

PlantDet: A benchmark for Plant Detection in the Three-Rivers-Source Region

Huanhuan Li¹[0009-0004-1825-5025], Yu-an Zhang^{1*}[0000-0002-7450-9876],
Xuechao Zou^{1**}, Zhiyong Li¹, Jiangcai Zhaba², Guomei Li², and Lamao
Yongga²

¹ Department of Computer Technology and Applications, Qinghai University, Xining, China <https://www.qhu.edu.cn/>

² Forestry and Grassland Comprehensive Service Center of Yushu Prefecture, Yushu, China

Abstract. The Three-River-Source region is a highly significant natural reserve in China that harbors a plethora of botanical resources. To meet the practical requirements of botanical research and intelligent plant management, we construct a dataset for **Plant** detection in the **Three-River-Source** region (PTRS). It comprises 21 types, 6965 high-resolution images of 2160×3840 pixels, captured by diverse sensors and platforms, and featuring objects of varying shapes and sizes. The PTRS presents us with challenges such as dense occlusion, varying leaf resolutions, and high feature similarity among plants, prompting us to develop a novel object detection network named PlantDet. This network employs a window-based efficient self-attention module (ST block) to generate robust feature representation at multiple scales, improving the detection efficiency for small and densely-occluded objects. Our experimental results validate the efficacy of our proposed plant detection benchmark, with a precision of 88.1%, a mean average precision (mAP) of 77.6%, and a higher recall compared to the baseline. Additionally, our method effectively overcomes the issue of missing small objects.

Keywords: Object Detection · Plant Recognition · Transformer.

1 Introduction

The Three-Rivers-Source region is located in the hinterland of the Qinghai-Tibet Plateau, in the southern part of Qinghai Province. It is the largest nature reserve in China, containing extremely rich wild plant resources. In recent years, the conservation of flora and fauna in the Three-Rivers-Source region has become a focus of attention. However, due to its remote geographical location, underdeveloped information technology, people’s awareness of vegetation protection in the Three-Rivers-Source region is relatively low. Therefore, conducting a survey of

* Corresponding author.

** Huanhuan Li and Xuechao Zou have contributed equally to this work.

plant resources in the Three-Rivers-Source region, especially in plant detection, is of great significance for achieving intelligent plant management and protection.

In recent years, with the rapid development of artificial intelligence and computer vision, many convolutional neural network models based on deep learning, such as AlexNet[1], ResNet[2], and VGGNet[3], have emerged. They have propelled the development of object detection algorithms. The introduction of algorithms such as SSD[4], YOLO series[5,6,7,8,9,10], and the algorithms based on the R-CNN[11], has expanded the promotion and application of object detection in the agricultural field. Numerous experimental results have shown that algorithmic models based on convolutional neural networks perform well in plant recognition research. Therefore, utilizing artificial intelligence and deep learning technology to detect plants in the Three-Rivers-Source region is feasible.

In essence, we have made the following contributions:

- We collected 6965 plant images of 21 categories from the Three-River-Source region, and manually annotated them to establish a large-scale dataset called PTRS for plant detection. This dataset lays the foundation for precise and modern plant detection in the Three-River-Source region.
- We proposed a novel object detection benchmark called PlantDet on PTRS to tackle the challenges of uneven leaf sizes and high feature similarity of diverse plant species. This method consists of three parts: Backbone, Neck, and Head. We introduced an efficient self-attention module based on sliding windows to enhance the feature extraction ability of the backbone and obtained robust feature representation of different scales through efficient feature fusion strategies.
- Experimental results on PTRS demonstrated that our benchmark (PlantDet) surpasses the baseline (YOLOv5), achieves a precision of 88.1% and mAP of 77.6%, and mitigates the problem of missed detection and false positives for small objects.

2 Related work

2.1 Object Detection

Compared to image classification, object detection not only identifies the category of various objects in the image but also determines their location. Object detection can be divided into two types : one is the two-stage algorithm represented by R-CNN[11]. The principle of such methods is to generate candidate boxes, search for prospects, and adjust bounding boxes through specialized modules. Although this candidate region-based detection method has relatively high accuracy, it runs slowly and does not meet the demand for real-time detection.

To tackle the crucial issue of slow detection speed, one-stage object detection algorithms such as SSD and YOLO series algorithms have emerged. They consider the detection task as a regression problem and directly classify and locate objects in the image through a single neural network. Due to the usage of a single

network, they are relatively faster and can meet the real-time detection requirements in the industry. The YOLO series of algorithms have been widely applied in the agricultural field such as detecting diseases and pests[12], maturity[13], and growth stages[14], among others.

2.2 Visual Transformer

In 2017, the Google research team proposed the transformer architecture based on the self-attention mechanism, which achieved tremendous success in the field of natural language processing. The rapid development of the transformer in natural language processing has attracted widespread attention in the field of computer vision. The advantage of a transformer lies in its explicit modeling of long-range dependencies between contextual information, so many researchers have attempted to apply the transformer to computer vision in order to enhance the overall perceptual ability of images. In 2020 Carion et al.[15] proposed the first end-to-end transformer-based object detection model. That same year, the proposal of the image classification model ViT[16] led to the rapid development of visual transformers.

Today, the visual transformer is widely used in various computer vision fields, such as image classification, object detection, image segmentation, and object tracking. So far, many algorithm models based on the visual transformer have emerged: 1) Transformer-based object detection and segmentation models, such as Swin Transformer[17] and Focal Transformer[18] which replace CNN-based backbone networks for feature extraction and combine classic object detection and segmentation networks to complete detection and segmentation tasks; 2) Transformer-based object tracking tasks, such as TrSiam[19] for single-object tracking tasks, TransTrack[20] for multi-object tracking. The rapid development of the transformer in computer vision is mainly due to its ability to extract the relevance of contextual information to obtain global receptive fields, which improves the performance of the model compared to CNN-based models.

3 Method

3.1 Overall Pipeline

Plant detection is an application of object detection technology in botany. A deep learning-based task takes an image with plants as input and outputs the plant's category and bounding box location of its leaves. The Three-River-Source region has diverse flora, to achieve real-time detection, we use YOLOv5 as the baseline for the plant detection pipeline.

Having an efficient model structure is one of the most critical issues in designing a real-time object detector. Our proposed method, PlantDet, uses CSP-DarkNet and CSPPAFPN composed of the same building units for multi-scale feature fusion, and finally inputs the features into different detection heads. The overall model structure of PlantDet is shown in Fig. 1. PlantDet consists of

three parts: 1) Backbone: it mainly performs feature extraction in the main part and effectively extracts crucial feature information of the feature map through downsampling; 2) Neck: this part consists of FPN[21] and PAN[22], respectively performing upsampling and downsampling to achieve the transmission of object feature vectors of different scales and fusion of multiple feature layers; 3) Head: which is made up of three multi-scale detectors and performs object detection on feature maps of different scales using grid-based anchors.

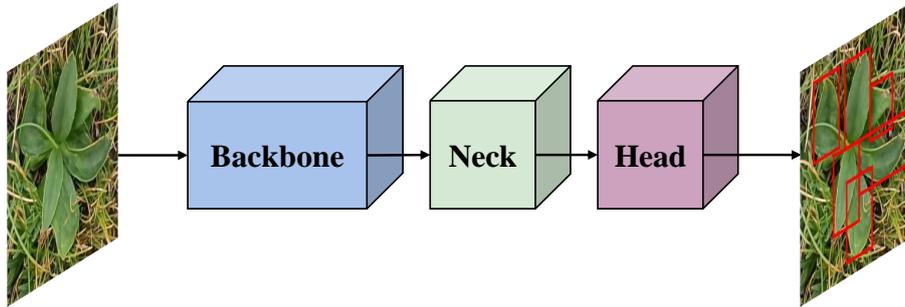


Fig. 1. A pipeline of our one-stage plant detection methods (like the YOLO series).

3.2 Detection Backbone

We use YOLOv5 as the baseline for the plant detection pipeline. YOLOv5 consists of Input, Backbone, Neck and Head. The backbone (refer to Fig. 2(a)) mainly includes C3 and SPPF, where C3 consists of a CBS layer with x residual connections for Concat operation, which improves the feature extraction ability and retains more feature information. SPPF first performs the extracted feature map for multiple maximum pooling operations, and then the results after each maximum pooling are summed for Concat operation, i.e., feature fusion.

In response to the issue of difficult detection caused by varying distributions of leaf sizes in different plants, severe occlusion, and high feature overlap, we have introduced a sliding window module based on self-attention and embedded it into the backbone (refer to Fig. 2) to obtain a robust feature representation with multi-scale resolution.

Specifically, we have re-designed the C3 module in the YOLOv5 backbone, which has the most significant impact on feature extraction. The C3 module primarily acquires feature representation through two parallel convolution branches and introduces residual connections, but does not consider modeling global contextual information. Therefore, we have introduced a self-attention module named "ST block" (refer to Fig. 2(c)) to obtain a global receptive field and more robust representation.

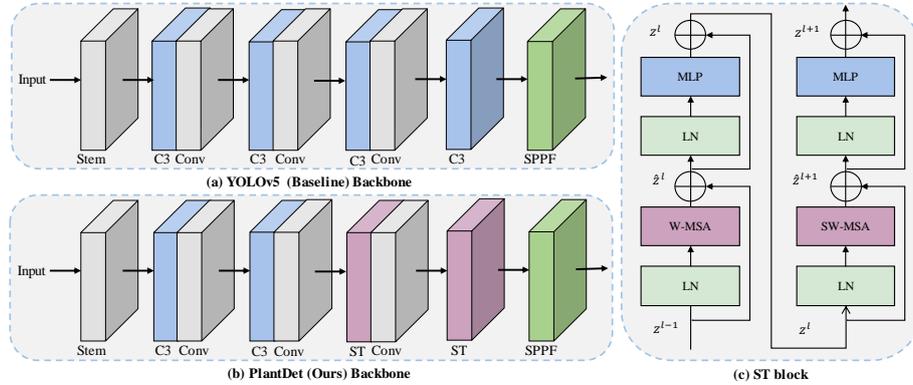


Fig. 2. Backbone of our proposed PlantDet and details of the ST block. (a) Structure of the original YOLOv5’s backbone. (b) Structure of our proposed PlantDet’s backbone. (c) Details of the ST block in PlantDet. The specially designed ST block is used for extracting global contextual information, mainly composed of W-MSA and SW-MSA, for information exchange within and between windows, respectively.

The ST block includes sliding window operation with hierarchical design. It consists of LayerNorm and a shifted window-based MSA with two layers of MLP. Firstly, input features are normalized using Layer Normalization (LN) to expedite model convergence. Subsequently, global feature representation is obtained through the multi-head self-attention mechanism. Furthermore, the features are further enhanced and their expression ability is strengthened through the use of MLP. Finally, residual connection is employed for feature reuse. In addition, a window mechanism is utilized to reduce the additional overhead resulting from the calculation of self-attention matrices.

As is well known, Convolutional Neural Networks (CNNs) perform exceptionally well in local feature extraction due to their inductive bias, while transformer networks based on self-attention mechanisms are effective in modeling long-range global contextual information. Taking into account the superiority of both convolutional and self-attention mechanisms, we have designed a robust backbone feature extractor for plant detection, as shown in Fig. 2(b). It consists of two C3 modules for local feature extraction and two ST blocks for global feature extraction. Finally, the SPPF module fuses the features extracted by both modules to obtain a robust feature representation.

3.3 Loss Function

The task of object detection involves the regression of bounding boxes in addition to classification. Consequently, the training loss function comprises three parts: 1) bounding box regression loss; 2) confidence prediction loss; 3) category prediction loss. These three loss functions are jointly optimized to achieve the

goal of object detection

$$\mathcal{L} = \lambda_1 \mathcal{L}_{reg} + \lambda_2 \mathcal{L}_{obj} + \lambda_3 \mathcal{L}_{cls}, \quad (1)$$

where $\lambda_1, \lambda_2, \lambda_3$ represent the weights of the three loss functions, respectively.

Bounding Box Regression Loss. To account for the large variation in scale among different plant leaves and to balance the impact of objects of different sizes on detection performance, we use the Complete-IoU(CIoU)[23] to calculate the bounding box regression loss

$$\mathcal{L}_{reg} = CIoU = IoU - \frac{\rho^2}{c^2} - \alpha v, \quad (2)$$

Where ρ , c , and v respectively represent the distance, the diagonal length and the similarity of aspect ratio between the centers of the predicted and ground-truth bounding boxes, and α represents the impact factor of v .

Confidence Loss. The loss function for confidence prediction is computed by matching positive and negative samples. Firstly, it involves the predicted confidence within the bounding box. Secondly, it uses the Intersection over Union (IoU) value between the predicted bounding box and its corresponding ground-truth bounding box as the ground-truth value. These two values are then used to calculate the final loss for the confidence prediction, which is obtained through binary cross-entropy

$$\mathcal{L}_{obj}(p_o, p_{iou}) = BCE_{obj}^{sig}(p_o, p_{iou}; w_{obj}), \quad (3)$$

where p_o and p_{iou} represent the predicted confidence and ground truth confidence, respectively, w_{obj} demonstrates the weight of positive samples.

Classification Loss. The category prediction loss is similar to the confidence loss. It involves predicting the category score within the bounding box and using the ground-truth one-hot encoding of the category for the corresponding ground-truth bounding box. The category prediction loss is computed using the following formula

$$\mathcal{L}_{cls}(c_p, c_{gt}) = BCE_{cls}^{sig}(c_p, c_{gt}; w_{cls}), \quad (4)$$

where c_p and c_{gt} represent the predicted values for the corresponding categories.

4 Experiments

4.1 Dataset

Data Collection. The research object of this experiment is the vegetation in the grassland plots distributed in the Three-River-Source region. The plant species image data were taken between July and August 2022 using a handheld camera, approximately 20 centimeters away from the plot, and recorded at a certain speed. After processing, 6965 grassland images were obtained, involving 21 plant species. The plant images involved in the experiment and their corresponding Latin names are shown in Fig. 3. These plants were all identified by experienced experts in the field.



Fig. 3. Samples and corresponding Latin names of our dataset for **Plant** detection in the **Three-River-Source** region (PTRS).

Data Annotations. 6965 images of 21 plant species involved in this experiment were annotated by experienced experts in the field. Initially, the Make Sense (<https://www.makesense.ai/>) labeling tool provided by YOLO was utilized to generate label files containing information about plant categories and target plant coordinates. Subsequently, the above data was organized into VOC format datasets for **Plant** detection in the **Three-River-Source** region (PTRS). Finally, PTRS was divided into training, validation, and testing sets in an 8:1:1 ratio. In addition, comparing our dataset PTRS with other plant detection datasets, the detailed comparison results of these existing datasets are shown in Table 1.

Table 1. Comparison among PTRS (Ours) and other plant detection datasets in agriculture. “-” indicates that this metric is not revealed in the original paper.

Dataset	Annotation way	Classes	Instances	Images	Image Size
AT[24]	Oriented Bounding Box	1	1000	1000	410*410
GHL D[25]	Horizontal Bounding Box	1	-	300	416*416
TDAP[26]	Horizontal Bounding Box	1	-	5000	-
TFP[27]	Oriented Bounding Box	1	-	814	-
GPSD[28]	Horizontal Bounding Box	4	-	1200	-
PTRS (Ours)	Horizontal Bounding Box	21	122300	6965	2160*3840

4.2 Implementation Details

Training Settings. The important training parameters for the model in this experiment were set as follows: training epoch of 300, uniform resizing of input

images to 640×640 resolution, training batch-size of 32, an initial learning rate of 0.01 with Stochastic Gradient Descent (SGD) optimizer. The model was trained on a device with a GPU of 1xNVIDIA A100 and 80GB memory, and the deep learning framework PyTorch was used for implementation.

Evaluation Metrics. In these experiments, Precision (P), Recall (R), and mean of Average Precision (mAP@0.5) are used as evaluation metrics.

4.3 Ablation Studies

Transformer Backbone. To investigate the efficacy of self-attention mechanisms and determine the optimal mechanism applicable to plant detection, we conducted experiments using YOLOv5 as the baseline as shown in Table 2.

Table 2. Ablation experiments of self-attention mechanisms.

Self-Attention	Precision	Recall	mAP@0.5
Baseline	88.0	71.9	76.6
Baseline+MSA	86.1	66.2	72.1
Baseline+W-MSA	88.1	72.9	77.6

The MSA represents the original implementation of self-attention, whereas the W-MSA is a window-based self-attention mechanism. The experimental results demonstrate that compared to the model combined with MSA, the combination of W-MSA module yields better results on the PTRS dataset. Specifically, the Precision, Recall, and mAP were improved by 2.0%, 6.7%, and 5.5%, respectively. This improvement is attributed to the fact that the W-MSA is constructed based on the image resolution hierarchy, which not only achieves feature connections across different windows but also enhances information exchange among different windows, allowing for the extraction of more effective multi-scale feature information to exhibit superior detection performance.

Strategy for Combining Global and Local Modules. In the original YOLOv5, the feature extraction network of the backbone consists of four C3 modules. We conducted ablation experiments to explore the impact of different module combination strategies (C3 and ST block) on the detection results, and the results are shown in Table 3.

The results indicate that the best performance in feature extraction is achieved by using two C3 modules and two ST blocks in the backbone. This is because the C3 module based on the convolutional network can extract local features, while the ST block based on self-attention can extract global features, and the fusion of the two types of features can obtain a more robust feature representation. Therefore, we use two C3 modules and two ST blocks for feature extraction, aiming to improve model performance while minimizing model parameters and computation time complexity.

Table 3. Ablation experiments of module combination strategies.

Number of Module		Precision	Recall	mAP@0.5
C3	ST block			
0	4	87.3	70.8	75.9
1	3	84.3	72.5	75.8
2	2	88.1	72.9	77.6
3	1	85.7	72.3	76.0
4	0	88.0	71.9	76.6

5 Comparison with the State-of-the-Arts

Quantitative Comparison. We conducted experiments to quantitatively compare PlantDet with currently popular object detection algorithms on our self-made PTRS dataset. The results are shown in Table 4. The results indicate that comparing to the baseline YOLOv5, the Recall and mAP of PlantDet increased by 1%, and achieves SOTA results. The outstanding performance of PlantDet is due to the robust detection backbone we have proposed, which integrates global and local information to obtain a more robust multi-scale feature representation. In addition, the numerical evaluation results of Precision, Recall and mAP of baseline (YOLOv5) and PlantDet on the PTRS dataset are shown in Table 5.

Table 4. Quantitative comparison between our and existing models on the dat-aset.

Methods	Precision	Recall	mAP@0.5
SSD [4]	46.6	18.6	48.9
FCOS [29]	-	71.8	57.4
CornerNet[30]	11.0	51.9	38.1
Fast R-CNN [31]	-	40.0	56.3
YOLOF [32]	-	69.7	54.6
YOLOv7 [10]	84.9	72.7	76.0
YOLOv5	88.0	71.9	76.6
PlantDet (Ours)	88.1	72.9	77.6

Qualitative Comparison. In order to further verify the superiority of our proposed PlantDet for plant detection, we conducted qualitative experiments to compare the detection performance of PlantDet and other models (FCOS, YOLOv5 and YOLOv7). The specific visualization results are shown in Fig. 4.

The visualization results show that comparing with other models, our PlantDet can effectively prevent the occurrence of missed inspections and reduce the false detection rate while ensuring detection accuracy. In summary, our PlantDet has better performance for plant detection in the Three-River-Source region, and can meet the needs of botanical Studies and intelligent plant management.

Table 5. Numerical results of YOLOv5 and our PlantDet in 21 categories of our PTRS dataset. Size represents the size of plant leaves, obtained by calculating the bounding box size of all class instances. It can be easily observed that PlantDet enhances the detection performance of small and medium-sized targets, and has a superior effect.

Plant	Size	YOLOv5 (Baseline)			PlantDet (Ours)		
		Precision	Recall	mAP@0.5	Precision	Recall	mAP@0.5
01	Small	86.9	66.6	73.6	85.4	65.4	72.8
02	Medium	88.9	77.1	82.3	88.9	76.7	82.1
03	Medium	89.2	74.9	80.8	89.5	74.9	80.8
04	Large	90.7	75.7	81.5	88.6	76.3	81.8
05	Medium	85.5	66.4	70.3	82.4	63.8	70.2
06	Medium	87.7	73.6	78.0	82.5	74.3	77.2
07	Small	90.0	73.3	77.8	92.9	68.0	75.0
08	Large	95.5	77.9	81.1	96.2	77.0	79.1
09	Small	82.8	67.2	69.8	84.2	67.5	71.5
10	Small	87.7	72.7	81.3	98.8	90.9	90.6
11	Medium	93.5	76.3	81.1	92.3	75.2	81.9
12	Small	89.7	76.0	84.5	89.5	84.0	88.2
13	Medium	86.8	69.7	74.6	85.0	69.7	73.8
14	Medium	91.1	77.1	81.0	90.0	77.9	81.8
15	Small	77.2	53.4	57.6	73.2	55.1	59.6
16	Small	93.7	77.4	81.4	93.6	80.2	85.4
17	Small	82.6	68.2	69.7	81.7	68.1	70.7
18	Small	85.7	69.7	73.7	86.1	69.7	75.0
19	Large	79.8	66.7	66.0	80.4	61.9	63.7
20	Medium	93.3	86.0	92.7	97.3	86.0	89.7
21	Large	88.9	63.0	70.3	91.8	68.5	79.6
Avg.	-	88.0	71.9	76.6	88.1	72.9	77.6

6 Conclusion

To address the problem of varying leaf resolution, severe occlusion, and high feature similarity in plant species, we proposed a novel object detection benchmark called PlantDet. Experimental results show that our PlantDet achieves SOTA detection performance and effectively prevents false detection and missed detection. In addition, we construct a large-scale dataset for plant detection in the Three-River-Source region, which provides data foundation and technical support for the informatization of grassland resources and the construction of a smart ecological animal husbandry new model of "reducing pressure and increasing efficiency" for ecological protection of the Three-Rivers-Source region.

7 Acknowledgments

This study is supported by the Science and Technology Plan of Qinghai Province (2020-QY-218), and China Agriculture Research System of MOF and MARA (CARS-37).



Fig. 4. Visualization results on our PTRS dataset.

References

1. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. *Communications of the ACM* **60**(6), 84–90 (2017)
2. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
3. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
4. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: Ssd: Single shot multibox detector. In: *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*. pp. 21–37. Springer (2016)
5. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 779–788 (2016)
6. Redmon, J., Farhadi, A.: Yolo9000: better, faster, stronger. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 7263–7271 (2017)
7. Redmon, J., Farhadi, A.: Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767* (2018)
8. Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M.: Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934* (2020)
9. Li, C., Li, L., Jiang, H., Weng, K., Geng, Y., Li, L., Ke, Z., Li, Q., Cheng, M., Nie, W., et al.: Yolov6: A single-stage object detection framework for industrial applications. *arXiv preprint arXiv:2209.02976* (2022)
10. Wang, C.Y., Bochkovskiy, A., Liao, H.Y.M.: Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. *arXiv preprint arXiv:2207.02696* (2022)
11. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 580–587 (2014)
12. Liu, J., Wang, X.: Plant diseases and pests detection based on deep learning: a review. *Plant Methods* **17**, 1–18 (2021)

13. Mohammadi, V., Kheiralipour, K., Ghasemi-Varnamkhasti, M.: Detecting maturity of persimmon fruit based on image processing technique. *Scientia Horticulturae* **184**, 123–128 (2015)
14. Tian, Y., Yang, G., Wang, Z., Wang, H., Li, E., Liang, Z.: Apple detection during different growth stages in orchards using the improved yolo-v3 model. *Computers and electronics in agriculture* **157**, 417–426 (2019)
15. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I* 16. pp. 213–229. Springer (2020)
16. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Deghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
17. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 10012–10022 (2021)
18. Yang, J., Li, C., Zhang, P., Dai, X., Xiao, B., Yuan, L., Gao, J.: Focal self-attention for local-global interactions in vision transformers. arXiv preprint arXiv:2107.00641 (2021)
19. Wang, N., Zhou, W., Wang, J., Li, H.: Transformer meets tracker: Exploiting temporal context for robust visual tracking. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 1571–1580 (2021)
20. Sun, P., Cao, J., Jiang, Y., Zhang, R., Xie, E., Yuan, Z., Wang, C., Luo, P.: Transtrack: Multiple object tracking with transformer. arXiv preprint arXiv:2012.15460 (2020)
21. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2117–2125 (2017)
22. Wang, W., Xie, E., Song, X., Zang, Y., Wang, W., Lu, T., Yu, G., Shen, C.: Efficient and accurate arbitrary-shaped text detection with pixel aggregation network. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 8440–8449 (2019)
23. Zheng, Z., Wang, P., Liu, W., Li, J., Ye, R., Ren, D.: Distance-iou loss: Faster and better learning for bounding box regression. In: *Proceedings of the AAAI conference on artificial intelligence*. vol. 34, pp. 12993–13000 (2020)
24. Buzzy, M., Thesma, V., Davoodi, M., Mohammadpour Velni, J.: Real-time plant leaf counting using deep object detection networks. *Sensors* **20**(23), 6896 (2020)
25. Oh, S., Chang, A., Ashapure, A., Jung, J., Dube, N., Maeda, M., Gonzalez, D., Landivar, J.: Plant counting of cotton from uas imagery using deep learning-based object detection framework. *Remote Sensing* **12**(18), 2981 (2020)
26. Fuentes, A., Yoon, S., Kim, S.C., Park, D.S.: A robust deep-learning-based detector for real-time tomato plant diseases and pests recognition. *Sensors* **17**(9), 2022 (2017)
27. Reckling, W., Mitasova, H., Wegmann, K., Kauffman, G., Reid, R.: Efficient drone-based rare plant monitoring using a species distribution model and ai-based object detection. *Drones* **5**(4), 110 (2021)
28. Basavegowda, D.H., Mosebach, P., Schleip, I., Weltzien, C.: Indicator plant species detection in grassland using efficientdet object detector. 42. GIL-Jahrestagung, Künstliche Intelligenz in der Agrar-und Ernährungswirtschaft (2022)

29. Tian, Z., Shen, C., Chen, H., He, T.: Fcos: Fully convolutional one-stage object detection. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 9627–9636 (2019)
30. Law, H., Deng, J.: Cornernet: Detecting objects as paired keypoints. In: Proceedings of the European conference on computer vision (ECCV). pp. 734–750 (2018)
31. Girshick, R.: Fast r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 1440–1448 (2015)
32. Chen, Q., Wang, Y., Yang, T., Zhang, X., Cheng, J., Sun, J.: You only look one-level feature. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 13039–13048 (2021)