using linear interpolation.
 398 The resized maps from each scale are then aggregated by combining them to-
 399 gether, followed by averaging to produce a final feature map that integrates in-
 400 formation from all scales. This final map highlights the regions that the model
 401 considers most relevant during the prediction process, with higher values in-
 402 dicating areas of greater focus. The output feature discriminability map is a
 403 valuable tool for evaluating the model’s interpretability and its ability to cor-
 404 rectly identify disease-affected regions in the image.

405 Let feats be a list of multi-scale feature maps, each with dimensions $L \times B \times$
 406 $C \times H \times W$. L is the number of layers or scales, $B = 1$ is the batch size, C
 407 is the number of channels, and $H \times W$ are the spatial dimensions. Based on
 408 the multi-scale feature maps, the feature discriminability map attn_map can
 409 be mathematically represented as follows:

$$\text{attn_map} = \frac{1}{L} \sum_{i=1}^L \text{resize} \left(\frac{\|\text{feat}[i]\|_2}{\max(\|\text{feat}[i]\|_2 + \varepsilon)}, H_{\text{img}}, W_{\text{img}} \right) \quad (10)$$

410 where $\|\text{feat}[i]\|_2$ represents the L2-norm of the feature map at the i -th scale,
 411 calculated along the channel dimension for generating a feature map of size $H \times$
 412 W . The term $\max(\|\text{feat}[i]\|_2)$ denotes the maximum value in the normed feature
 413 map, which is used to normalize the map. The function $\text{resize}(\cdot, H_{\text{img}}, W_{\text{img}})$
 414 interpolates the normalized feature map to match the dimensions $H_{\text{img}} \times W_{\text{img}}$ of
 415 the input image. The summation aggregates the resized feature discriminability
 416 maps from all scales, and the division by L averages the aggregated map.

417 In Figure 6, the feature discriminability map extraction of an input image
 418 highlights how FD-TR effectively focuses on disease-affected regions using multi-
 419 scale features from the encoder. The map illustrates the DETR-based model's
 420 ability to precisely target the main regions showing disease symptoms. This
 421 demonstrates the model's robustness and accuracy in detecting various fruit
 422 diseases.

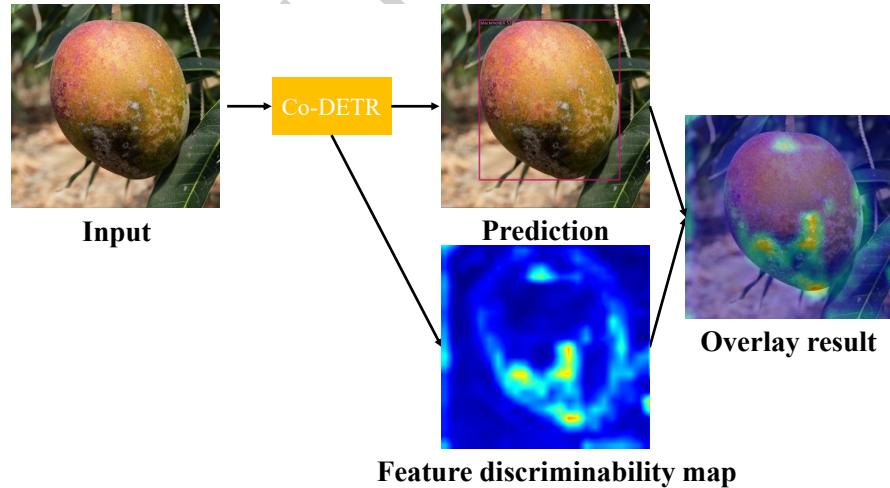


Figure 6: Visualization of the feature discriminability map prediction process.

423 *4.6. Implementation Details*

424 The fruit disease detection framework was developed using the MMDetection
 425 library v2.25.3, built on PyTorch 1.11.0. To ensure consistent and fair
 426 experimentation, all detection models in the study utilized ResNet-50 and Swin
 427 backbone pre-trained on the ImageNet dataset. The training process was con-
 428 ducted on an Nvidia A100 GPU with 40 GB of memory.

429 We integrate our Co-DETR into existing DETR-like pipelines while main-
 430 taining similar training settings with the baseline models. For $K = 2$, we im-
 431 plement both ATSS and Faster-RCNN as auxiliary heads, whereas for $K = 1$,
 432 we use only the ATSS head. In addition, the number of learnable object queries
 433 is set to 300, and the weight coefficients $\{\lambda_1, \lambda_2\}$ are set to their default values
 434 of $\{1.0, 2.0\}$.

435 For all transformer-based experiments (FD-TR and DETR variants), each
 436 model is trained for up to 15 epochs with validation process perform at the
 437 end of each epoch. Early stopping is applied to the validation bounding-box
 438 loss with a patience of three epochs and a minimum improvement threshold
 439 $\Delta = 10^{-3}$. If the bounding-box loss fails to decrease by at least Δ for three
 440 consecutive epochs, training halts and the model reverts to the weights from
 441 the epoch with the lowest validation loss.

442 *4.7. Evaluation Protocols*

443 In this section, we comprehensively evaluate the fruit disease recognition
 444 framework using several standard metrics, including mAP, precision, and recall.
 445 These metrics are computed based on the three elements of the confusion matrix:
 446 true positive (TP), false positive (FP), and false negative (FN). Precision reveals
 447 the ratio of correctly predicted positive instances out of all predicted positives,
 448 while recall captures the proportion of true positives among all actual positives
 449 in the dataset. The formulation of these metrics is as follows:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$
(11)

450 To evaluate the overall detection accuracy of multiple disease classes, a
 451 standard average precision metric was calculated. In particular, we adopt
 452 $AP@[IoU = 0.50 : 0.95]$, which measures the detection performance at IoU
 453 thresholds from 0.50 to 0.95. This threshold is used to evaluate the model's
 454 ability to localize fruit diseases by calculating the area under the precision-
 455 recall curve at the specified IoU threshold. The AP for each class is determined
 456 from this curve, and the mAP is then computed as the average of the AP values
 457 on all disease types. The mAP is expressed as follows:

$$\text{mAP} = \frac{1}{N} \sum_{i=1}^N \text{AP}_i \quad (12)$$

458 where N represents the number of disease types, and AP_i denotes the average
 459 precision for the i -th disease class. AP_i is calculated based on the precision-recall
 460 curve for that disease type.

461 5. Results

462 5.1. Comparison of Transformer Models

463 In this experiment, a comprehensive comparison of fruit disease detection
 464 performance is conducted by applying Co-DETR on various DETR-based mod-
 465 els, including Deformable DETR [33] and DINO [54]. Moreover, two different
 466 backbones, Swin Transformer and ResNet-50, are employed and compared, re-
 467 sulting in a total of four model variants. The models include Co-DETR on the
 468 Deformable DETR with the ResNet-50 backbone (co_deformable_detr_r50),
 469 Co-DETR on the Deformable DETR with the Swin backbone (co_deformable_detr_swin),
 470 Co-DETR on the DINO model with the ResNet-50 backbone using 5-scale fea-
 471 ture processing (co_dino_5scale_r50), and Co-DETR on the DINO model with
 472 the Swin backbone using 5-scale feature processing (co_dino_5scale_swin).
 473 The performance comparison is shown in Figure 7.

474 Overall, the co_dino_5scale_swin model demonstrates the highest perfor-
 475 mance with a detection mAP starting from around 0.6 and steadily improving

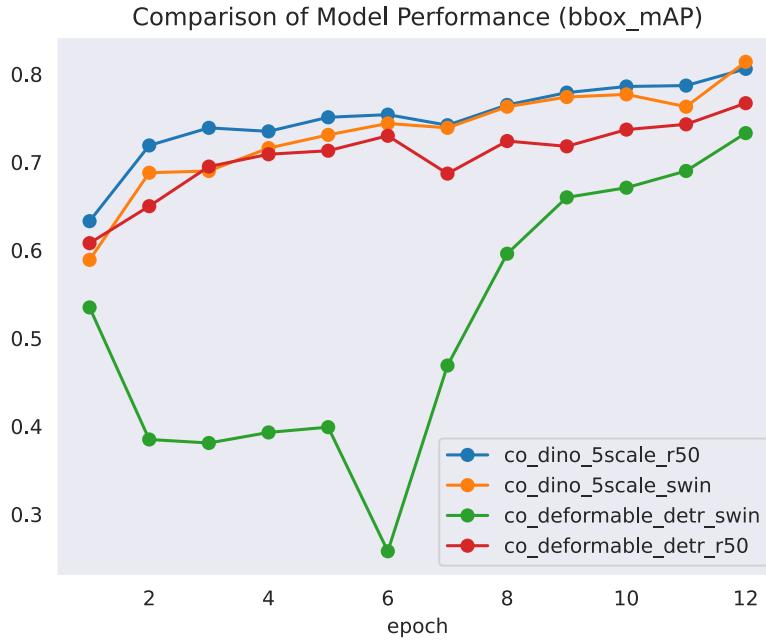


Figure 7: Comparison of fruit disease detection performance using Co-DETR applied to two baseline DETR models: Deformable DETR and DINO.

to around 0.81 by the 12th epoch. This indicates that the Swin backbone combined with 5-scale feature extraction is particularly effective in detecting fruit diseases. The `co_dino_5scale_detr` model also performs well, closely following `co_dino_5scale_r50` while maintaining a high performance at around 0.79 at the 12th epoch. The `co_deformable_detr_r50` model shows relatively stable performance but with lower performance than the DINO models. In contrast, the `co_deformable_detr_swin` model exhibits significant fluctuations in its performance, particularly between epochs 5 and 7, where it experiences a sharp drop to around 0.25 mAP. However, the model recovers rapidly from epoch 8th and reaches a comparable mAP of approximately 0.72 by the 12th epoch. These fluctuations suggest that while the Deformable DETR architecture may be more sensitive to certain training conditions, it is capable of eventually reaching a

488 competitive performance.

489 Given that Co-DETR on the DINO model with the Swin backbone demon-
 490 strated the highest fruit disease detection performance, we selected this configu-
 491 ration as the default model for subsequent experiments (referred to as FD-TR).
 492 This extension was chosen because it delivered robust and stable detection ac-
 493 curacy during training and validation. FD-TR was then used to evaluate the
 494 effects of additional enhancements, such as data augmentation techniques, hy-
 495 perparameter tuning, and its deployment in real-world environments.

496 *5.2. Preprocessing Module Analysis*

497 This section examines the impact of data augmentation on the proposed
 498 FD-TR model by comparing its results with the one trained on raw data. As
 499 shown in Table 5, FD-TR trained with augmented images outperformed the
 500 one trained on raw data. For example, the mAP increased by 0.05 from 0.76
 501 to 0.81, indicating better overall accuracy in detection. The data augmenta-
 502 tion approach also reduced the false positive detection (higher precision) and
 503 increased the rate of correctly identifying true positives (higher recall).

Table 5: Comparison of FD-TR model performance on original and augmented data.

	mAP	Precision	Recall
Original data	0.76	0.75	0.78
Data augmentation	0.81	0.79	0.82

504 The observed performance improvement suggests that data augmentation
 505 plays a crucial role in boosting FD-TR model’s ability to detect fruit diseases
 506 with higher detection accuracy. By introducing variations in the training data,
 507 augmentation not only boosts detection precision but also significantly improves
 508 the model’s robustness.

509 *5.3. PD-TR Performance Evaluation*

510 Figure 8 provides a detailed performance evaluation of FD-TR model, which
 511 consists of two charts.

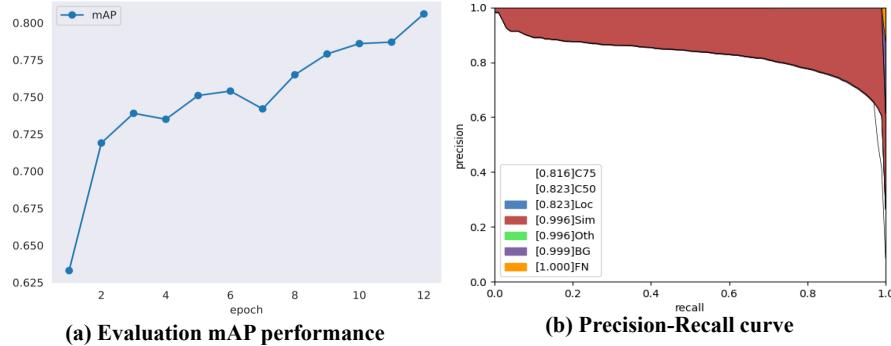


Figure 8: Detailed performance evaluation of FD-TR model using different evaluation metrics.

- The evaluation mAP performance (a) plots FD-TR’s mAP over 12 epochs of training. The mAP started at approximately 0.625 and steadily increased. It peaked at around 0.8 by the 12th epoch. This consistent improvement in mAP indicated that the model was learning effectively and becoming increasingly better at detecting diseases as training progresses. The gradual increase suggested that the model generalized well and converged to high performance, especially in the later epochs.
- The precision-recall curve (b) represents the trade-off between precision and recall for different thresholds. This curve can be used to evaluate how well FD-TR performs on different confidence levels. Overall, the model accurately detected diseases with minimal false positives because the curve showed high precision for most recall values. Key metrics like C75, C50, and Loc revealed precise localization and detection capabilities, with precision values around 0.816 to 0.823, suggesting that the model performed well even under challenging IoU thresholds. The model also excelled in distinguishing between similar diseases (Sim) and avoiding background errors (BG), with a precision near 1.0 in both cases. The curve’s slight decline at very high recall indicated that while the model maintained accuracy on most conditions, it introduced minor false positives when recall was pushed to its limit. Finally, a good false negative (FN) showed that

532 the model had a very low rate of missing diseased fruits.

533 Table 6 describes the experimental results of FD-TR framework in detect-
 534 ing six different fruit diseases, including anthracnose (d1), bacterial fruit blotch
 535 (d2), broad mite (d3), weevil (d4), thrips (d5), and fungal infection (d6). In
 536 general, FD-TR framework showed consistent performance in detecting all dis-
 537 ease classes with an average mAP of 0.81, precision of 0.79, and recall of 0.82.
 538 The model achieved the highest performance for detecting d3 and d6 with the
 539 mAP scores of 0.88 and 0.86, respectively. These classes also obtained strong
 540 precision (0.85 and 0.84) and recall (0.89 and 0.88). On the other hand, the
 541 detection performance for d2 and d5 was slightly lower, with mAP values of
 542 0.74 and 0.77. The low detection performance of d2 and d5 could be due to sev-
 543 eral factors: 1) fewer labeled instances in the training data, which limited the
 544 framework’s ability to extract distinct features for these diseases, and 2) visual
 545 similarities between d2 and d5 made it challenging for the model to effectively
 546 differentiate between these diseases and others.

Table 6: Evaluation results of the proposed model on different fruit disease classes.

	d1	d2	d3	d4	d5	d6	Average
mAP	0.78	0.74	0.88	0.83	0.77	0.86	0.81
Precision	0.76	0.73	0.85	0.8	0.77	0.84	0.79
Recall	0.79	0.77	0.89	0.82	0.79	0.88	0.82

547 5.4. Analysis of the Feature Discriminability Analysis

548 Table 7 reports the mean and standard deviation of the normalized L₂-norm
 549 discriminability scores for each disease class over the test set. The scores confirm
 550 that the model concentrates more strongly on classes with more distinct lesion
 551 features, such as anthracnose, broad mite, and fungal infection.

552 Figure 9 provides a detailed description of the proposed framework in ef-
 553 fectively detecting six distinct fruit diseases. Each row in the figure serves a

Table 7: Mean (\pm std) of feature discriminability scores per disease class.

Disease class	Mean (\pm std) score
(d1) Anthracnose	0.86 \pm 0.05
(d2) BFB	0.65 \pm 0.10
(d3) Thrips	0.72 \pm 0.03
(d4) Weevil	0.68 \pm 0.08
(d5) Broad mite	0.80 \pm 0.01
(d6) Fungal infection	0.83 \pm 0.03

554 distinct purpose. Row (a) displays the original images of fruits affected by dis-
 555 eases such as anthracnose, BFB, thrips, weevil, broad mite, and fungal infection.
 556 The second row (b) demonstrates the model’s detection results by highlighting
 557 the areas where the model has identified disease presence with BB and predicted
 558 labels.

559 Overall, the model correctly predicted and localized the fruit diseases pre-
 560 cisely. In order to explain the model’s prediction process, the third row (c) fur-
 561 ther shows attention-weight visualizations from FD-TR model. The extracted
 562 attention map reveals where the model is focusing its attention on the images.
 563 Warmer color areas indicate higher focus, which is typically around spots show-
 564 ing visible symptoms of the disease. It can be concluded by observing the at-
 565 tention maps that the model focused on disease regions but also provided visual
 566 explanations for its predictions. Moreover, the attention analysis also enhanced
 567 trust and understanding in its diagnostic capabilities.

568 Figure 10 demonstrates FD-TR model’s performance on some challenging
 569 fruit disease detection cases, such as lighting variations, image blurring, and
 570 low contrast. The top row (a) displays the input images, while the second
 571 row (b) shows the detection results, including the predicted BB, disease name
 572 and confidence score. The attention map visualization in the bottom row (c)
 573 indicates how FD-TR model focuses on specific regions of the image for its
 574 predictions.

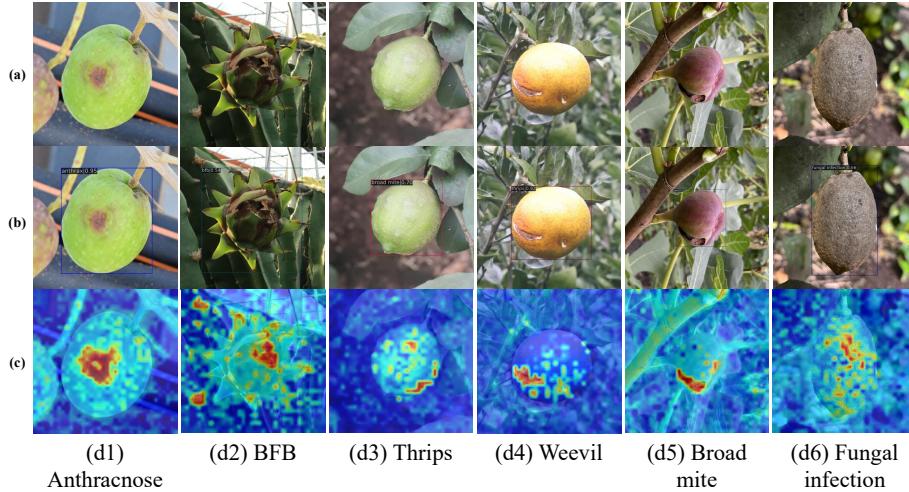


Figure 9: The proposed model’s outputs for each fruit disease, including (a) input images, (b) detection results, and (c) feature discriminability visualizations.

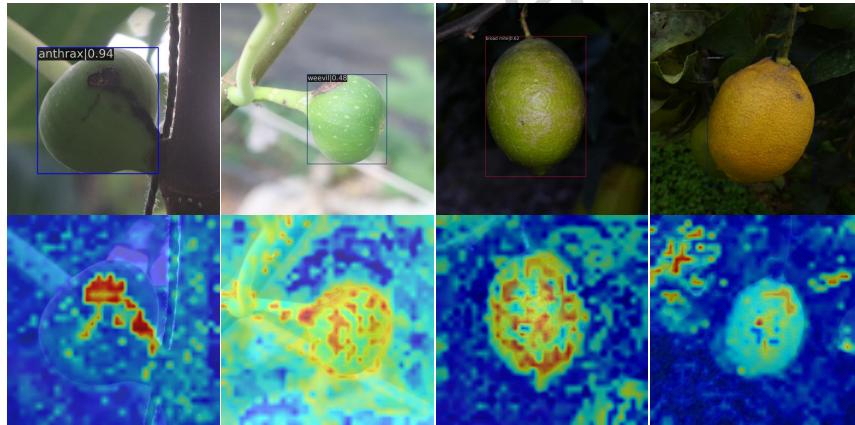


Figure 10: The proposed model’s outputs for challenging cases, including (a) input images, (b) detection results, and (c) feature discriminability visualizations.

575 FD-TR model demonstrates strong disease prediction performance in real-
 576 world conditions. This is important for practical deployment in agricultural
 577 environments where image quality may vary. For instance, the model demon-
 578 strates its robustness by accurately detecting anthracnose (first column) and
 579 black mold (fourth column) with high confidence scores of 0.94 and 0.77, re-
 580 spectively. In these cases, the model focuses effectively on the infected areas

581 with well-defined and concentrated regions in the attention maps.

582 In contrast, for more challenging cases such as weevil (second column) and
 583 broad mite (third column), the attention maps appear more diffuse, with less
 584 sharply defined focus areas. Factors such as image blurring and uneven lighting
 585 seem to affect the model’s ability to identify the diseased regions accurately.
 586 This results in a lower confidence score for weevil detection (0.48), indicating the
 587 model’s difficulty in isolating the specific features of the disease. Nevertheless,
 588 FD-TR model manages to generate reasonable predictions.

589 *5.5. Analysis of the Effectiveness of Customized Components to the Perfor-
 590 mance of FD-TR Model*

591 This section reports the effectiveness of important components of the pro-
 592 posed fruit disease detection model’s performance. Table 8 summarizes the
 593 ablation study’s results of each component of FD-TR model.

Table 8: Ablation analysis for evaluating the effects of different components on the perfor-
 mance of FD-TR model.

Configuration	CIoU+L1 loss	LAMB optimizer	BatchFormerV2	mAP
Baseline	–	–	–	0.812
+ CIoU only	✓	–	–	0.847
+ LAMB only	–	✓	–	0.818
+ BatchFormerV2 only	–	–	✓	0.853
+ CIoU & LAMB	✓	✓	–	0.834
+ CIoU & BatchFormerV2	✓	–	✓	0.882
+ LAMB & BatchFormerV2	–	✓	✓	0.838
Full integration	✓	✓	✓	0.894

594 The baseline configuration, without any of the proposed components, achieved
 595 an mAP of 0.812. When added individually, CIoU+L1 loss improved the mAP
 596 to 0.847, which demonstrated its significant contribution to the model perfor-
 597 mance. The LAMB optimizer showed a marginal improvement to 0.818, while

598 BatchFormerV2 alone boosted the mAP to 0.853. Further analysis of pairwise
 599 combinations revealed additional insights. The combination of LAMB optimizer
 600 with CIoU+L1 loss or BatchFormerV2 yielded lower mAP compared to using
 601 CIoU+L1 loss or BatchFormerV2 alone. However, these configurations achieved
 602 an average of 13% faster convergence and reduced training time. Meanwhile,
 603 the integration of both CIoU+L1 loss with BatchFormer V2 led to a substan-
 604 tial increase to 0.882, which suggested a stronger interaction between these two
 605 components. Finally, the full integration of all three components achieved the
 606 highest performance at 0.894, which highlighted their effectiveness on enhancing
 607 the model’s capabilities.

608 *5.6. Comparison with Other Models*

609 Table 9 presents a performance comparison between the proposed FD-TR
 610 model and five other state-of-the-art detection models (YOLOv8 [60], SSD [51],
 611 DETR [32], Deformable DETR [33], DINO [54]). When evaluated on the val-
 612 idation dataset, FD-TR consistently outperformed the others on all metrics.
 613 Specifically, FD-TR significantly outperformed the next best model by 9% with
 614 an mAP of 0.89. In addition, with high precision and recall values, FD-TR
 615 demonstrated its ability to accurately identify and localize objects. In contrast,
 616 SSD exhibited the lowest performance, with an mAP of 0.69, and a precision
 617 and recall of 0.66 and 0.70, respectively.

Table 9: Model performance evaluation between the proposed model and five state-of-the-art DL models on the validation dataset.

Model name	mAP	Precision	Recall
SSD [51]	0.69	0.66	0.70
YOLOv8 [60]	0.8	0.81	0.83
DETR [32]	0.72	0.71	0.73
Deformable DETR [33]	0.74	0.74	0.77
DINO [54]	0.75	0.74	0.76
FD-TR (Ours)	0.89	0.86	0.87

618 Moreover, while other transformer-based models like DETR, Deformable
 619 DETR, and DINO demonstrated higher performance over SSD, they were con-
 620 sistently outperformed by FD-TR. For example, Deformable DETR showed an
 621 mAP of 0.74, precision of 0.74, and recall of 0.77, while DINO achieved slightly
 622 better precision and recall but a comparable mAP. YOLOv8, well-known for
 623 its performance, performed well with an mAP of 0.80 but was outperformed by
 624 FD-TR in all metrics. The results highlight that FD-TR model provides the
 625 most accurate and reliable predictions for fruit disease detection due to several
 626 enhancements such as the Co-DETR scheme and effective integration of other
 627 components.

628 *5.7. Comparison on Various Benchmark Datasets*

629 Table 10 describes the performance of FD-TR on four publicly available
 630 datasets compared to the baseline model (Co-DETR). This table includes two
 631 agricultural datasets (PlantVillage [61] and Pest-D2Det [62]) and widely used
 632 general benchmarks (COCO [52] and VOC2012 [63]). The variation in domain
 633 complexity, class count, and dataset size provides a comprehensive evaluation
 634 of the model’s adaptability.

Table 10: FD-TR performance and gains compared to baseline methods. Note: pp stands for absolute gain in percentage points

Dataset	Domain	# classes	# images	Baseline mAP	FD-TR mAP	pp (%)
PlantVillage	Agriculture	38	54,308	0.407 (YOLOv8 [64])	0.594	18.7
Pest-D2Det	Agriculture	10	9,472	0.704 (D2Det [65])	0.731	2.7
COCO	General	80	118,287	0.659 (Co-DETR [22])	0.589	-7
VOC2012	General	20	11,540	0.804 (CoupleNet [66])	0.812	0.8

635 In the agricultural domain, FD-TR demonstrates significant advancements,
 636 particularly on PlantVillage, where it achieves an mAP of 0.594, an 18.7%
 637 gain over the YOLOv8 baseline (0.407). This improvement highlights FD-TR’s
 638 effectiveness in handling high-class diversity (38 classes) and complex disease
 639 manifestations. Similarly, on Pest-D2Det, FD-TR obtains an mAP of 0.731, a
 640 2.7% increase over the D2Det baseline (0.704), which confirms its strength in
 641 pest detection tasks with fewer classes (10). These results indicate that FD-
 642 TR performs well in agricultural context, where precise feature learning and
 643 optimization are critical for real-world applications like crop monitoring.

644 For general-domain datasets, FD-TR exhibits robust but context-dependent
 645 performance. On VOC2012 (20 classes), it achieves an mAP of 0.812, a modest
 646 0.8% improvement over the baseline CoupleNet (0.804). However, on COCO
 647 (80 classes), FD-TR records an mAP of 0.589, approximately 7.0% below Co-
 648 DETR’s reported 0.659. This gap does not undermine FD-TR’s efficiency but
 649 rather reflects key architectural and training differences. Co-DETR leverages
 650 a large ViT-Large backbone and extensive pre-training on Objects365 (opti-
 651 mized for large-scale benchmarks like COCO). In contrast, FD-TR prioritizes
 652 lightweight efficiency using Swin as backbone, and targets agricultural special-
 653 ization without target pre-training. FD-TR’s modifications (BatchFormerV2 for
 654 enhanced feature representation, CIoU for improved box learning, and LAMB
 655 for training stabilization) emphasize domain-specific adaptability over maximiz-
 656 ing COCO accuracy. Despite the lower score, FD-TR remains competitive with
 657 many transformer-based detectors and aligns with its goal of balancing per-
 658 formance, efficiency, and specialization. Overall, these results confirm FD-TR’s
 659 contributions, particularly in agricultural contexts, while maintaining versatility
 660 across domains.

661 *5.8. Real-world Robustness Analysis*

662 To evaluate the model’s ability to distinguish healthy fruits, which is a crit-
 663 ical requirement for real-world agricultural applications, an independent test
 664 dataset comprising 500 images of healthy fruits was collected. These images

665 were curated from a publicly available agricultural image repository and veri-
 666 fied by domain experts to confirm the absence of disease symptoms. This dataset
 667 was excluded from training and reserved solely for evaluating the model’s per-
 668 formance in real-world scenarios. An image was classified as “healthy” if no
 669 disease-related BB were predicted. The model correctly identified 431 out of
 670 500 healthy images, leading to a false positive rate of 13.8%. This demonstrates
 671 that FD-TR can effectively differentiate healthy fruits from unhealthy ones in
 672 most cases. Figure 11 highlights three failure modes where natural fruit fea-
 673 tures were mistakenly classified as disease symptoms. In these cases, the model
 674 misinterpreted natural variations in fruit appearance, such as blemishes, color
 675 gradients, or developmental traits, as pathological indicators:

- 676 • Case (A): A healed scar on a citrus fruit (red arrow) was misclassified as
 677 a fungal infection (confidence: 0.47). The model failed to distinguish the
 678 scar’s shallow, textured appearance from active fungal lesions.
- 679 • Case (B): A young dragon fruit exhibiting natural tip browning (red arrow)
 680 was incorrectly flagged as infected with BFB, despite lacking characteristic
 681 water-soaked lesions.
- 682 • Case (C): A faint reddish patch on a young fig (red arrow) was predicted
 683 as a fungal spot, even though the coloration was uniform and confined to
 684 healthy epidermal tissue.

685 These examples revealed that the model’s false positives occurred not from
 686 complex background clutter or extreme lighting artifacts, but from everyday
 687 morphological and variations traits of healthy fruits that were not included in
 688 the training set. Such improvements would enhance the model’s robustness to
 689 real-world variability and reduce overfitting to disease-centric features.

690 6. Discussion

691 FD-TR model improves fruit disease detection by combining the Co-DETR
 692 training scheme with the DINO transformer model, multi-scale feature extrac-

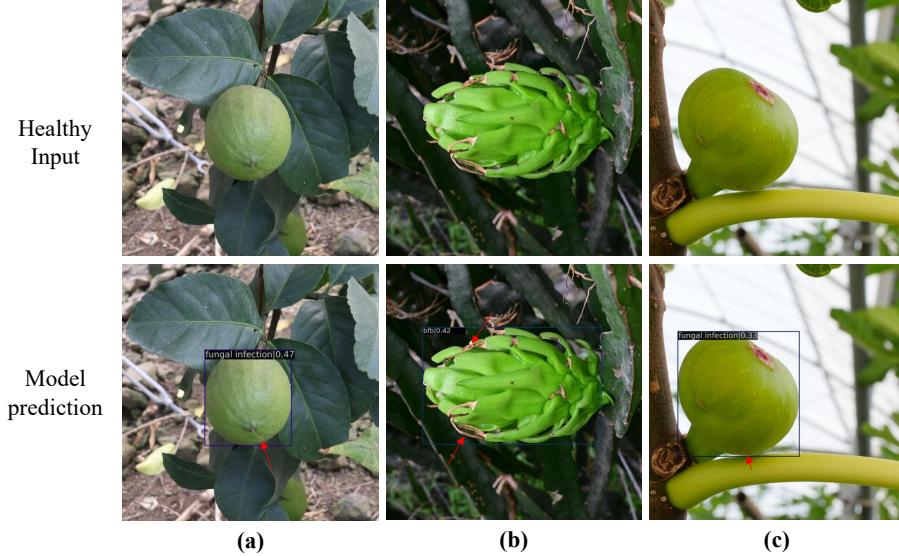


Figure 11: Samples of false positive prediction by the model for healthy fruit images.

tion, and attention mechanisms. Key model customization, including CIoU loss for precise BB, the LAMB optimizer for faster convergence, and BatchFormerV2 for scalable training, enhanced detection performance and efficiency for six fruit disease classes. FD-TR’s end-to-end design and integrated data augmentation improved robustness under diverse real-world scenarios, such as lighting and angles.

The experiment results showed that targeted customization improved detection mAP from 0.81 to 0.89. FD-TR also outperformed YOLOv8 (0.80) and Deformable DETR (0.74). With precision and recall rates of 0.86 and 0.87, respectively, it demonstrated robust generalization across diverse disease symptoms, scales, and environmental conditions. These capabilities are crucial for real-world agricultural settings, where early and accurate detection is crucial for effective intervention and crop protection. Furthermore, its attention-based interpretability via feature discriminability scores and deformable attention weights provided transparent insights into decision-making. The evaluation on healthy fruit images, as introduced in Section 5.8, demonstrated the model’s

⁷⁰⁹ potential to operate effectively in real-world settings where both diseased and
⁷¹⁰ healthy fruits are present. Although, a false positive rate of 13% on healthy
⁷¹¹ samples was promising, the misclassifications highlighted a limitation in the
⁷¹² current training data, which lacked explicit healthy examples.

⁷¹³ 7. Conclusions

⁷¹⁴ This research introduces an enhanced end-to-end transformer-based fruit
⁷¹⁵ disease recognition model that can be applied to real-life disease management
⁷¹⁶ systems. The dataset used to train the model consists of 81,000 images of
⁷¹⁷ six different fruit diseases. The proposed FD-TR model demonstrates high
⁷¹⁸ detection performance on the dataset compared to state-of-the-art models such
⁷¹⁹ as YOLOv8, DINO, and Deformable DETR. FD-TR is based on the DINO
⁷²⁰ transformer model with an improved Co-DETR training scheme and additional
⁷²¹ components like CIoU loss, the LAMB optimizer, and BatchFormerV2. These
⁷²² improvements contribute to the model's improved detection capabilities and
⁷²³ faster convergence during training. Therefore, FD-TR model not only improves
⁷²⁴ the accuracy of predictions but also achieves robust performance in various
⁷²⁵ experiments.

⁷²⁶ Moreover, FD-TR model's ability to maintain high performances on diverse
⁷²⁷ testing scenarios demonstrates its generalization ability and reliability. Even in
⁷²⁸ challenging cases, such as images affected by poor lighting or blurring, the model
⁷²⁹ provides correct and robust predictions. The attention mechanism of the trans-
⁷³⁰ former allows the model to focus on relevant disease features, which reduces false
⁷³¹ predictions. In addition, the unique multi-scale attention map extracted from
⁷³² the transformer offers experts/farmers valuable insights into how the model de-
⁷³³ tects and highlights disease-related areas. FD-TR model represents a significant
⁷³⁴ advancement in automated disease detection and offers substantial potential to
⁷³⁵ improve agricultural productivity and disease management in modern farming.

⁷³⁶ While FD-TR model demonstrates strong performance in detecting fruit dis-
⁷³⁷ eases, several limitations persist. One of the main limitations is the reliance on

738 a dataset with a limited number of disease classes, which fails to capture the full
 739 diversity of fruit diseases and environmental conditions. Moreover, the model's
 740 performance could be further optimized in challenging environmental conditions,
 741 where it occasionally struggles to detect diseases accurately. In the future, the
 742 dataset can be expanded to include more diverse conditions and disease types to
 743 improve the model's generalizability. In addition, techniques like multi-modal
 744 data integration, which analyze data from sensors such as infrared cameras or
 745 spectroscopy, can be considered for further development and improvement. Fi-
 746 nally, the model optimization on edge/mobile devices is a critical future work
 747 to enable real-time, on-field disease detection, especially in resource-constrained
 748 environments. This would involve exploring lightweight backbones and model
 749 compression techniques to reduce computational demands for edge devices.

750 **Author contributions**

751 **Yanfen Li:** Writing – original draft. **Muhammad Fayaz:** Data cura-
 752 tion. **Sufyan Danish:** Visualization. **Lilia Tightiz:** Investigation. **Hanxi-**
 753 **ang Wang:** Conceptualization, Methodology. **Tan N. Nguyen:** Supervision,
 754 Validation. **L. Minh Dang:** Methodology, Writing – review & editing.

755 **References**

- 756 [1] FAO, How to Feed the World in 2050, 2009. https://www.fao.org/fileadmin/templates/wsfs/docs/expert_paper/How_to_Feed_the_World_in_2050.pdf, accessed 2023-07-15.
- 759 [2] O. Turnbull, M. Homer, H. Ensaff, Food insecurity: Its prevalence and
 760 relationship to fruit and vegetable consumption, Journal of human nutrition
 761 and dietetics 34 (2021) 849–857.
- 762 [3] M. Dang, H. Wang, Y. Li, T.-H. Nguyen, L. Tightiz, N. Xuan-Mung, T. N.
 763 Nguyen, Computer vision for plant disease recognition: A comprehensive
 764 review, The Botanical Review (2024) 1–61.

- 765 [4] V. Singh, N. Sharma, S. Singh, A review of imaging techniques for plant
 766 disease detection, *Artificial Intelligence in Agriculture* 4 (2020) 229–242.
- 767 [5] L. C. Ngugi, M. Abelwahab, M. Abo-Zahhad, Recent advances in image
 768 processing techniques for automated leaf pest and disease recognition—a
 769 review, *Information processing in agriculture* 8 (2021) 27–51.
- 770 [6] S. A. Gaikwad, K. S. Deore, M. K. Waykar, P. R. Dudhane, G. Sorate,
 771 Fruit disease detection and classification, *International Research Journal
 772 of Engineering and Technology* 4 (2017) 1151–1154.
- 773 [7] B. J. Samajpati, S. D. Degadwala, Hybrid approach for apple fruit diseases
 774 detection and classification using random forest classifier, in: 2016 Interna-
 775 tional conference on communication and signal processing (ICCSP), IEEE,
 776 2016, pp. 1015–1019.
- 777 [8] I. M. Nasir, A. Bibi, J. H. Shah, M. A. Khan, M. Sharif, K. Iqbal, Y. Nam,
 778 S. Kadry, Deep learning-based classification of fruit diseases: An applica-
 779 tion for precision agriculture, *Comput. Mater. Contin* 66 (2021) 1949–1962.
- 780 [9] N. Yao, F. Ni, M. Wu, H. Wang, G. Li, W.-K. Sung, Deep learning-
 781 based segmentation of peach diseases using convolutional neural network,
 782 *Frontiers in Plant Science* 13 (2022) 876357.
- 783 [10] L. M. Dang, S. I. Hassan, I. Suhyeon, A. kumar Sangaiah, I. Mehmood,
 784 S. Rho, S. Seo, H. Moon, Uav based wilt detection system via convolu-
 785 tional neural networks, *Sustainable Computing: Informatics and Systems*
 786 28 (2020) 100250.
- 787 [11] H. Sun, J. Xue, Y. Song, P. Wang, Y. Wen, T. Zhang, Detection of fruit tree
 788 diseases in natural environments: A novel approach based on stereo camera
 789 and deep learning, *Engineering Applications of Artificial Intelligence* 137
 790 (2024) 109148.

- 791 [12] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, M. Shah, Trans-
 792 formers in vision: A survey, *ACM computing surveys (CSUR)* 54 (2022)
 793 1–41.
- 794 [13] A. Vaswani, Attention is all you need, *Advances in Neural Information
 795 Processing Systems* (2017).
- 796 [14] C. Zhu, W. Ping, C. Xiao, M. Shoeybi, T. Goldstein, A. Anandkumar,
 797 B. Catanzaro, Long-short transformer: Efficient transformers for language
 798 and vision, *Advances in neural information processing systems* 34 (2021)
 799 17723–17736.
- 800 [15] N. Kitaev, L. Kaiser, A. Levskaya, Reformer: The efficient transformer, in:
 801 International Conference on Learning Representations, 2020.
- 802 [16] C. Wang, Y. Chen, S. Zhang, Q. Zhang, Stock market index prediction us-
 803 ing deep transformer model, *Expert Systems with Applications* 208 (2022)
 804 118128.
- 805 [17] T. Kehinde, O. J. Adedokun, A. Joseph, K. M. Kabirat, H. A. Akano,
 806 O. A. Olanrewaju, Helformer: an attention-based deep learning model for
 807 cryptocurrency price forecasting, *Journal of Big Data* 12 (2025) 81.
- 808 [18] H. Chen, Z. Qi, Z. Shi, Remote sensing image change detection with trans-
 809 formers, *IEEE Transactions on Geoscience and Remote Sensing* 60 (2021)
 810 1–14.
- 811 [19] M. Orabi, K. P. Tran, P. Egger, S. Thomassey, Anomaly detection in smart
 812 manufacturing: An adaptive adversarial transformer-based model, *Journal
 813 of Manufacturing Systems* 77 (2024) 591–611.
- 814 [20] K. M. Choromanski, V. Likhoshesterov, D. Dohan, X. Song, A. Gane, T. Sar-
 815 los, P. Hawkins, J. Q. Davis, A. Mohiuddin, L. Kaiser, et al., Rethinking
 816 attention with performers, in: International Conference on Learning Rep-
 817 resentations, 2020.

- 818 [21] P. S. Thakur, P. Khanna, T. Sheorey, A. Ojha, Vision transformer for plant
 819 disease detection: Plantvit, in: International Conference on Computer
 820 Vision and Image Processing, Springer, 2021, pp. 501–511.
- 821 [22] Z. Zong, G. Song, Y. Liu, Detrs with collaborative hybrid assignments
 822 training, in: Proceedings of the IEEE/CVF international conference on
 823 computer vision, 2023, pp. 6748–6758.
- 824 [23] X. Xie, Y. Ma, B. Liu, J. He, S. Li, H. Wang, A deep-learning-based real-
 825 time detector for grape leaf diseases using improved convolutional neural
 826 networks, *Frontiers in plant science* 11 (2020) 751.
- 827 [24] J. You, K. Jiang, J. Lee, Deep metric learning-based strawberry disease
 828 detection with unknowns, *Frontiers in Plant Science* 13 (2022) 891785.
- 829 [25] S. F. Syed-Ab-Rahman, M. H. Hesamian, M. Prasad, Citrus disease detec-
 830 tion and classification using end-to-end anchor-based deep learning model,
 831 *Applied Intelligence* 52 (2022) 927–938.
- 832 [26] Z. Huang, X. Jiang, S. Huang, S. Qin, S. Yang, An efficient convolutional
 833 neural network-based diagnosis system for citrus fruit diseases, *Frontiers*
 834 in *Genetics* 14 (2023) 1253934.
- 835 [27] K. Nosiba Arifin, S. Akter Rupa, M. M. Anwar, I. Jahan, Lemon and or-
 836 ange disease classification using cnn-extracted features and machine learn-
 837 ing classifier, in: Proceedings of the 3rd International Conference on Com-
 838 puting Advancements, 2024, pp. 154–161.
- 839 [28] S. Aksoy, P. Demircioglu, I. Bogrekci, Web-based ai system for detecting
 840 apple leaf and fruit diseases., *AgriEngineering* 7 (2025).
- 841 [29] D. Faye, I. Diop, N. Mbaye, D. Dione, M. M. Diedhiou, Mango fruit dis-
 842 eases severity estimation based on image segmentation and deep learning,
 843 *Discover Applied Sciences* 7 (2025) 1–12.

- 844 [30] Y. He, N. Zhang, X. Ge, S. Li, L. Yang, M. Kong, Y. Guo, C. Lv, Passion
 845 fruit disease detection using sparse parallel attention mechanism and
 846 optical sensing, Agriculture 15 (2025) 733.
- 847 [31] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner,
 848 M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An
 849 image is worth 16x16 words: Transformers for image recognition at scale,
 850 in: International Conference on Learning Representations, 2020.
- 851 [32] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, S. Zagoruyko,
 852 End-to-end object detection with transformers, in: European conference
 853 on computer vision, Springer, 2020, pp. 213–229.
- 854 [33] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, J. Dai, Deformable detr: Deformable
 855 transformers for end-to-end object detection, in: International Conference
 856 on Learning Representations, 2020.
- 857 [34] F. Li, H. Zhang, S. Liu, J. Guo, L. M. Ni, L. Zhang, Dn-detr: Accelerate
 858 detr training by introducing query denoising, in: Proceedings of
 859 the IEEE/CVF Conference on Computer Vision and Pattern Recognition,
 860 2022, pp. 13619–13627.
- 861 [35] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, L. Zhang, Dab-
 862 detr: Dynamic anchor boxes are better queries for detr, in: International
 863 Conference on Learning Representations, 2022.
- 864 [36] D. Singh, N. Jain, P. Jain, P. Kayal, S. Kumawat, N. Batra, Plantdoc: A
 865 dataset for visual plant disease detection, in: Proceedings of the 7th ACM
 866 IKDD CoDS and 25th COMAD, 2020, pp. 249–253.
- 867 [37] D. Hughes, M. Salathé, et al., An open access repository of images on
 868 plant health to enable the development of mobile disease diagnostics, arXiv
 869 preprint arXiv:1511.08060 (2015).

- 870 [38] B. Pakruddin, R. Hemavathy, A comprehensive standardized dataset of
 871 numerous pomegranate fruit diseases for deep learning, Data in Brief 54
 872 (2024) 110284.
- 873 [39] H. T. Rauf, B. A. Saleem, M. I. U. Lali, M. A. Khan, M. Sharif, S. A. C.
 874 Bukhari, A citrus fruits and leaves dataset for detection and classification
 875 of citrus diseases through machine learning, Data in brief 26 (2019) 104340.
- 876 [40] NIA, AI Hub dataset, 2023. <https://www.aihub.or.kr>, accessed 2023-05-
 877 11.
- 878 [41] A. Ciofini, F. Negrini, R. Baroncelli, E. Baraldi, Management of post-
 879 harvest anthracnose: current approaches and future perspectives, Plants
 880 11 (2022) 1856.
- 881 [42] J. Daley, T. C. Wehner, Screening for bacterial fruit blotch resistance in
 882 watermelon fruit, Crop Science 61 (2021) 1228–1240.
- 883 [43] M. Cabedo-López, J. Cruz-Miralles, D. Peris, M. V. Ibáñez-Gual, V. Flors,
 884 J. A. Jaques, The response of citrus plants to the broad mite polyphago-
 885 tarsonemus latus (banks)(acari: Tarsonemidae), Agricultural and Forest
 886 Entomology 23 (2021) 411–419.
- 887 [44] J. Haran, G. J. Kergoat, B. A. de Medeiros, Most diverse, most neglected:
 888 weevils (coleoptera: Curculionoidea) are ubiquitous specialized brood-site
 889 pollinators of tropical flora, Peer Community Journal 3 (2023).
- 890 [45] A. Gallardo-Ferrand, L. A. Escudero-Colomar, J. Avilla, D. Bosch-Serra,
 891 Thrips (thysanoptera: Terebrantia) in nectarine orchards in north-east
 892 spain: species diversity and fruit damage, Insects 15 (2024) 699.
- 893 [46] A. Goudarzi, S. Samavi, M. Amiri Mazraie, Z. Majidi, Fungal pathogens
 894 associated with pre-and postharvest fruit rots of mango in southern iran,
 895 Journal of Phytopathology 169 (2021) 545–555.

- 896 [47] K. Wada, Labelme: Image Polygonal Annotation with Python, 2019. URL:
 897 <https://github.com/wkentaro/labelme>. doi:10.5281/zenodo.5711226.
- 898 [48] L. Zhang, G. Zhou, C. Lu, A. Chen, Y. Wang, L. Li, W. Cai, Mmdgan:
 899 A fusion data augmentation method for tomato-leaf disease identification,
 900 Applied Soft Computing 123 (2022) 108969.
- 901 [49] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask r-cnn, in: Proceedings
 902 of the IEEE international conference on computer vision, 2017, pp. 2961–
 903 2969.
- 904 [50] M. Zhang, S. Xu, W. Song, Q. He, Q. Wei, Lightweight underwater ob-
 905 ject detection based on yolo v4 and multi-scale attentional feature fusion,
 906 Remote Sensing 13 (2021) 4706.
- 907 [51] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A. C.
 908 Berg, Ssd: Single shot multibox detector, in: Computer Vision–ECCV
 909 2016: 14th European Conference, Amsterdam, The Netherlands, October
 910 11–14, 2016, Proceedings, Part I 14, Springer, 2016, pp. 21–37.
- 911 [52] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dol-
 912 lár, C. L. Zitnick, Microsoft coco: Common objects in context, in: Com-
 913 puter Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland,
 914 September 6–12, 2014, Proceedings, Part V 13, Springer, 2014, pp. 740–755.
- 915 [53] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman,
 916 The pascal visual object classes (voc) challenge, International journal of
 917 computer vision 88 (2010) 303–338.
- 918 [54] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. Ni, H.-Y. Shum, Dino:
 919 Detr with improved denoising anchor boxes for end-to-end object detection,
 920 in: The Eleventh International Conference on Learning Representations,
 921 2022.
- 922 [55] Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song,
 923 J. Demmel, K. Keutzer, C.-J. Hsieh, Large batch optimization for deep

- learning: Training bert in 76 minutes, arXiv preprint arXiv:1904.00962 (2019).
- [56] Z. Hou, B. Yu, D. Tao, Batchformer: Learning to explore sample relationships for robust representation learning, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 7256–7266.
- [57] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, S. Savarese, Generalized intersection over union: A metric and a loss for bounding box regression, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 658–666.
- [58] K.-a. Tessera, S. Hooker, B. Rosman, Keep the gradients flowing: Using gradient flow to study sparse network optimization, arXiv preprint arXiv:2102.01670 (2021).
- [59] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, H. Lipson, Understanding neural networks through deep visualization, arXiv preprint arXiv:1506.06579 (2015).
- [60] G. Jocher, A. Chaurasia, J. Qiu, Ultralytics YOLO, 2023. URL: <https://github.com/ultralytics/ultralytics>.
- [61] Kaggle, PlantVillage for object detection YOLO, 2024. <https://www.kaggle.com/datasets/sebastianpalaciob/plantvillage-for-object-detection-yolo>, accessed 2025-06-25.
- [62] H. Wang, Y. Li, L. M. Dang, H. Moon, An efficient attention module for instance segmentation network in pest monitoring, Computers and Electronics in Agriculture 195 (2022) 106853.
- [63] S. Zhang, L. Wen, X. Bian, Z. Lei, S. Z. Li, Single-shot refinement neural network for object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 4203–4212.

- 951 [64] Kaggle, PlantVillage for object detection, 2023.
- 952 [https://www.kaggle.com/datasets/sebastianpalaciob/
953 plantvillage-for-object-detection-yolo/](https://www.kaggle.com/datasets/sebastianpalaciob/plantvillage-for-object-detection-yolo/), accessed 2023-05-11.
- 954 [65] J. Cao, H. Cholakkal, R. M. Anwer, F. S. Khan, Y. Pang, L. Shao, D2det:
955 Towards high quality object detection and instance segmentation, in: Pro-
956 ceedings of the IEEE/CVF conference on computer vision and pattern
957 recognition, 2020, pp. 11485–11494.
- 958 [66] Y. Zhu, C. Zhao, J. Wang, X. Zhao, Y. Wu, H. Lu, Couplenet: Coupling
959 global structure with local parts for object detection, in: Proceedings of the
960 IEEE international conference on computer vision, 2017, pp. 4126–4134.

Declaration of interests

- The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.
- The author is an Editorial Board Member/Editor-in-Chief/Associate Editor/Guest Editor for *[Journal name]* and was not involved in the editorial review or the decision to publish this article.
- The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

