

# A prompt based approach for apple disease classification

xing sheng

Shandong Normal University

Yangyang Zhang

Shandong Normal University

Chen Lyu (✉ [lvchen@sdu.edu.cn](mailto:lvchen@sdu.edu.cn))

Shandong Normal University <https://orcid.org/0000-0002-5044-1459>

---

## Research Article

**Keywords:** prompt, Disease classification, Transfer Learning, WPM, 0-shot

**Posted Date:** February 21st, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1308089/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

## RESEARCH

# A prompt based approach for apple disease classification

Xing Sheng, Yangyang Zhang and Chen Lyu\*

\*Correspondence:  
lvchen@sdnu.edu.cn  
School of Information Science and  
Engineering, Shandong Normal  
University, Jinan, China  
Full list of author information is  
available at the end of the article

## Abstract

**Background:** Apples occupy a large part of agricultural production as a fruit with a high yield and also a high nutritional value. However, diseases of the fruit and leaves of apples seriously affect the quality and yield of apples. In the past, people had to rely on their own experience to control apple diseases, however, this approach was poorly accurate, inefficient and did not meet the requirements of fruit farmers. Many current methods are based on convolutional neural networks, but convolutional neural networks usually require a large amount of labelled data to train the network, and datasets in the agricultural field can hardly meet this requirement.

**Results:** To solve this problem, this paper introduces zero-times learning, which can achieve equally good results even if the test object is a dataset that has never been seen before. Specifically, we give a short description of an image as a prompt word according to its category in the public dataset, form a graphical pair of the image to be trained and the corresponding prompt word, feed them into a deep convolutional neural network, and then pre-train it on a large public dataset, and migrate it directly to our dataset after the training is completed, saving a lot of resources. In addition, we also propose a new attention module WPM (Weighted-Pooling Module) to deeply mine feature vectors by combining weighted pooling operations with fully connected operations and activation functions. Through extensive experiments, we validate the effectiveness of the proposed approach of combining zero-times learning with prompt words and achieve good results on our own collected field dataset.

**Conclusions:** Our work provides new ideas and resource savings for disease classification tasks in agriculture.

**Keywords:** prompt; Disease Classification; Transfer Learning; WPM; 0-shot

1  
2

## Introduction

The apple is one of the four most important fruits in the world and the most important deciduous fruit tree in China; It ranks first in the country in terms of area and production. In Shandong province alone, the apple area has reached 800,000 hectares with a production of over 4 million tonnes, surpassing that of the USA and the Southern Hemisphere. However, it has been reported that there are as many as 100 disease problems of apples, which can be divided into leaf, branch, fruit and root diseases according to their location, the most serious and common of which are diseases of the leaves and fruit of apples. Only apple fruit ring rot and fruit anthracnose, the number of apples lost each year is above 5%, and in serious years it is as high as 60%.

In the past, the identification of diseases of apple fruit and leaves relied heavily on farmers conducting field visits and judging the type of disease based on experience. However, the lack of appropriate knowledge of some diseases among farmers has led to poor identification and assessment of the type

14

15 and severity of apple diseases and vague diagnostic criteria for diseases, resulting in failure to con-  
16 trol diseases in a timely and reasonable manner. It is therefore important to identify and control  
17 apple diseases effectively, rationally and accurately, and to research and apply modern identifica-  
18 tion techniques. The use of computer vision-based technology for disease identification is a major  
19 way of achieving this goal. Not only can computer vision technology quickly and accurately obtain  
20 information about apple diseases, but it can also select the appropriate control method according  
21 to the severity of the disease. This will greatly save manpower and material resources, improve the  
22 efficiency of production and reduce costs.

23 Machine learning based methods have been commonly used for disease identification of fruits.  
24 Mohan *et al.* [1] used KNN and SVM to classify brown spot disease, leaf blast disease and bacterial  
25 blight disease of paddy plants with good results. Mokhtar *et al.* [2] used wavelet transform technique  
26 in combination with support vector machine and alternating kernel function to detect and identify  
27 diseases of tomato leaves and finally achieved 99.5% accuracy. Sindhuja *et al.* [3] used principal  
28 component analysis (PCA) on citrus yellow dragon disease pretreatment dataset and then used  
29 linear discriminant analysis, quadratic discriminant analysis and K-nearest neighbor methods to  
30 model and classify the dataset, and finally obtained an overall accuracy of 98%. Arivazhagan *et al.*  
31 [4] first built a color transformation structure of HIS on the input RGB images, then used  
32 specific thresholds to mask and remove green pixels, then performed a segmentation process to  
33 count their texture information, and finally used SVM for classification, achieving an accuracy of  
34 94% of accuracy, and experimental results on a database of approximately 500 plant leaves confirmed  
35 the robustness of the method. However, although the research on plant pest recognition based on  
36 traditional image processing techniques has achieved certain results and the recognition accuracy of  
37 diseases is high, there are also shortcomings and limitations: the research process is tedious, relies too  
38 much on manually designed feature extraction methods, is highly subjective and time-consuming,  
39 etc. It cannot be adapted to practical application scenarios with more complex backgrounds and  
40 cannot meet the complex situations in practical applications, etc.

41 With the continuous development of deep learning, more and more research has applied deep  
42 learning to agricultural disease detection. Oppenheim *et al.* [5] collected 400 potato photos of different  
43 sizes, shapes and tones under different lighting conditions indoors, and by adding several new dropout  
44 layers behind the VGG [6] network to deal with the overfitting problem, they finally obtained the  
45 best Yusuke [7] proposed a convolutional neural network-based plant disease detection system using  
46 800 cucumber leaf images taken in the field to train the convolutional neural network for detecting  
47 disease infection in two cucumber plants, and eventually achieved an average accuracy of 94.9%  
48 under a 4-fold cross-validation strategy. Fuentes *et al.* [8] proposed a deep learning-based approach  
49 to detect diseases in tomato plants using images taken by cameras of different resolutions, combining  
50 detectors such as Faster R-CNN with deep feature extractors (VGG and RESNet [9]) and proposing  
51 a method based on local and global class annotation and data enhancement to improve accuracy and  
52 reduce the number of false positives during training, ultimately achieving good results. Ferentinos *et al.*  
53 [10] developed a specialized deep learning model based on a specific convolutional neural network  
54 architecture and tested it on a publicly available dataset (87,848 images in the dataset, with photos  
55 taken both in controlled laboratory conditions and in the field), achieving 99.53% accuracy. Liu *et al.*  
56 [11] proposed a deep learning model based on two networks, VGG16 and ResNet50, to identify  
57 the species of large chrysanthemums, trained on a balanced dataset constructed from 14,000 images

58 of 103 cultivars, and ultimately achieved a top-5 accuracy of 98%. Although the convolutional neural  
59 network-based approach achieved very good results, the dataset and the time and resources required  
60 to train the convolutional neural network were too large for a plant disease dataset in agriculture to  
61 meet the requirements of a deep convolutional neural network.

62 To address these issues, in this paper we introduce zero-times learning. The basic idea of zero-  
63 times learning is to give machines the ability to reason so that the models we train can classify  
64 models that have never been seen before, achieving true “artificial intelligence”. The definition of  
65 zero learning can be expressed as follows: given a labelled training instance  $D$  belonging to a seen  
66 category  $C$ , the goal of zero learning is to learn a classifier  $f(-):X \rightarrow U$  that classifies a test instance  
67  $X$  into a category  $U$  that it has not seen before. The label spaces covered by the training and test  
68 instances are disjoint. Therefore, zero-times learning is a sub-domain of migration learning [12]. In  
69 transfer learning, the source domain and the knowledge contained in the source task are transferred  
70 to the target domain to learn the model in the task goal. According to [12, 13], transfer learning  
71 can be classified as homogeneous transfer learning and heterogeneous transfer learning depending  
72 on whether the feature space and label space in the source and target domains/tasks are the same.  
73 In zero-times learning, the original label space is the visible class set, while the target label space  
74 is the invisible class set. Therefore, zero-times learning belongs to heterogeneous transfer learning.  
75 Since there are no labelled instances available in the invisible class set, some auxiliary information  
76 is needed to solve this problem. This auxiliary information should contain information about all  
77 invisible classes, which is also to ensure that the corresponding auxiliary information is provided for  
78 each invisible class. At the same time, the auxiliary information should be related to the instances  
79 in the feature space, this is to ensure that the auxiliary information is available.

80 In the existing work, the approach to incorporating supporting information is influenced by the  
81 way humans learn about the world. Humans can learn zero times with the help of some semantic  
82 background knowledge. For example, with the a priori knowledge that “wolves look like dogs, but  
83 their tails are short and thick and often hang back between their hind limbs”, we can recognize a  
84 wolf even if we have not seen one, provided we know what a dog looks like and what a dog’s tail  
85 looks like. The available auxiliary information is usually semantic, containing visible classes and  
86 some invisible classes. The approach taken in this paper is based on such semantic information in  
87 the form of prompt words. In particular, it should be noted that our approach differs from [14] in  
88 that the textual content in the method proposed in [14] for combining with text is described in  
89 detail and is far removed from our prompt words, and that it processes text and images separately  
90 using two models, only putting them together when the final output is produced. In contrast, our  
91 approach is to input the prompts and images into the model together as an image-text pair during  
92 training.

## 93 **Material and methods**

### 94 **Image acquisition and material**

95 The training set in the dataset we use is a public dataset and the test set is our own collected dataset.  
96 The training sets include ImageNet [15], PlantVillage [16] and PlantDoc [17]. The ImageNet dataset  
97 started in 2009 and was created by Professor Feifei Li and others. It has a total of 14,197,122 images  
98 and a total of 21,841 categories, with large categories including animals, birds, machines, flowers,  
99 food, fruits and 1000 other species, with roughly 1000 images per category. PlantVillage is a publicly

100 available dataset for testing machine learning plant disease detection algorithms. It contains 38 crop-  
 101 disease pairs of 14 plant species, with a total of 54,299 images. All images were taken indoors with  
 102 a single background image. The PlantDoc dataset is a dataset of crop disease images, manually  
 103 annotated with images acquired online, and includes 27 disease categories (10 healthy types, 17  
 104 disease types) for 13 plants, with a total of 2598 images for image classification and target detection.  
 105 The provider of the data for our test set is the Shandong Academy of Agricultural Sciences in Jinan,  
 106 Shandong Province, China. The test set contains a total of 1204 images of apple leaves and fruits,  
 107 including 474 images of apple fruits and 730 images of apple leaves. The images were taken with  
 108 a mobile phone under real production conditions, at a resolution of 3456\*4608, with non-uniform  
 109 light intensity (none of the images were taken with the flash on), at different angles and distances,  
 110 etc. The images contain a lot of extraneous background information such as other apples, leaves,  
 111 branches and sky. The dataset used in this paper includes normal leaves and fruits, of which there  
 112 are six types, all in .jpg format. As shown in Table 1, there are three types of apple fruit: healthy  
 113 fruit, apple fruit ring rot and apple fruit anthracnose, and three types of apple leaves: healthy leaves,  
 114 apple anthracnose leaf blight and apple leaf rust. Some of the images in the dataset are shown in  
 Figure 1.

Dataset.pdf

Figure 1: Image of part of the data in the dataset.

115  
 116 The overall flow of our method is shown in Figure 2. First, we collect disease samples of leaves  
 117 and fruits of apples in orchards under the guidance of an expert to label and classify them. The  
 118 collected images were then subjected to a series of processes, such as image cropping, image contrast  
 119 stretching, grey level slicing, dynamic range compression, etc. They were set to a uniform size of  
 120 224\*224 during training, and then the images were subjected to normalization operations, etc. A  
 121 prompt word was then added to each category, such as “A photo of a healthy apple.”, and placed  
 122 together with the corresponding image as a picture-text pair, trained in the same way as Clip. Once  
 123 the training was completed, the weights were tested directly on our own dataset without any changes  
 124 to the weights file. Inspired by [25], we trained the model under a total of two prompt words: “A  
 125 photo of a {label}” and “A photo of a {label}, a type of XX. We trained on three common datasets:  
 126 ImageNet, PlantVillage, and PlantDoc. Table S1 and S2 show the model trained on “A photo of a  
 127 {label}” and “A photo of a {label}, a type of XX.” prompt. In particular, because the image scenes  
 128 in ImageNet are complex and cannot be described simply by the two prompt words mentioned

Table 1: Apple fruit and leaf dataset display.

ClassID	Type	Number of Sample
0	Healthy Apple	157
1	Ring Rot Apple	172
2	Anthracnose Apple	145
3	Healthy Leaf	67
4	Anthracnose Leaf	435
5	Rust Leaf	228

129 above, in practice we use more complex prompt words, such as “A dark photo of a { }.” or “A black  
 130 and white photo of a { }.”. The model was trained and tested on a Linux server with an Intel(R)  
 131 Core(TM) i9-10980XE CPU (128GB RAM) and accelerated with two Nvidia GeForce 2080Ti GPUs  
 132 (24G RAM). The model was implemented in the pytorch 1.8 open source framework with python  
 133 version 3.7. The learning rate was set to 4e-6,eps to 1e-6 for Clip-RN50 and 5e-6,eps to 2e-6 for both  
 134 Clip-ViT-B/16 and Clip-ViT-B/32, and the weight decay was 0.1, and the optimization functions  
 used were all AdamW.

FrameWork.pdf

Figure 2: Overall flow chart of the method.

135

### 136 *Evaluation Metric*

137 We use accuracy, precision, recall and F1-score as our evaluation metrics.

### 138 *Accuracy.*

Accuracy refers to the number of correctly classified samples as a proportion of the number of samples determined to be positive by the classifier, and is often used to evaluate the quality of the results, which can be expressed as:

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

139 Where TP represents true cases, FP represents false positive cases, TN represents true negative  
 140 cases and FN represents false negative cases.

### 141 *Precision.*

Precision is often used to evaluate the completeness of the results and it can be expressed as:

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

### 142 *Recall.*

Recall refers to the proportion of correctly classified positive samples to true positive samples and is often used to indicate coverage and can be expressed as:

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

*F1-score.* The F1-score is a weighted summed average of precision and recall proposed to balance precision and recall, which is as suggestive of precision and recall as possible while expecting the difference between them to be as small as possible, and can be expressed as:

$$\text{F1} = 2 \frac{P * R}{P + R} \quad (4)$$

143 *Prompt*

144 With the continued research and development of deep convolutional networks, the learning approach  
 145 to natural language processing has begun to gradually shift from a fully supervised approach to a  
 146 pre-trained-fine-tuned model [18, 19, 20], in which models with fixed architectures are pre-trained  
 147 as language models that predict the probability of observed textual data. Due to the very large  
 148 and rich amount of raw text on which the language model is trained, the model can learn powerful  
 149 generic features of the language it models during the training process. The trained language model  
 150 is then applied to downstream tasks by fine-tuning it with additional parameters and using task-  
 151 specific objective functions. In this model, the focus shifts to goal engineering, designing training  
 152 goals for the pre-training and fine-tuning phases, which can lead to better pre-trained models for  
 153 text summary pre-training [21]. And by now (2022), the pre-training-fine-tuning model is being  
 154 replaced by a pre-training-prompt word-fine-tuning model. In this model, instead of adapting the  
 155 pre-trained language model to downstream tasks through goal engineering, the downstream task is  
 156 reformulated so that it solves these downstream tasks in the original language model with the help  
 157 of textual prompts. For example, when completing a fill-in-the-blank task, “I ate a fruit today”, we  
 158 can prompt “it was very \_” and have the language model fill in the blank with an adjective. With  
 159 appropriate prompts, we can get the desired output from the pre-trained language model [22, 23, 24].  
 This model has also been applied to the image domain.

Clip.pdf

Figure 3: Clip’s overall framework diagram.

160  
 161 Alec Radford et al. [25] proposed to use image-text pairs for pre-training (shown in Figure 3) and  
 162 migrate the completed training model to a new task to achieve zero learning. Specifically: Firstly,  
 163 a large number of unpurged pairs of text and image pairs are searched on the Internet, they are  
 164 labelled and encoded with text and image respectively, and the similarity between two pairs of text  
 165 and image in a batch is calculated by dot product to obtain a batch size by batch size similarity  
 166 matrix, and the similarity on the diagonal is the similarity value of the positive sample. Therefore, in  
 167 the training process, the optimization goal is to make the similarity value of the positive samples as  
 168 large as possible. In the inference of the image classification task, the model first needs to convert the  
 169 category labels into the same sentences as in the pre-training (where prompt is used), then get the  
 170 prompt words corresponding to the different categories and then input them into the test network to  
 171 form an image-text pair for zero-learning prediction, and output the corresponding categories after  
 172 the prediction is completed. After it was proposed, many studies have applied Clip to the image and  
 173 video domains, such as CoOp [26], Action Clip [27], Clip4Caption [28], Clip4Clip [29], etc.

174 **Deep learning based classification framework**

175 *Based on ResNet*

176 The proposed deep residual network is a milestone event in the history of convolutional neural  
 177 network images, which solves the problem of difficult training of deep convolutional neural network  
 178 models. Specifically, as the depth of the network increases, the network can perform more complex  
 179 feature extraction operations and theoretically achieve better results, but in the actual training  
 180 process, the deep network will suffer from degradation problems: the network accuracy saturates

181 or even decreases when the depth of the network increases. To solve this problem, He Keming et  
 182 al. [9] proposed ResNet, which solves the degradation problem through residual learning. When the  
 183 input is  $x$  its learned features are noted as  $H(x)$ , and residual learning is performed by means of  
 184  $H(x) = F(x) + x$  (as shown in Figure 4). The ResNet network improves on the network of VGG19  
 185 by adding residual units through a short-circuiting mechanism, and achieves good results. One of  
 the experimental models baseline we used was the RN50 model.

ResNet.pdf

Figure 4: Residual learning units.

186

### 187 *Based on Vision Transformer*

188 A detailed framework diagram of the Vision Transformer used in this paper is shown in Figure 5.  
 189 Alexey et al. [30] first applied the Transformer to the field of computer vision and proposed the  
 190 Vision Transformer model. Specifically, it first chunks the images and then spreads each image into  
 191 a one-dimensional vector; next, a linear transformation (i.e., a fully connected layer) is done on  
 192 each vector, called Patch Embedding. A vector is artificially added to the input vector as the final  
 193 classification vector. They are then fed into the Transformer's encoder for mask computation using  
 194 multi-headed self-attentiveness and feature mapping using an FFN (Feed-Forward Network), and the  
 195 processed feature vectors are fed into a Multilayer Perceptron (MLP) for classification. We used the  
 196 Vision Transformer trained ViT-B/16 and ViT-B/32 (B denotes Base, a relatively small amount of  
 197 data; 16 and 32 denote the input patch size of  $16 \times 16$  and  $32 \times 32$  respectively) models in our baseline  
 experiments.

ViT.pdf

Figure 5: Architecture diagram of Vision Transformer.

198

### 199 **Weighted-pooling Module**

200 As shown in Figure 6, a novel attention mechanism is proposed in order to reduce the number of  
 201 parameters and to facilitate the modelling of the information in the feature vector. This attention  
 202 mechanism uses a simple combination of a pooling layer with additional weights, a fully connected  
 layer, a batch normalization (BN) layer and a Sigmoid activation function layer.

WPM.pdf

Figure 6: Proposed weighted-pooling Module.

203

Specifically, for the feature vector  $Feature$ , the overall algorithmic process can be divided into two  
 parts: the first part is the feature extraction operation and the second part is the feature fusion  
 operation. The feature extraction operation can be formulated as follows:

$$Feature_1 = S(\text{ReLU}(FC((\text{Cat}(\gamma \text{Avg}(Feature), \beta \text{Std}(Feature)))))) \quad (5)$$

where Avg and Std denote AvgPool and StdPool operations respectively, Cat stands for connected operation and FC stands for fully connected operation, ReLU and S stand for ReLU activation function and Sigmoid function respectively,  $\gamma$  and  $\beta$  are two weights.

$$\text{Fused} = \text{Feature}_1 \otimes \text{Feature} \oplus \text{Feature} \quad (6)$$

204 where  $\otimes$  and  $\oplus$  represents the multiplication of elements and the addition of elements respectively,  
205 Fused denotes the feature vector of the total post-fusion.

## 206 Results and discussions

207 *Comparison using different datasets.* From Figure 7, we can see that ImageNet outperforms PlantVil-  
208 lage and PlantDoc in classifying the fruit of apples.

209 **Results and reasons:** This is because the ImageNet dataset has more diverse images, including  
210 a variety of fruits including apples, so it is better able to distinguish between the types of apples.  
211 The images in the PlantVillage and PlantDoc datasets are of plant leaves and lack images of fruit,  
so they are less generalizable to apple fruit.

Apple.pdf

Figure 7: Comparison of average accuracy on pre-trained models on different datasets for classification of apple fruits..

212 From Figure 8, we can see that the results on PlantDoc are much better than those on PlantVillage  
213 and ImageNet in terms of leaf classification of apples.

Leaf.pdf

Figure 8: Average accuracy comparison of pre-trained models on different datasets for the classification of apple leaves.

214 **Results and reasons:** This is due to the lack of plant leaf image data on ImageNet, whereas  
215 PlantDoc has leaf images of a wide range of plant diseases in the field, including apple leaf disease  
216 images, and is therefore better able to classify diseased apple leaves. As for PlantVillage, although  
217 it has a variety of plant disease foliage in its dataset, including apple disease foliage, its foliage was  
218 taken under controlled indoor conditions and lacks good generalization to images taken in the field.  
219

220 *Comparison of different prompt words.* As shown in Figure 9, the performance of our two trained  
221 prompt words on our apple fruit and leaf datasets after pre-training on both datasets on the four  
222 metrics illustrates that “A photo of a {}, a type of {}.” is better than “A photo of a {}.” overall  
223 better performance. The quantitative analysis tables using the different prompt words are shown in  
224 Table S1 and S2.

225 **Result and reason:** “A photo of a {}, a type of {}”. The phrase “a type of {}” contains more  
226 semantic information than the other prompt, which contains the category of the image being trained.  
227 This makes it easier for the model to combine the semantic information of the text with the feature  
228 information of the image, and the training results are better.

## Prompt.pdf

Figure 9: Plot of the average accuracy of the two prompt words on the two datasets (PlantVillage and PlantDoc) compared.

229 *Effectiveness of the prompt word method.* As we can see in Figure 10, the models trained on  
 230 ImageNet, PlantVillage and PlantDoc using the prompt word approach were tens of times more  
 231 effective in identifying diseases on apple fruit and leaves than the approach without the prompt  
 232 word. The quantitative analysis of the models without the use of prompt words is tabulated in Table  
 233 S3.

## Prompt Compare.pdf

Figure 10: Histogram of average accuracy on the three models with and without the prompt words.

234 **Results and reasons:** When training with prompt words, the model will first calculate the  
 235 similarity between the prompt word and the image, and will get the image and text with the  
 236 greatest similarity during training, and will be able to make good predictions when inferring with  
 237 the semantic information in the text, whereas a directly pre-trained model without this semantic  
 238 information from the text of the prompt word will be ineffective when faced with an image that it  
 239 has not seen before.

240 *Effectiveness of the proposed attention module.* As we can see in Table S4, after using our proposed  
 241 attention module, all of the pre-trained models showed varying degrees of improvement in the four  
 242 evaluation metrics.

243 **Results and reasons:** The effectiveness of our proposed attention module is illustrated by its  
 244 ability to better extract the associated features from the feature vector, improving its effectiveness  
 245 in classification tasks.

246 *Comparison of the validity of different models.* As shown in Figure 11, we trained a total of three  
 247 models, and from their comparison we can see that ViT-B/16 performed the best, ViT-B/32 the  
 second best, and then RN50.

## Model-Compare.pdf

Figure 11: Histogram comparing the average accuracy of each model tested on the Apple dataset.

248 **Results and reasons:** Firstly, the performance of the RN50-based model is less effective than  
 249 that of the Vision Transformer-based model, because the Vision Transformer divides the image into  
 250 multiple patches during training, which can learn the feature information between adjacent pixels  
 251 and help the network to distinguish the spatial information in the image, while the CNN-based  
 252 The CNN-based RN50 does not learn this information very well, and is therefore less effective. ViT-B/16  
 253 is more granular and learns more features than ViT-B/32, so it is more effective than ViT-B/32.  
 254

255 Table 2 shows the histogram analysis of the results of the three deep convolutional neural network  
 256 models on the Apple dataset after pre-training on the three dataset pre-training models, and provides  
 257 the mean  $\mu$  and standard deviation  $\delta$  of the four evaluation metrics on the Apple dataset.

Table 2: Box line plots of the four evaluation metrics (the four metrics in the box line plots, accuracy, recall, precision and f1-score, are indicated by red, orange, green and indigo blue respectively) on the ImageNet, PlantVillage and PlantDoc datasets experimented with the three models (RN50, ViT-16, ViT-32) and are given corresponding to the box line plots. The mean ( $\mu$ ) and standard deviation ( $\delta$ ) are given.

Models	Boxplot for Datasets			Region	Evaluation Metric				
	ImageNet	PlantVillage	PlantDoc		Accuracy	precision	recall	F1-score	
RN50	RN50-ImageNet.pdf	RN50-PlantVillage.pdf	RN50-PlantDoc.pdf	L	$\mu$	51.27	19.53	36.11	22.88
					$\delta$	35.79	15.8	6.81	15.86
				D	$\mu$	43.26	14.43	33.33	17.59
					$\delta$	37.27	12.42	6.35	11.91
				S	$\mu$	50.86	19.68	36.11	21.55
					$\delta$	44.87	19.26	6.81	19.04
ViT-16	ViT-16-ImageNet.pdf	ViT-16-PlantVillage.pdf	ViT-16-PlantDoc.pdf	L	$\mu$	54.98	24.89	38.89	29.39
					$\delta$	28.93	14.58	8.61	13.7
				D	$\mu$	66.12	25.32	38.89	28.29
					$\delta$	37.14	14.55	8.61	13.7
				S	$\mu$	59.71	24.97	33.89	28
					$\delta$	36.34	18.37	8.61	17.98
ViT-32	ViT-32-ImageNet.pdf	ViT-32-PlantVillage.pdf	ViT-32-PlantDoc.pdf	L	$\mu$	57.28	20.61	36.11	19.6
					$\delta$	27.61	9.76	6.81	10.88
				D	$\mu$	64.73	23.15	36.11	19.6
					$\delta$	29.34	14.15	6.81	10.88
				S	$\mu$	54.35	29.23	44.44	31.7
					$\delta$	42.29	36.68	27.22	35.25

258 As can be seen from the 9 images (RN50-ImageNet.pdf, RN50-PlantVillage.pdf, RN50-PlantDoc.pdf, ViT-  
259 16-ImageNet.pdf, ViT-16-PlantVillage.pdf, ViT-16-PlantDoc.pdf, ViT-32-ImageNet.pdf, ViT-32-  
260 PlantVillage.pdf, ViT-32-PlantDoc.pdf), the range of accuracy of the test results on all three datasets  
261 is very wide, which indicates that the results of the prompt word method fluctuate widely across  
262 the images, suggesting that it is not stable, and the results are sometimes good and bad. In terms of  
263 stability, the RN50 network has the lowest outlier value, but it also has the smallest median value,  
264 indicating that it is not as sensitive to the feature information of the images in the test set as the  
265 other two networks and is less effective, which is consistent with the reasons analyzed in the previ-  
266 ous section. Compared to ViT16, ViT32 has more outliers and the results vary more between the  
267 pre-trained models of different datasets, suggesting that ViT32 is more susceptible to the influence  
268 of the training images during training and lacks good generalization.

269 As for the recall, it is much more stable compared to the accuracy, but the mean and median  
270 values are significantly smaller than the accuracy, probably because the model did not learn enough  
271 feature information from the training dataset, a result that is consistent with the large gap between  
272 the pre-training and testing datasets we used.

273 The F1-score values used to balance accuracy and recall are more representative of the character-  
274 istics of the three models: RN50 has the smallest value while ViT-16 has a narrower F1-score width,  
275 indicating a smaller fluctuation range and a more stable effect and better generalization performance  
276 than ViT-32.

## 277 Conclusions

278 In this paper, we introduce prompt words into the agricultural disease classification task for the  
279 first time for zero learning in order to address the problem of small number of datasets and diffi-

280 cult collection in the agricultural domain. We used three deep convolutional neural network models:  
281 RN50, ViT-16 and ViT-32, using two prompt words: “A photo of a {}.” and “A photo of a {}, a  
282 type of {}.” Extensive experiments were done on three public datasets: ImageNet, PlantVillage and  
283 PlantDoc. Through these experiments, we verified the feasibility and effectiveness of the approach  
284 using prompt words (the approach using prompt words is tens of times better than training directly  
285 from public datasets such as ImageNet), and compared and analyzed the effectiveness of the two  
286 prompt words. In addition, in order to deeply mine the relationships between feature vectors and  
287 improve the effectiveness of deep convolutional neural networks, we also proposed a new attention  
288 module, WPM. Through a weighted pooling operation and some other operations, an average accu-  
289 racy improvement of 5.2% on apple fruits and 5.7% on apple leaves was achieved. Our work provides  
290 new ideas and resource savings for disease classification tasks in agriculture.

## 291 Appendix

292 Table S1. Results of experiments with RN50 and Vision Transformer-based models on apple fruit  
293 and leaves, with the prompt “A photo of a {label}.”.

294 Table S2. Results of experiments with RN50 and Vision Transformer-based models on apple fruit  
295 and leaves with the prompt “A photo of a {label}, a type of a .”.

296 Table S3. Quantitative analysis of pre-trained models on our dataset without the use of prompt  
297 words.

298 Table S4. Quantitative analysis on our dataset using the pre-trained model after attention (prompt:  
299 “A photo of a .”).

### 300 Ethics approval and consent to participate

301 Not applicable.

### 302 Consent for publication

303 Not applicable.

### 304 Availability of data and materials

305 Not applicable.

### 306 Competing interests

307 The authors declare no competing interests.

### 308 Funding

309 Not applicable.

### 310 Authors' contributions

311 XS carried out the simulations and generated the results. XS and YZ curated and annotated the datasets. CL conceived, designed, and  
312 supervised the work. XS and YZ wrote the manuscript. CL supervised the annotation and classification of diseases. CL supervised the data  
313 management and testing.

### 314 Acknowledgements

315 Not applicable.

### 316 Authors' information

317 Shandong Normal University, Jinan, China.

### 318 Author details

319 School of Information Science and Engineering, Shandong Normal University, Jinan, China.

## 320 References

- 321 1. Mohan, K.J., Balasubramanian, M., Palanivel, S.: Detection and recognition of diseases from paddy plant leaf images. *International*  
322 *Journal of Computer Applications* **144**(12) (2016)
- 323 2. Mokhtar, U., Ali, M.A., Hassenian, A.E., Hefny, H.: Tomato leaves diseases detection approach based on support vector machines. In:  
324 2015 11th International Computer Engineering Conference (ICENCO), pp. 246–250. IEEE

- 325 3. Sankaran, S., Mishra, A., Maja, J.M., Ehsani, R.: Visible-near infrared spectroscopy for detection of huanglongbing in citrus orchards.  
326 Computers and electronics in agriculture **77**(2), 127–134 (2011)
- 327 4. Arivazhagan, S., Shebiah, R.N., Ananthi, S., Varthini, S.V.: Detection of unhealthy region of plant leaves and classification of plant  
328 leaf diseases using texture features. Agricultural Engineering International: CIGR Journal **15**(1), 211–217 (2013)
- 329 5. Oppenheim, D., Shani, G.: Potato disease classification using convolution neural networks. Advances in Animal Biosciences **8**(2),  
330 244–249 (2017)
- 331 6. Simonyan, K., Zisserman, A.: Very Deep Convolutional Networks for Large-scale Image Recognition. (2014)
- 332 7. Kawasaki, Y., Uga, H., Kagiwada, S., Iyatomi, H.: Basic study of automated diagnosis of viral plant diseases using convolutional neural  
333 networks. In: International Symposium on Visual Computing, pp. 638–645 (2015). Springer
- 334 8. Fuentes, A., Yoon, S., Kim, S.C., Park, D.S.: A robust deep-learning-based detector for real-time tomato plant diseases and pests  
335 recognition. Sensors **17**(9), 2022 (2017)
- 336 9. He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition. In: Proceedings of the IEEE Conference on  
337 Computer Vision and Pattern Recognition, pp. 770–778 (2016)
- 338 10. Ferentinos, K.P.: Deep learning models for plant disease detection and diagnosis. Computers and Electronics in Agriculture **145**,  
339 311–318 (2018)
- 340 11. Ramcharan, A., Baranowski, K., McCloskey, P., Ahmed, B., Legg, J., Hughes, D.P.: Deep learning for image-based cassava disease  
341 detection. Frontiers in plant science **8**, 1852 (2017)
- 342 12. Torrey, L., Shavlik, J.: Transfer learning, pp. 242–264. IGI global, ??? (2010)
- 343 13. Day, O., Khoshgoftaar, T.M.: A survey on heterogeneous transfer learning. Journal of Big Data **4**(1), 1–42 (2017)
- 344 14. Wang, C., Zhou, J., Zhao, C., Li, J., Teng, G., Wu, H.: Few-shot vegetable disease recognition model based on image text  
345 collaborative representation learning. Computers and Electronics in Agriculture **184**, 106098 (2021)
- 346 15. Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L.: Imagenet: A Large-scale Hierarchical Image Database. In: 2009 IEEE  
347 Conference on Computer Vision and Pattern Recognition, pp. 248–255 (2009). Ieee
- 348 16. Hughes, D., Salathé, M., et al.: An open access repository of images on plant health to enable the development of mobile disease  
349 diagnostics. arXiv preprint arXiv:1511.08060 (2015)
- 350 17. Singh, D., Jain, N., Jain, P., Kayal, P., Kumawat, S., Batra, N.: PlantDoc: a Dataset for Visual Plant Disease Detection. In:  
351 Proceedings of the 7th ACM IKDD CoDS and 25th COMAD, pp. 249–253 (2020)
- 352 18. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.: Improving language understanding by generative pre-training (2018)
- 353 19. Dong, L., Yang, N., Wang, W., Wei, F., Liu, X., Wang, Y., Gao, J., Zhou, M., Hon, H.-W.: Unified language model pre-training for  
354 natural language understanding and generation. arXiv preprint arXiv:1905.03197 (2019)
- 355 20. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L.: Bart: Denoising  
356 sequence-to-sequence pre-training for natural language generation, translation, and comprehension. arXiv preprint arXiv:1910.13461  
357 (2019)
- 358 21. Zhang, J., Zhao, Y., Saleh, M., Liu, P.: Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In:  
359 International Conference on Machine Learning, pp. 11328–11339 (2020). PMLR
- 360 22. Petroni, F., Rocktäschel, T., Lewis, P., Bakhtin, A., Wu, Y., Miller, A.H., Riedel, S.: Language models as knowledge bases? arXiv  
361 preprint arXiv:1909.01066 (2019)
- 362 23. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.:  
363 Language models are few-shot learners. arXiv preprint arXiv:2005.14165 (2020)
- 364 24. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer  
365 learning with a unified text-to-text transformer. arXiv preprint arXiv:1910.10683 (2019)
- 366 25. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning  
367 Transferable Visual Models from Natural Language Supervision. (2021)
- 368 26. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. arXiv preprint arXiv:2109.01134 (2021)
- 369 27. Wang, M., Xing, J., Liu, Y.: Actionclip: A new paradigm for video action recognition. arXiv preprint arXiv:2109.08472 (2021)
- 370 28. Tang, M., Wang, Z., Liu, Z., Rao, F., Li, D., Li, X.: CLIP4Caption: CLIP for Video Caption. In: Proceedings of the 29th ACM  
371 International Conference on Multimedia, pp. 4858–4862 (2021)
- 372 29. Luo, H., Ji, L., Zhong, M., Chen, Y., Lei, W., Duan, N., Li, T.: Clip4clip: An empirical study of clip for end to end video clip retrieval.  
373 arXiv preprint arXiv:2104.08860 (2021)
- 374 30. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly,  
375 S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)

#### 376 Additional Files

377 Additional file 1 — Sample additional file title

378 Additional file descriptions text (including details of how to view the file, if it is in a non-standard format or the file extension). This might  
379 refer to a multi-page table or a figure.

380 Additional file 2 — Sample additional file title

381 Additional file descriptions text.

# Figures



(a) Healthy Apple

(b) Ring Rot Apple

(c) Anthracnose Apple



(d) Healthy Leaf

(e) Anthracnose Leaf

(f) Rust Leaf

Figure 1

Image of part of the data in the dataset.

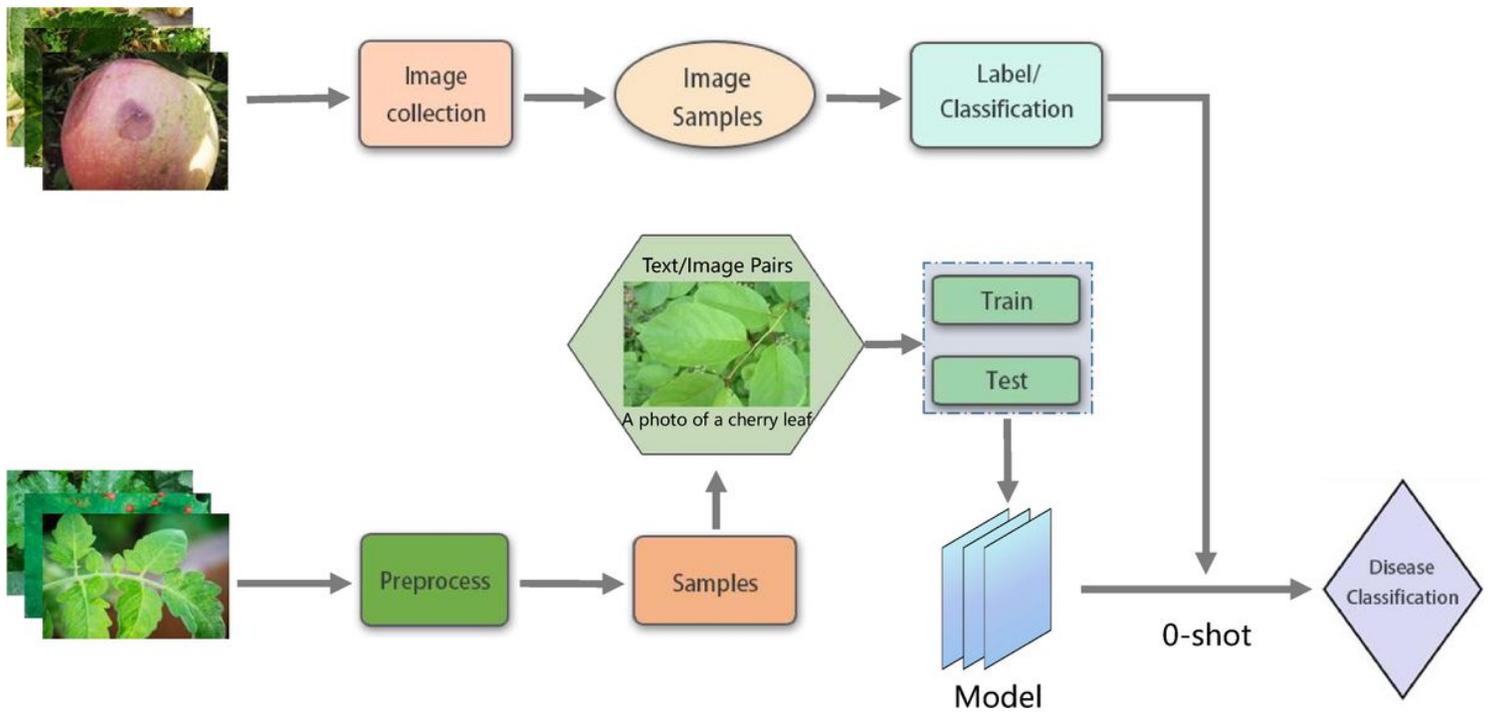
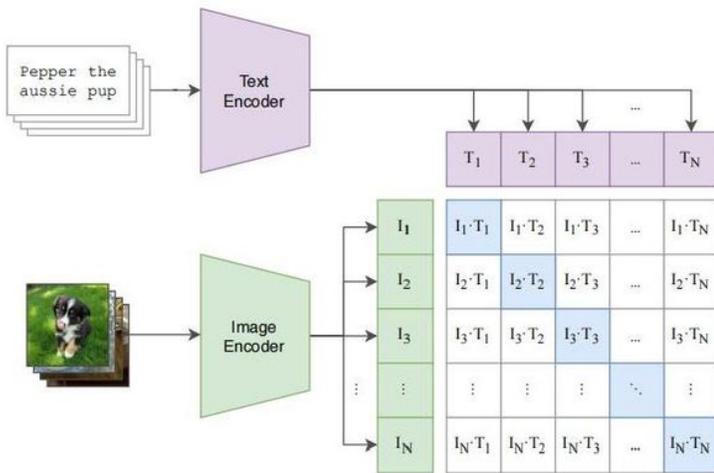


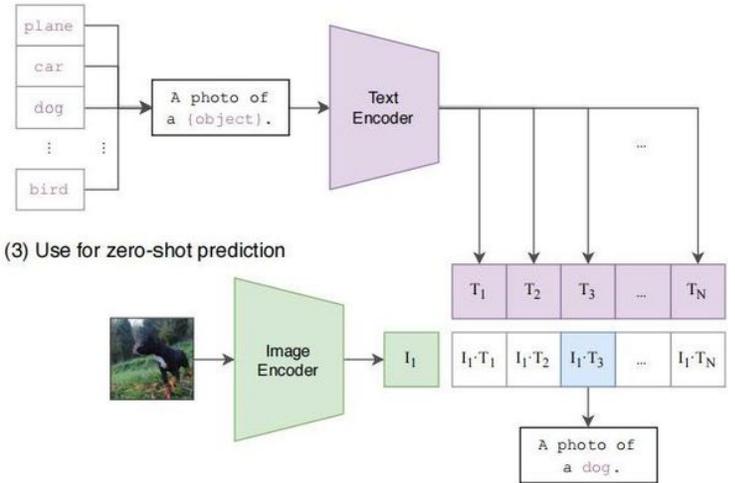
Figure 2

Overall flow chart of the method.

(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

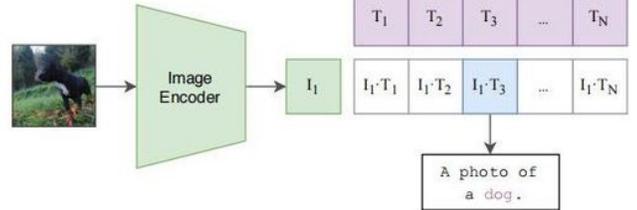


Figure 3

Clip's overall framework diagram.

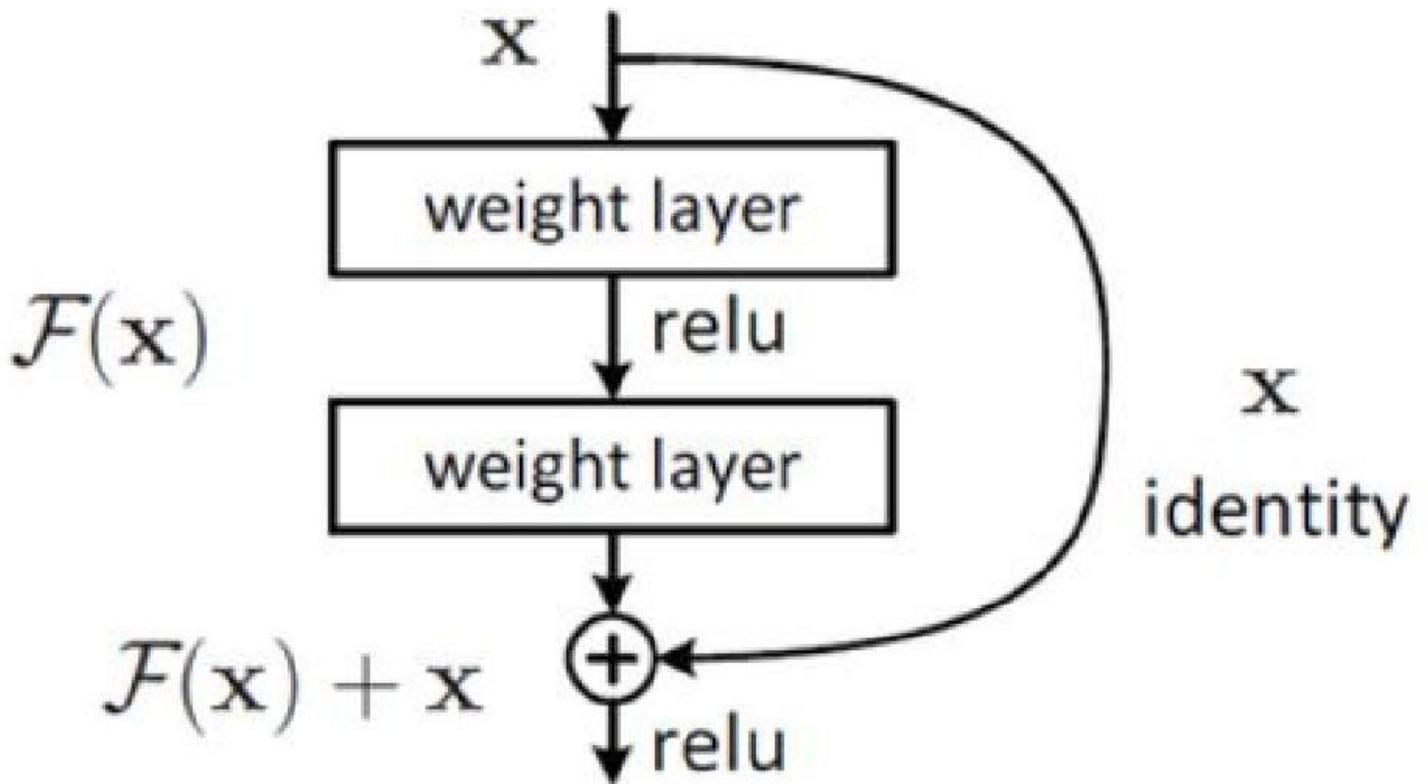


Figure 4

Residual learning units.

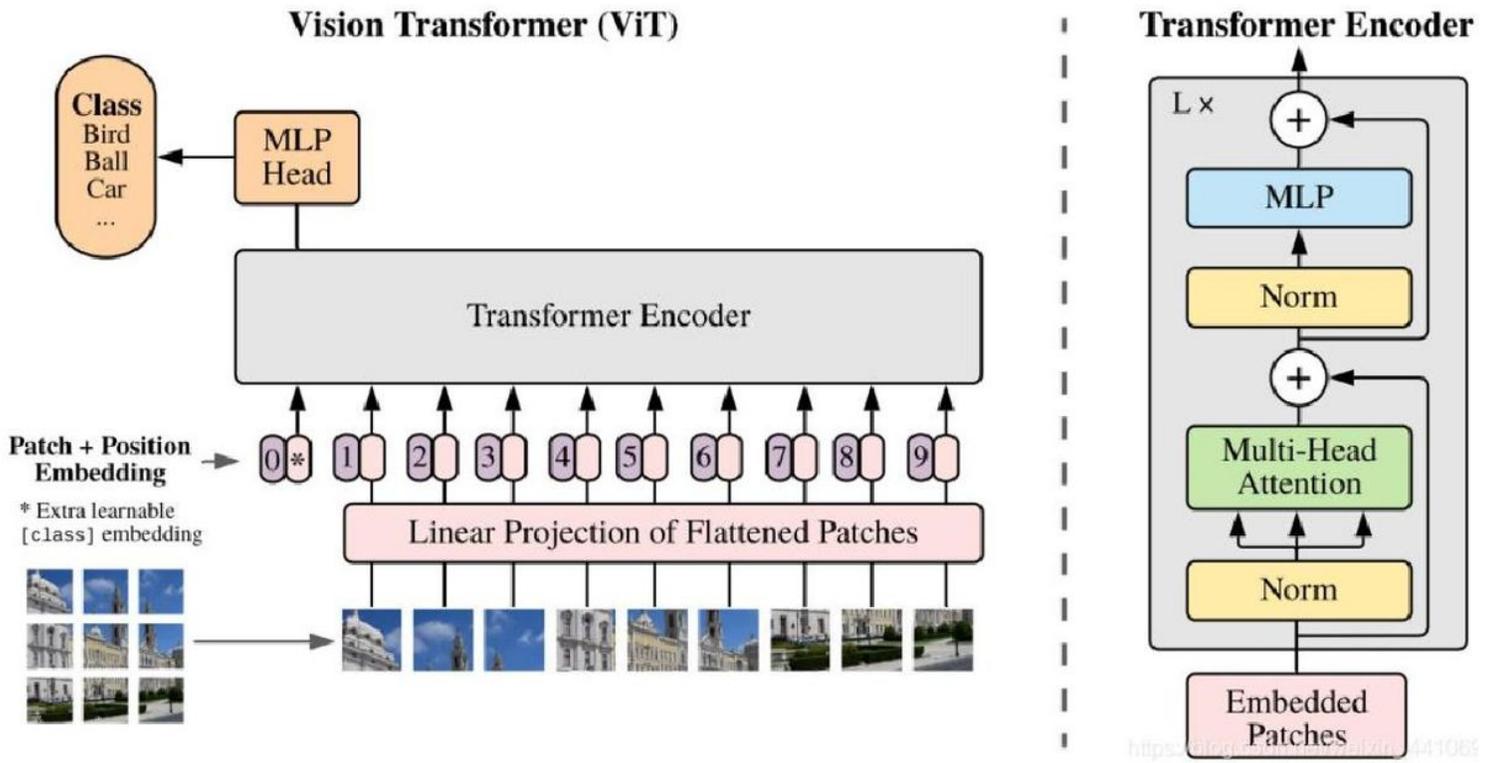


Figure 5

Architecture diagram of Vision Transformer.

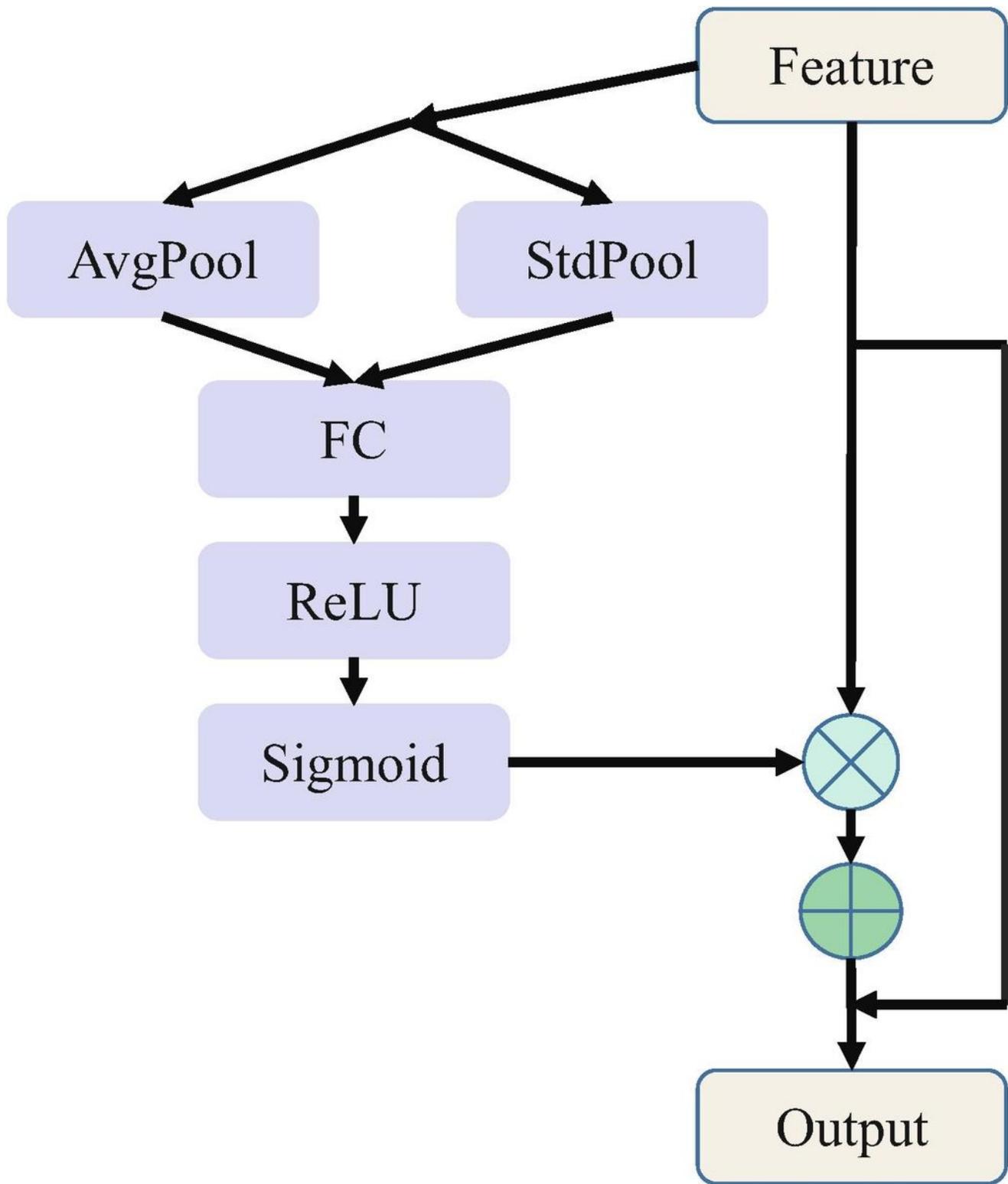
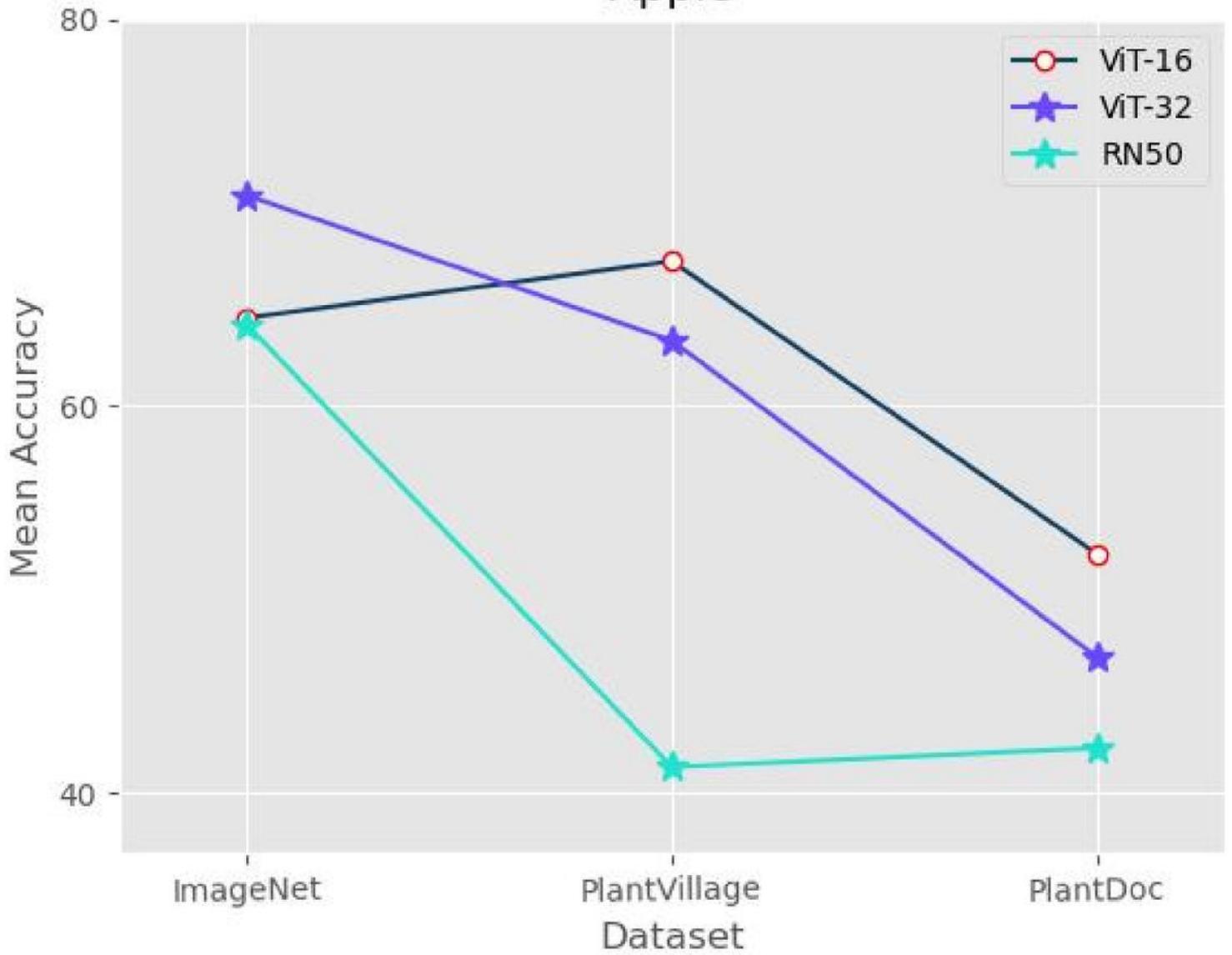


Figure 6

Proposed weighted-pooling Module.

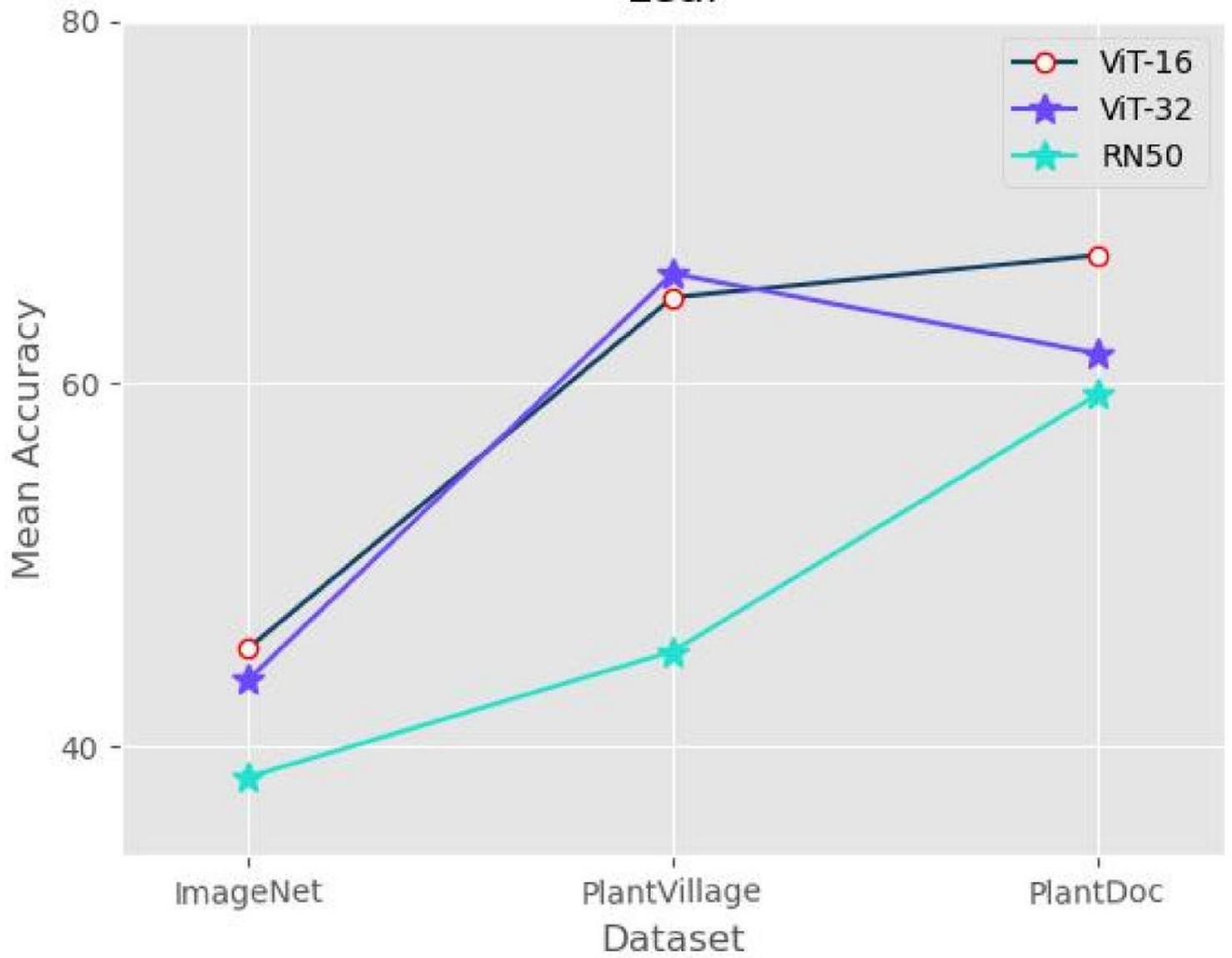
# Apple



**Figure 7**

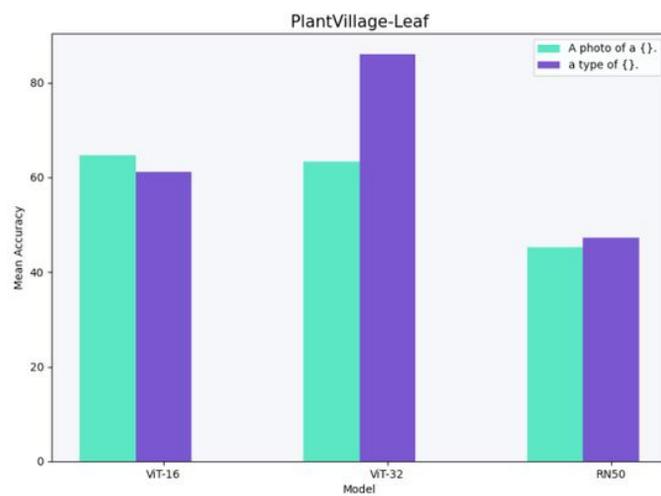
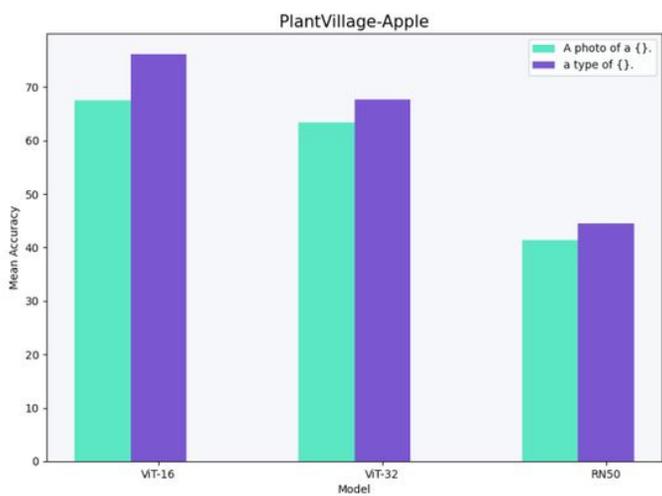
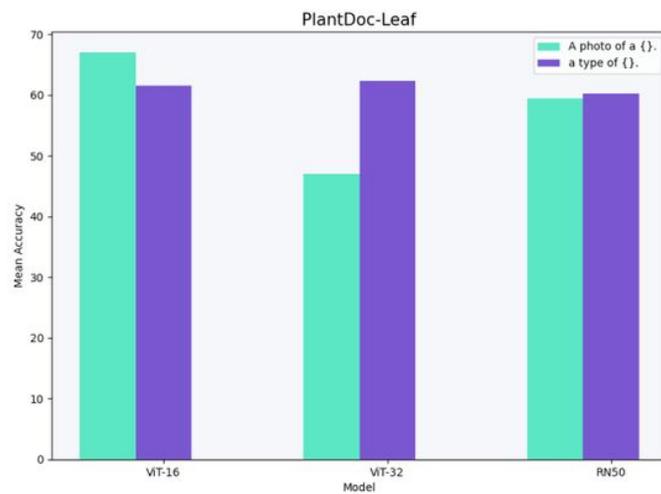
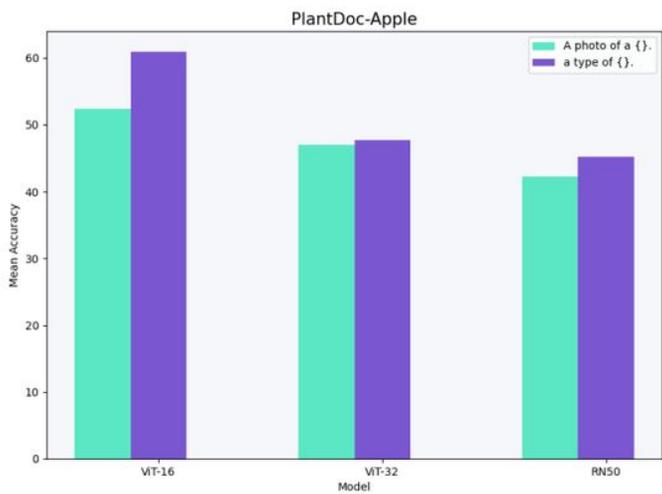
Comparison of average accuracy on pre-trained models on diferent datasets for classi-fication of apple fruits..

# Leaf



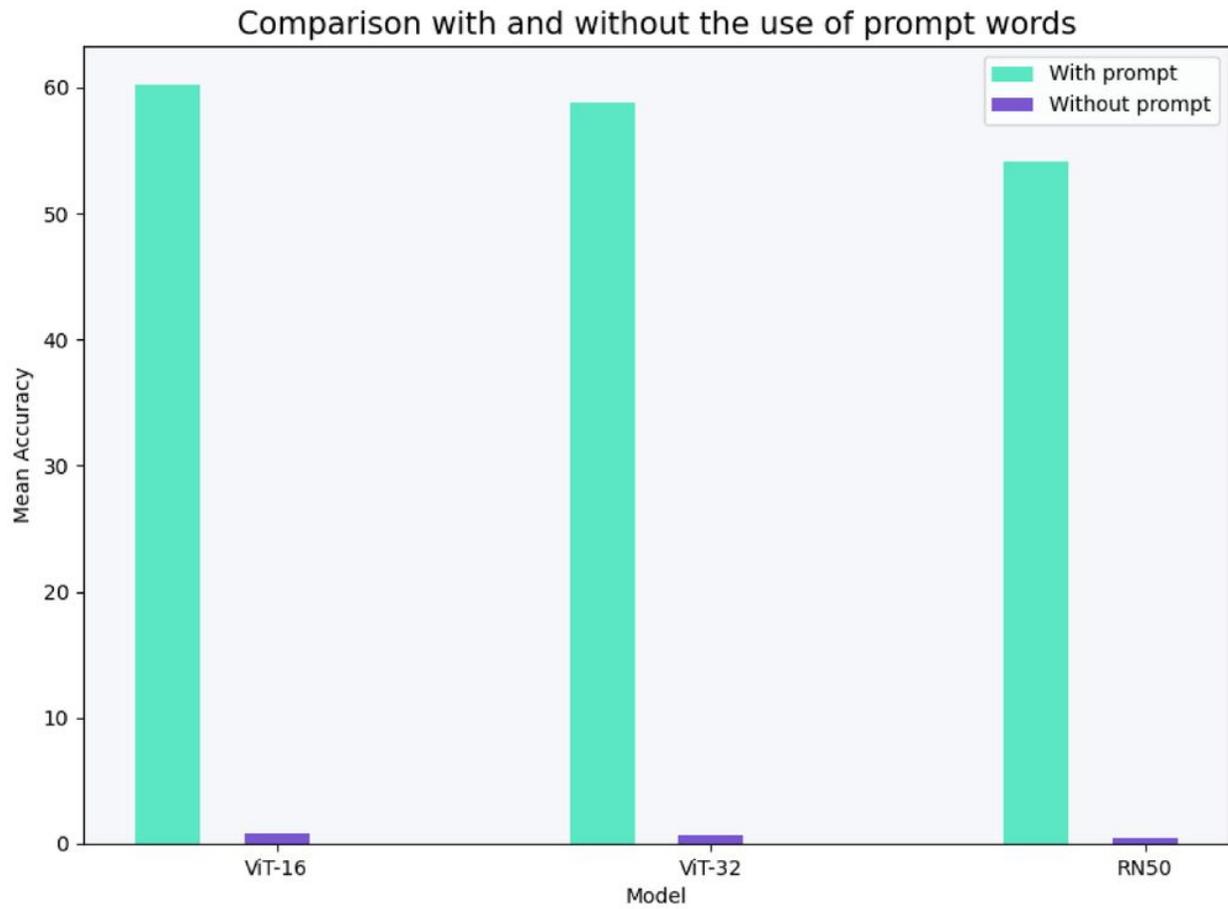
**Figure 8**

Average accuracy comparison of pre-trained models on different datasets for the classification of apple leaves.



**Figure 9**

Plot of the average accuracy of the two prompt words on the two datasets (PlantVillage and PlantDoc) compared.



**Figure 10**

Histogram of average accuracy on the three models with and without the prompt words.

Validity of using the model

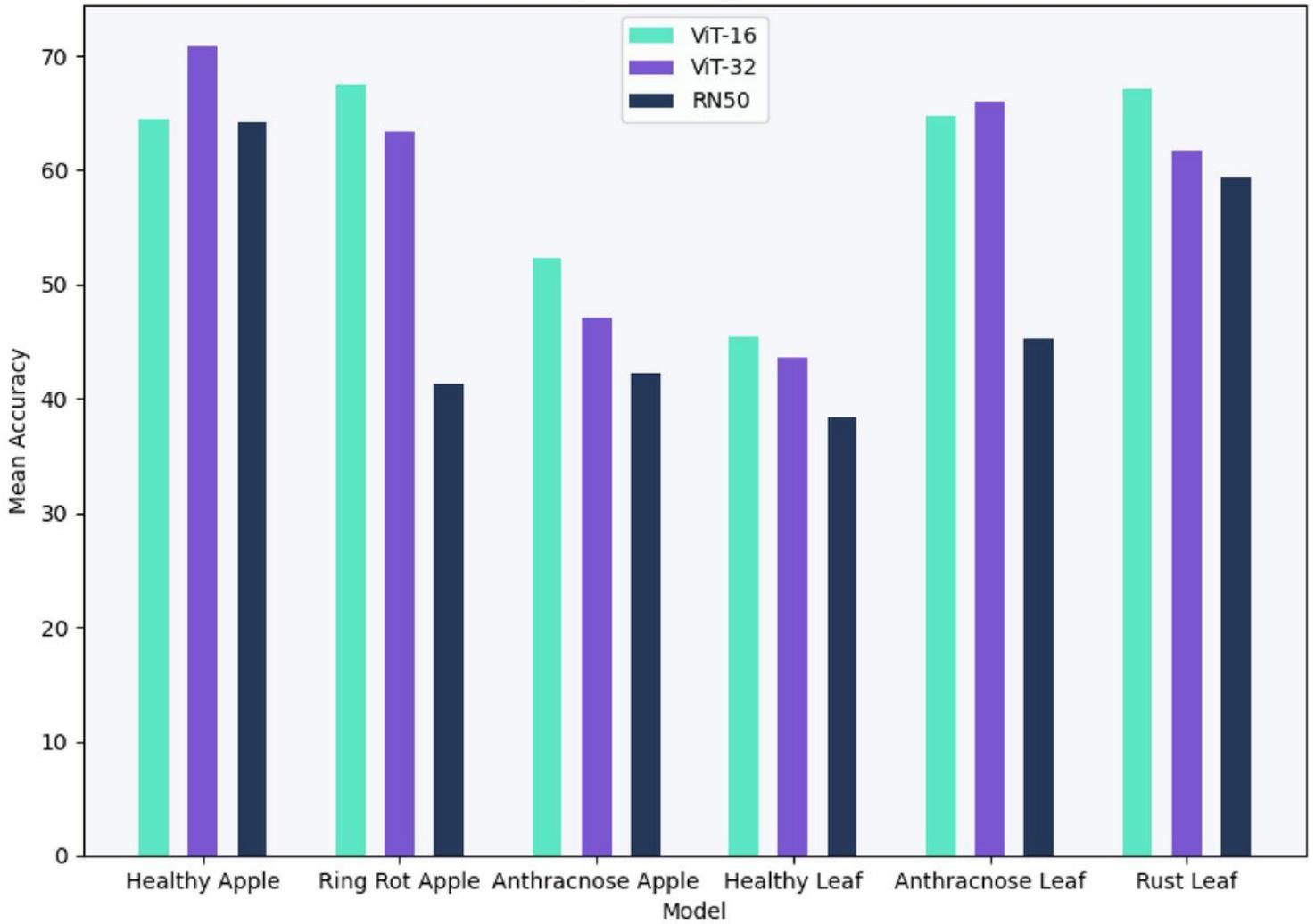


Figure 11

Histogram comparing the average accuracy of each model tested on the Apple dataset.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [TableS1.xlsx](#)
- [TableS2.xlsx](#)
- [TableS3.xlsx](#)
- [TableS4.xlsx](#)