# ORCA – Online Research @ Cardiff

# ODL Net: Object Detection and Location Network for Small Pears around the Thinning Period

**Yuqi Lu[1], Shuang Du[1], Ze Ji[2], Xiang Yin[3*], Weikuan Jia[1,4*]**

[1] School of Information Science and Engineering, Shandong Normal University, Jinan 250358, China

[2] School of Engineering, Cardiff University, Cardiff CF24 3AA, UK

[3] School of Agricultural Engineering and Food Science, Shandong University of Technology, Zibo 255000, China

[4] Key Laboratory of Facility Agriculture Measurement and Control Technology and Equipment of Machinery Industry, Zhenjiang 212013, China

**Abstract:** In the process of efficient management of intelligent orchards, due to the short cycle and high intensity of fruit thinning, it is urgent to realize the automatic operation of fruit thinning in orchards. However, affected by the complex orchard environment, the color of fruit and the background are similar, and the more important problem is that the fruit is small-scale. These factors bring great challenges to fruit detection before and after the thinning period. For this reason, a detection algorithm for fruits of small green objects is proposed, namely, ODL Net. By integrating the semantic enhancement module and label assignment Center-Box, the small size problem of the target fruit is alleviated. The feature enhancement module and position enhancement module are constructed to enhance the fusion effect of features and improve the detection accuracy. To better verify the performance of the algorithm, this study takes a pear orchard as an example to produce two datasets before and after pear thinning. The experimental results show that the detection accuracy of ODL Net can reach 56.2% and 65.1% before and after the fruit thinning period, respectively, and the recall rate can reach 61.3% and 70.8%, respectively, which are significantly higher than those of other mainstream algorithms at present. The new algorithm can effectively assist the orchard automatic fruit

23     thinning operation and provide the basis for orchard yield measurement after the fruit thinning period. This study

24     can provide a theoretical basis for the scientific management of intelligent orchards.

26

## 1. Introduction

28     The development of cutting-edge theories and technologies such as artificial intelligence and

29     5G communication provides strong support for efficient agricultural production. Intelligent

30     agriculture (Patrício and Rieder, 2018; Wu and Tsai, 2019) and orchards (Xu et al., 2023; Maheswari,

31     et al., 2021) have gradually entered the public view, and agricultural production efficiency has been

32     greatly improved. In the orchard production process, due to the short operation cycle and high labor

33     intensity of fruit thinning, it is urgent to realize automatic fruit thinning in orchards. However, in

34     the fruit thinning period, orchards present a complex environment, and the fruit color is similar to

35     the background and is still small-scale and easily covered by branches and leaves. These factors

36     bring great challenges to the efficient recognition of fruit during this time. The realization of the

37     efficient detection of small fruits can assist automatic fruit thinning operations in orchards when the

38     fruit is clustered at the early stage of fruit thinning. It can also assist the fruit yield measurement to

39     realize the scientific management of the orchard when the fruit is in a single state in the late stage

40     of fruit thinning. In addition, it also helps fruit farmers recalculate irrigation and fertilizer supply

41     due to the change in fruit quantity after thinning. Taking the golden pear orchard as an example, this

42     study focuses on the detection accuracy of small target fruits before and after pear thinning and

43     constructs a high-precision small fruit detection algorithm.

44     In the orchard environment, object detection has been widely used in orchards (Gongal et al.,

45    2015; Fu et al., 2020; Tang et al., 2023), such as automatic driving (Yang et al., 2021; Tey and

46    Brindal, 2022), pest detection (Ebrahimi et al., 2017; Ngugi et al., 2021), and other operations. Its

47    detection accuracy also restricts the production efficiency of orchards. In complex orchard

48    environments, fruit detection has attracted many scholars' attention and has also achieved gratifying

49    research results. Sa (Sa et al., 2016) proposed a fruit detection algorithm based on Faster RCNN

50    that used images obtained from two modes, color (RGB) and near-infrared (NIR), to compose multi-

51    modal information; in this paper, the algorithm was applied to the detection task of seven kinds of

52    fruits, such as sweet pepper and rock melon. Bargoti (Bargoti, et al., 2017) proposed a tiling method

53    for images containing more than 100 target fruits; combined with image enhancement technology,

54    the F1-score of this new algorithm on apples and mangoes exceeded 0.9. Zhao (Zhao and Yan, 2021)

55    proposed CenterNet for fruit detection, which implemented three backbone networks and finally

56    confirmed CenterNet based on DLA-34 (Yu et al., 2018). In addition, Jia (Jia et al., 2021) presented

57    an algorithm with a transformer structure, which was popular in recent years, to detect green apples

58    in orchards. Hussain (Hussain et al., 2022) proposed a deep learning based framework for automatic

59    detection and recognition of fruits and vegetables in complex scenes. It can help sellers identify

60    vegetables and fruits with high similarity. Although its accuracy is as high as 96%, there is no special

61    design for detecting small-scale fruits. Most of the above algorithms were detection algorithms

62    proposed for specific orchard environments. These algorithms achieved relatively ideal detection

63    results for large-scale fruits, green fruits, etc. However, they ignored the detection effect of small-

64    scale target fruits.

65        The detection effect of small objects is easily affected by the external environment. For

66    example, the proportion of pixels is small, so the features are difficult to effectively represent. The

67   target itself is small-scale and easily occluded by the background, resulting in missing its recognition.

68   In addition, the color of a small target is similar to the background, leading to incorrect identification.

69   Small object detection is so challenging that it has attracted scholars' attention in many fields. Rabbi

70   (Rabbi et al., 2020) used small objects to over-sample images and enhanced each image by copying

71   and pasting small objects many times to achieve small object detection; the detection accuracy of

72   small objects on the MS COCO dataset was increased by 7.1 percentage points. Yang (Yang et al.,

73   2019) presented a novel multi-category rotation detector for small, cluttered and rotated objects,

74   namely, SCRDet, in which a sampling fusion network was devised that fused multi-layer features

75   with effective anchor sampling to improve the sensitivity to small objects. It was shown on the two

76   remote sensing public datasets and the COCO and VOC 2007 dataset. In the area of remote sensing,

77   Zhang (Zhang et al., 2018) proposed a network with deconvolution layers after the last convolution

78   layer of the basic network for small object detection in remote sensing data; in an experiment on a

79   remote sensing image dataset, the Deconv RCNN reached a much higher mean average precision

80   than the Faster RCNN. Inspired by these different fields, research on small fruit detection has also

81   made significant progress. Mai (Mai et al., 201) presented a multi-classifier fusion strategy for small

82   fruits, which used three different feature levels to learn three classifiers for object classification in

83   the proposal localization phase; at the same time, a new classifier correlation loss term was

84   introduced to improve the detection accuracy of small objects. Tu (Tu et al., 2020) proposed an

85   improved method based on multi-scale Faster RCNN, which used color and depth images acquired

86   by an RGB-D camera; it was improved by combining the feature map of the shallow convolution

87   maps from the region of interest (ROI) pool to detect small passions. Sun (Sun et al., 2022) proposed

88   a balanced feature pyramid network (BFP Net) for small apple detection; the network balanced the

89    information mapped to small apples from two perspectives and was verified on three fruit datasets.

90    The above algorithms achieved ideal results in solving small objects or small target fruit detection.

91    However, these small targets were mostly "small" due to the perspective that they were not as small

92    as the fruit in the fruit thinning period.

93    At present, there are relatively few studies on fruit recognition during fruit thinning. At this

94    time, the state of fruit appears as the target color is similar to the background, and the real volume

95    of fruit is relatively small and easily blocked. To solve the above problems, this study presents ODL

96    Net, a detection algorithm for small-scale fruit around the pear thinning period. The semantic

97    enhancement module (SEM) and the label assignment Center-Box in this algorithm can deal with

98    small-scale fruit detection well. In addition, the feature enhancement module (FEM) and the

99    positional enhancement module (PEM) for feature fusion also improve the detection accuracy. The

100    following is an explanation of the innovations in this paper:

101    (1) This study presents ODL Net, a novel detection algorithm for small fruits around the pear

102    thinning period. The detection accuracy in orchards is higher than that of most current detection

103    algorithms.

104    (2) Two pear datasets are prepared in this study, including image data before and after the pear

105    thinning period. In this way, the detection effect of ODL Net around this period in the orchard can

106    be accurately verified.

107    (3) Three modules, SEM, FEM, PEM  are constructed in the feature fusion network. The

108    modules enhance the information from different angles and provide it to the downstream detection

109    task of the ODL Net.

110    (4) ODL Net relies on a special label assignment, Center-Box, to accurately locate small fruits.

111     Center-Box eliminates the influence of object size on positive sample allocation, avoiding ignoring

112     small objects.

113     This study introduces the ODL Net, which aims to achieve precise detection of pear fruits in r

114     eal orchard environments both before and after the thinning period. The pre-

115     thinning detection of pears provides valuable insights to fruit farmers, allowing them to monitor ea

116     rly-

117     stage fruit growth. This not only facilitates the determination of optimal irrigation and fertilizer su

118     pply for orchards but also guides the thinning process. Similarly, post-

119     thinning detection of pear fruits remains crucial, providing ongoing recommendations for irrigatio

120     n and fertilizer supply to fruit farmers and enabling scientifically informed yield predictions for or

121     chards. In summary, the primary objective of this

122     study is to enable comprehensive monitoring of fruit growth stages, encompassing both pre

123     and post-thinning stages, thereby achieving intelligent management of orchards.

124     The organizational structure of this article is as follows: Section 1 describes this research

125     purpose and related work in the current field. The second section is the production process of the

126     two datasets. Section 3 details the composition of the ODL Net, as well as the structure and functions

127     of each component. The experimental details, data and results are shown in Section 4, including

128     contrast and ablation experiments. The summary and expectation of the overall research content is

129     presented in Section 5.
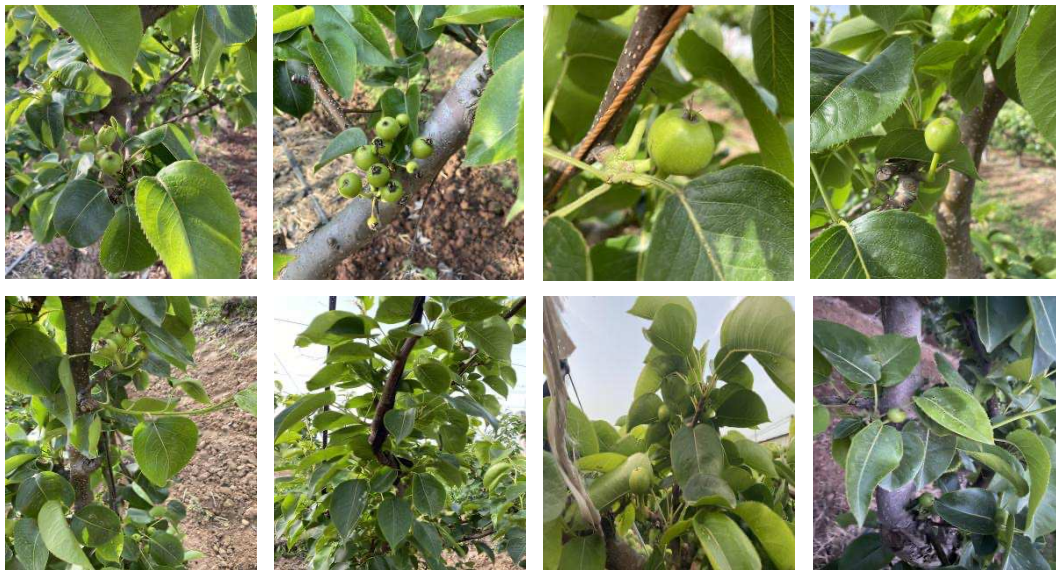
## 130   2. Datasets

131     In this study, two datasets were produced, corresponding to the period before and after pear

132     thinning. The object of the datasets was the golden pear around the fruit thinning period, which was

133    characterized by its small-scale and full green color. The following is an introduction to the data

134    collection and production process.

135    **2.1 Data Acquisition**

136        The objective of this study is to detect the fruit before and after the pear thinning period in

137    orchards to provide thinning guidance for fruit farmers and realize intelligent orchard management.

138    To achieve this goal, two pear datasets were made before and after the pear thinning period to test

139    the feasibility of the ODL Net. The datasets were all taken from the RiSheng Golden Pear

140    Professional Cooperative of Jiaozhou, Qingdao City, Shandong Province. The images taken were

141    saved as.jpg, 24-bit color. As shown in Figure 1, the fruit was characterized by its small-scale and

142    green color around the pear thinning period. It can be seen from the figure that the fruit density

143    before thinning is higher than that after thinning, so detection before thinning is more difficult.



a) images before thinning period                    b) images after thinning period

144                          Fig. 1. Images around the pear thinning period in datasets

145    **2.2 Data Processing**

146        LabelMe was used to process the images taken. It used boxes to mark the target, with the

147 marked closed part as the foreground, labeled "pear", and the remaining part as the background.

148 Annotated images were automatically generated into.json files containing coordinates and label

149 information. The final datasets were divided according to a ratio of 7:3. The dataset before pear

150 thinning included 1549 images, with 1084 images in the training set and 465 images in the test set.

151 The dataset after pear thinning included 891 images, with 623 images in the training set and 268

152 images in the test set. We also calculated statistics on the object scale of the dataset, and the
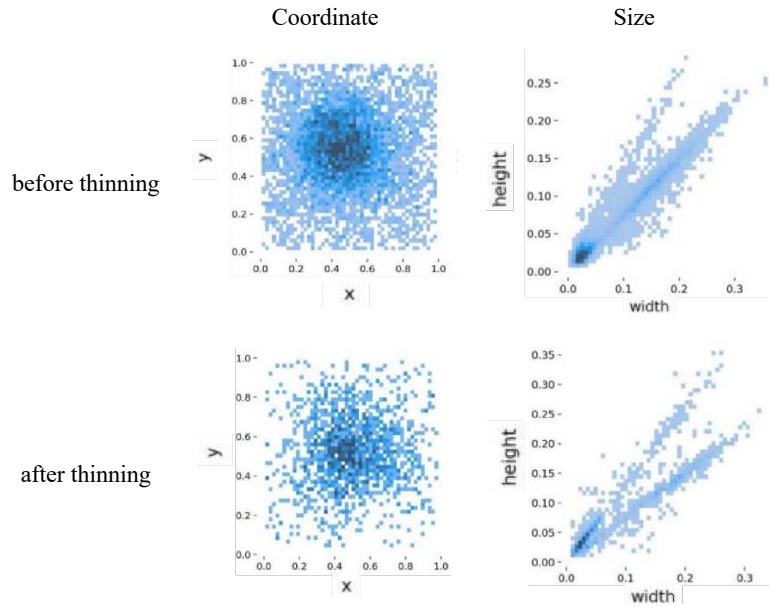
153 information is shown in Table 1.

154

Table 1. Statistics of pear fruit scale

|  | small-scale | middle-scale | large-scale | images |
|---|---|---|---|---|
| dataset before thinning | 4427 (48.41%) | 3251 (35.55%) | 1466 (16.04%) | 1549 |
| dataset after thinning | 972 (45.13%) | 641 (29.76%) | 541 (25.11%) | 891 |

155 It should be noted that COCO format datasets usually define objects with area pixels less than

156 32×32 as small-scale targets, objects with area pixels greater than 96×96 as large-scale targets, and

157 objects between them are defined as medium-scale targets. However, the image size of the pear

158 datasets we shot is 3024×4032 pixels, which is bigger than the image size of 640×640 pixels in the

159 COCO dataset. Therefore, this study takes the pixel area as the standard and redefines the scale

160 range according to the multiple relationships of areas. Objects with pixels less than 174×174 are

161 small-scale targets, objects with pixels greater than 523×523 are large-scale targets, and objects in

162 between are defined as medium-scale targets. As seen from Table 1, large-scale objects have the

163 least amount. Intermediate-scale objects are the most numerous, accounting for more than half of

164 the fruit. The fruit coordinates and sizes in the datasets are visualized as shown in Figure 2. The

165 coordinate diagram shows that the fruit density after thinning is much lower than that before
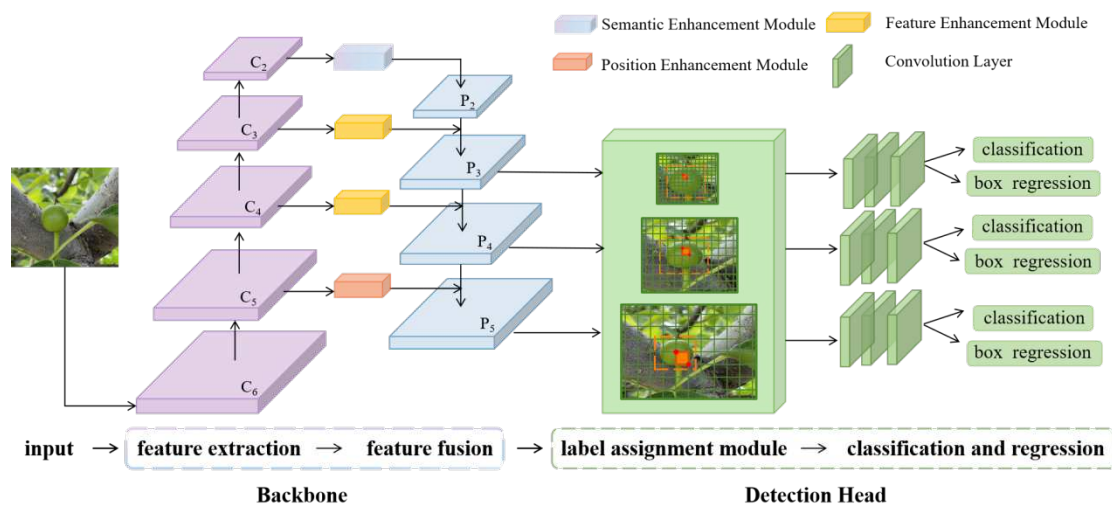
166    thinning. The fruit size before fruit thinning was smaller, concentrated within 0.1.



Fig. 2. Statistics of pears in datasets

167

## 3. ODL Net Detection Model

169    ODL Net includes two parts: the backbone and detection head. The backbone network consists

170    of feature extraction (bottom-up) and feature fusion (top-down). Three enhancement modules are

171    built in the backbone for more effective feature fusion and small fruit feature capture. In the

172    detection head, this study uses a label assignment that ignores fruit size to strengthen the detection

173    of small objects. The overall structure of the algorithm is shown in Figure 3.



174

Note: ODL Net mainly includes two parts: the backbone and the detection head.

## 3.1 Image Enhancement

To enhance the learning ability of the network, the algorithm enhances the image for the input. Before the input is used for training, they are first scaled to 640×640 pixels. The scaled images are enhanced in four ways: random rotation, saturation transformation, random affine transformation and random stitching. An example of data enhancement is shown in Figure 4. The random rotation operation randomly rotates the original image by 90°. The random saturation transform changes the hue and saturation value of the input to simulate different light conditions in the orchard. Random affine transformation includes random t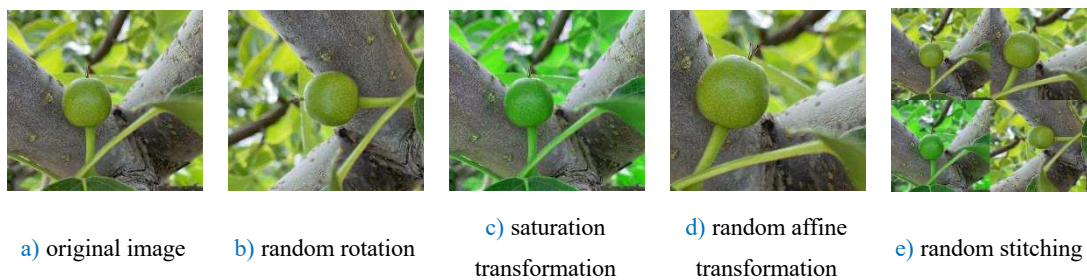ranslation, scaling and rotation. The translation, scaling, and rotation factors are set to 0.0625, 0.5, and 45 degrees, respectively. Finally, 4 images are randomly selected for mosaic processing. The image enhancement operation expands the learning range of the neural network to better learn the fruit feature.



a) original image    b) random rotation    c) saturation transformation    d) random affine transformation    e) random stitching

Fig. 4. Image enhancement display

## 3.2 Feature Extraction

The feature extraction network mainly includes CBS, CSP and SPPF modules, and its architecture is shown in Figure 5. The figure shows the extraction process of feature maps in each layer and the specific structure of each module.
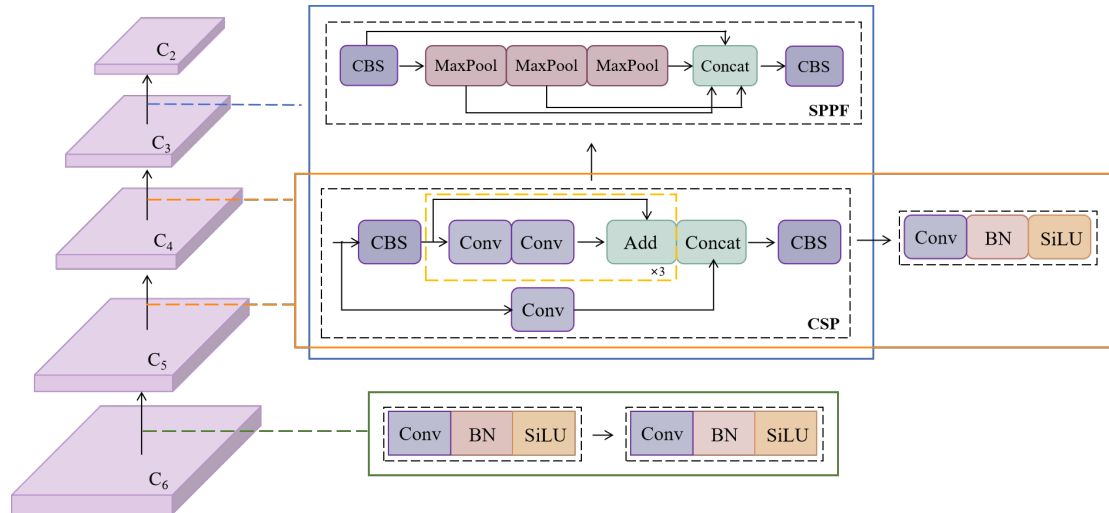
Fig. 5. Diagram of the feature extraction network

The convolution layer, batch normalization layer and activation function leaky ReLU are encapsulated in the CBL of YOLOV5. The ability of the activation function is nonlinear in the neural network, which is replaced by SiLU (Elfwing et al., 2018) in this study, and the module is named CBS. SiLU is defined as the activation of network function approximation in reinforcement learning. It is a weighted linear combination of sigmoid, whose function expression is $SiLU(x) = \frac{x}{1-e^x}$. Leaky ReLU solves the problem of zero ReLU output, but it is still nearly linear. As shown in Figure 6, unlike Leaky ReLU, SiLU is not monotonically increasing but has a minimum value. This makes it self-stable, thus inhibiting the learning of a large number of weights. It can nonlinear neural networks better than Leaky ReLU, thus improving the expression ability of networks to models and solving problems that linear models are not equipped to deal with.
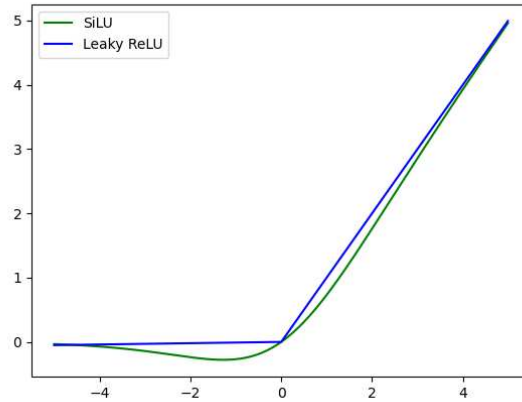
The CSP divides the input into two branches. The number of channels is halved by the convolution operation, and one of the branches is subject to a multi-layer residual operation (i.e., double-layer convolution residual component). Then, two branches are concatenated to make the input and output the same size. Finally, a CBS module is placed to further process feature information, which enables the feature extraction network to learn more fruit features. The SPPF pools the features passing through the CBS three times and concatenates the four groups of features. SPPF specifies one convolution kernel, and the output of each pooling layer is used as the input of the next pooling, which is faster than specifying three. Similar to the CSP, the last step of SPPF is still the CBS module. The SPPF increases the feature representation ability of the feature maps.

As shown in Figure 5, feature map $C_5$ is processed by $C_6$ through two stacked CBS modules, whose main step is convolution operations. Both operations of $C_4$ and $C_3$ are the same, passing through a CSP and a CBS module. The top-level feature map $C_2$ is obtained by $C_3$ through the CSP and SPPF modules in series, which enhances the expression of the algorithm for small objects. The feature extraction network generates six feature layers, denoted from bottom to top by $C_6$-$C_2$. However, only the upper five layers are used for feature fusion. The remaining layer is used to deepen the network and obtain richer feature information.

## 3.3 Feature Fusion

Feature fusion includes horizontal and vertical fusion. Horizontal fusion adds three different enhancement modules, which can also be used in the feature fusion phase of any other algorithm. Top-down fusion combines the CSP and CBS modules in the feature extraction network and uses the sampling and concatenation operations to fuse the adjacent feature maps. The following describes the overall architecture of the three modules and the feature fusion network.
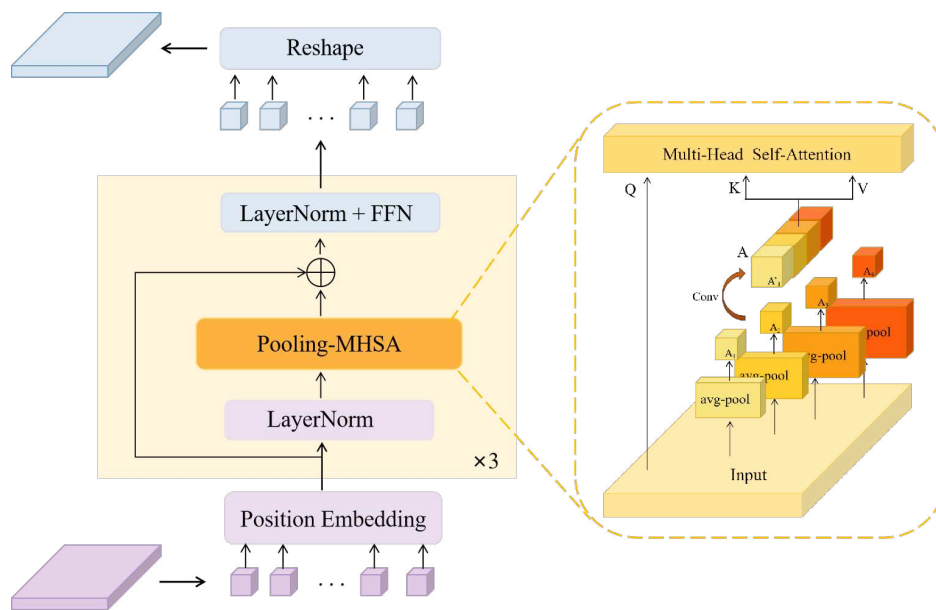
### 3.3.1 Semantic Enhancement Module (SEM)

Recently, there have been two main structures for image feature processing: Convolutional Neural Network (CNN) and Transformer, which have different core concepts. CNN focuses on the correlation between two-dimensional local data. With the deepening of layers, its focus area will be wider. This makes it suitable for image processing, especially layer-by-layer processing of images (Lin et al., 2017; Liu et al., 2018). However, it cannot capture long-distance information and is limited by the receptive field. A common solution to this problem is to increase the depth of the neural network. This approach can indeed obtain more global information, but it will lead to gradient instability, network degradation and other problems.

At this time, transformers are widely used in the field of computer vision by virtue of their excellent spatial modeling ability (Zhu et al., 2020; Liu et al., 2021; Liu et al., 2022). The multi-head attention mechanism in the visual transformer captures richer information and relationships of features. However, the limitation of Transformer is that it cannot take advantage of the prior knowledge of scale, translation invariance and feature locality of the image itself, which makes it necessary to use a large amount of data for training. In addition, the main reason why the transformer structure cannot replace CNN at present is computational efficiency due to its sequential input

245     format. In natural language processing, the sequence length of the WMT 2014 English-German

246     dataset containing 50 million words and 2 million sentences is only 25. The code length is increased

247     to 3136 when the image resolution of the common ImageNet dataset is 224, and the segmented

248     image block size is defined as 4×4.

249        In this study, to detect small-scale fruits, the pixels in the dataset will be higher, and the

250     corresponding coding length will be multiplied, which is difficult for computational memory. In

251     consideration of the above factors, this study only constructs a semantic enhancement module with

252     the help of a transformer structure. It does not involve the hierarchical association of feature maps

253     but is applied to a feature map itself to enhance its semantic information. The specific structure of

254     the SEM is shown in Figure 7.



255

256                     Fig. 7. Diagram of the semantic enhancement module structure

257        The feature map obtained through the feature extraction network is divided into many patches,

258     generating sequences for position embedding. Then, three structural layers consisting of the norm

259     layer, multihead attention and feed-forward network (FFN) are used to enrich the semantic

260     information. Finally, the sequence is restored to a feature map of the same size through the Reshape

261  operation. In the process above, the traditional multihead attention is replaced by pooling-MHSA,

262  whose structure is shown in the right dotted line box in Figure 7.

263      The input sequence X is reshaped into the feature map format when it enters the pooling-

264  MHSA. The reshaped feature map is still represented as X for the convenience of understanding the

265  input. Multiple average pooling layers of different sizes are applied to X to generate a feature map

266  $A_i$ of different contents:

267                                 $$A_i = \text{AvgPool}_j(X) \tag{1}$$

268  Where i=1, 2, 3, 4. J stands for pool size, $j = (\frac{H}{\text{ratios}} \times \frac{W}{\text{ratios}})$. H, W represents the size of the input

269  feature map, ratios=[1, 2, 5, 10]. Next, the feature map is sent to the convolution layer for relative

270  position coding:

271                                 $$A'_i = \text{Conv}(A_i) + A_i \tag{2}$$

272  The encoded feature maps are stacked:

273                          $$A = \text{LayerNorm}(\text{Concat}(A'_1, \ A'_2, A'_3, A'_4)) \tag{3}$$

274  The stacked feature maps A carry more context information in feature map X, which can replace X

275  as the input of the subsequent multihead self-attention. The size of pooled feature maps is smaller,

276  so the generation of K and V matrices is smaller than that of traditional ones, which means that the

277  Pooling-MHSA is more efficient. It can be expressed as Equation 4 and Equation 5:

278                           $$(Q, K, V) = (XW^q, AW^k, AW^v) \tag{4}$$

279                       $$X_{\text{Patt}} = \text{Softmax}(\frac{Q \times K^T}{\sqrt{d\ K}}) \times V \tag{5}$$

280      The feed-forward network is an important part of the SEM. The traditional transformer

281  structure uses the fully connected layer as the feed-forward network. To integrate the nearest

282  neighbor relationships between features, convolution structures are combined to process sequences.

283 First, the sequence after the cross-layer residual structure is reconstructed into a feature map by the

284 ToImage function:

$$X_{Patt}^I = ToImage(NormLayer(X_{Patt} + X)) \tag{6}$$

286 Then, through the two-level convolution matrix in Equation 7 and Equation 8:

$$X' = Hardswish(X_{Patt}^I W^1) \tag{7}$$

$$X^2 = Hardswish(Conv(X'))W^2 \tag{8}$$

289 Where $W^1, W^2$ represents the size of the 1×1 weight matrix, and Hardswish is the activation

290 function. Finally, the feature map through the feed-forward network is converted into a sequence

291 format by the function ToSeq and is given to the structure layer in series or the reshaping layer for
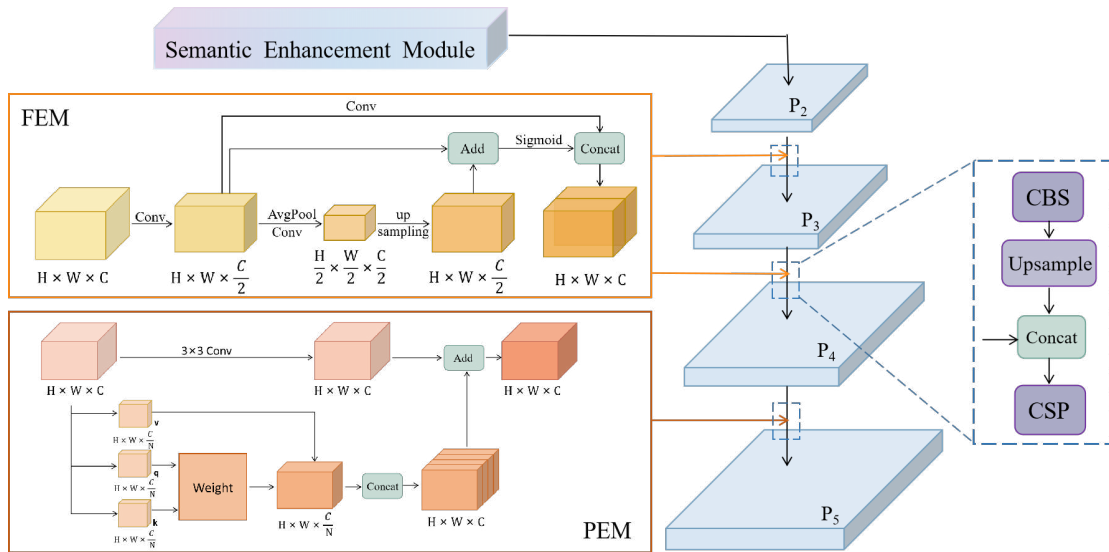
292 subsequent operations:

$$X^{out} = ToSeq(X^2) \tag{9}$$

294 The semantic enhancement module refers to the Transformer structure and constructs the

295 Pooling-MHSA and a new feed-forward neural network to enhance semantic information. The

296 structure composed of NormLayer, Pooling-MHSA, and FFN is stacked three layers deep. The SEM

297 is applied to the top-level feature map in the ODL Net. This is because the top-level feature map

298 often contains more abstract and semantic features, and applying semantic enhancement operations

299 to it can further extract higher-level semantic features. Moreover, applying SEM to the top-level

300 feature map can expand the receptive field, i.e., increase the observation range of each feature point

301 on the input image. Additionally, the self-attention mechanism in SEM helps the algorithm better

302 understand the contextual information and global structure of the targets. In summary, the

303 application of SEM to the top-level feature map can enhance the receptive field, feature extraction

304 capability, and object localization accuracy of ODL Net, thereby improving the performance and

305　　effectiveness of object detection. The experiments also demonstrate that when it is applied to the

306　　top-level feature map, ODL Net achieves the highest detection accuracy, as shown in section 3.2.2

307　　of the experimental results. Furthermore, the module can be independently applied and inserted at

308　　suitable positions, including downstream tasks, within the neural network.

309　　**3.3.2 Feature Fusion Network**

310　　　　The feature fusion network constructs three different enhancement modules to enhance the

311　　information of the feature map before fusion in different aspects. In addition, CBS and CSP modules

312　　are used for feature integration when adjacent feature maps are fused. Figure 8 shows the structure

313　　diagram of the feature fusion network, in which $H \times W$ represents the size of the feature map and

314　　C represents the number of channels.



316　　　　　　　　　　　　　　　　Fig. 8. Diagram of the feature fusion network

317　　　　As shown in Figure 8 above, the feature enhancement module consists mainly of two nested

318　　residual structures. The input feature map is first reduced by a 1×1 convolution to reduce the number

319　　of channels to half of the original number and then further processed by average pooling and

320　　convolution operations with sizes of 2×2 and 3×3, respectively. This step also makes the feature

321    map size half of the input, as well as the number of channels. Then, it restores the size through the

322    upsampling operation and adds it to the feature map before pooling. It is then restored to size by an

323    upsampling operation and added to the feature map before pooling. The last step is to restore the

324    number of channels by stacking the feature maps after addition and before pooling. In the process

325    of halving the resolution of the feature map, more object features will be amplified and extracted by

326    the network. When the image size is restored, the image information will be updated. The cross-

327    layer addition and stacking operation in the FEM effectively avoids the loss of information during

328    image size changes and achieves the function of feature enhancement as a whole.

329        There are also two branches in the Position Enhancement Module (PEM). As shown in Figure

330    8, the upper branch is set with a convolution layer to further extract features on the basis of keeping

331    the size and channel number unchanged. The convolution has a size of 3×3, with a stride and

332    padding of 1. The other branch sets the self-attention mechanism, where N represents the number

333    of attention heads. The PEM mainly enriches the location information through a self-attention

334    mechanism. It provides an effective modeling method through the triplet of Key, Query and Value

335    and obtains greater receptive field and context information by capturing global information. Finally,

336    the captured information will be added with the convolution features to achieve position

337    enhancement.

338        The following describes the position of modules in the feature fusion network. Obviously, the

339    lower-level feature map brings higher resolution, which means it carries more location information.

340    Therefore, the Position Enhancement Module is added in the fusion process of the lowest feature

341    map to supplement the context information. The feature maps in the middle of the two layers are

342    responsible for extracting features. The fusion of these two layers focuses on whether the features

343  are effectively extracted. Therefore, the Feature Enhancement Module is added to the fusion process

344  of $P_2$-$P_4$. Chen (Chen et al., 2021) proved through experiments that the top-level feature map carries

345  the most abundant semantic information among the feature maps generated by the feature extraction

346  network. In addition, the memory requirements of the semantic enhancement module also limit its

347  application scope, so it is added in the process of the top-level feature map C2-P2. To maximize the

348  function of the SEM, this study discusses its reasonable position in the feature fusion network. The

349  experimental data are shown in section 4.3, and the results show that it is best to add SEM to the

350  top layer. The effects of the FEM and PEM are also shown in the same section. The feature fusion

351  network enables ODL Net to fully capture the object features. The enhancement of location and

352  semantic information greatly improved the sensitivity of the algorithm to the feature, which is

353  conducive to the detection of fruit before and after the pear thinning stage.
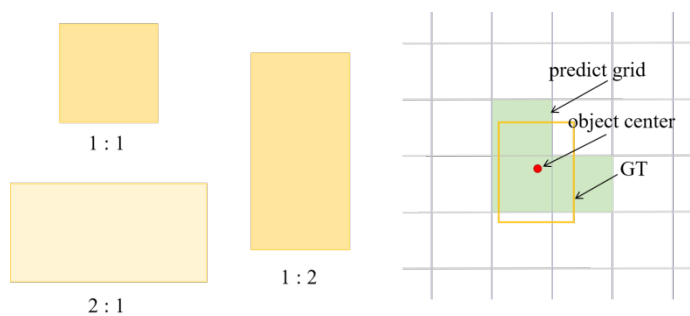
354  **3.4 Detection Head**

355  In this study, a detection head for small objects is constructed, which mainly relies on a special

356  label assignment to improve the detection accuracy. This assignment eliminates attention to the size

357  and shape of the object so that small-scale fruit will not be ignored. The detection head is mainly

358  composed of the label assignment Center-Box and convolution layers, which are shown in Figure

359  3. The following is a description of the label assignment, aiming at small-scale objects.

360  **3.4.1 Label Assignment**

361  Yolov5, as the baseline of this study, is an anchor-based algorithm whose sample selection

362  method increases the number of positive samples to a certain extent. In the feature map, the two

363  adjacent grids closest to the center point of the ground truth are selected as the prediction grids. In

364  addition to the grid where the ground truth is located, there are at most nine anchor boxes

365    corresponding to three grids that match it. In the matching process, the aspect ratio between the

366    ground truth and anchor is calculated twice. If the aspect ratio is less than the specified threshold,

367    the anchor is judged as a positive sample; otherwise, it is the background. For example, if the ground

368    truth is matched with the 1:1 and 1:2 size anchors corresponding to the current layer and its own

369    grid, then there are also two sizes of anchors in the nearest two grids. The number of positive

370    samples of this ground truth in the current layer is 6, while the range of possible anchors is [0, 9],

371    and the number of matching three feature maps is [0, 27]. The process above is shown in Figure 9.

372    Although this label assignment is relatively advanced, which increases the number of positive

373    samples to a certain extent, it cannot be used for small object detection.

374



375    Fig. 9. The label assignment diagram of YOLOV5

376    Note: The left side is the anchor box of three sizes corresponding to each grid, and the right side is the selection

377    diagram of the prediction grid.

378        To improve the detection accuracy of small objects, ODL Net uses a label assignment named

379    Center-Box without anchors, which is specifically described in section 3.4.2. The comparison with

380    other label assignments is shown in Figure 10. The dashed box represents the ground truth, and the

381    orange part represents the positive sample. The two columns on the right in Figure 10 show the

382    representative label assignments of the two types of detection algorithms.
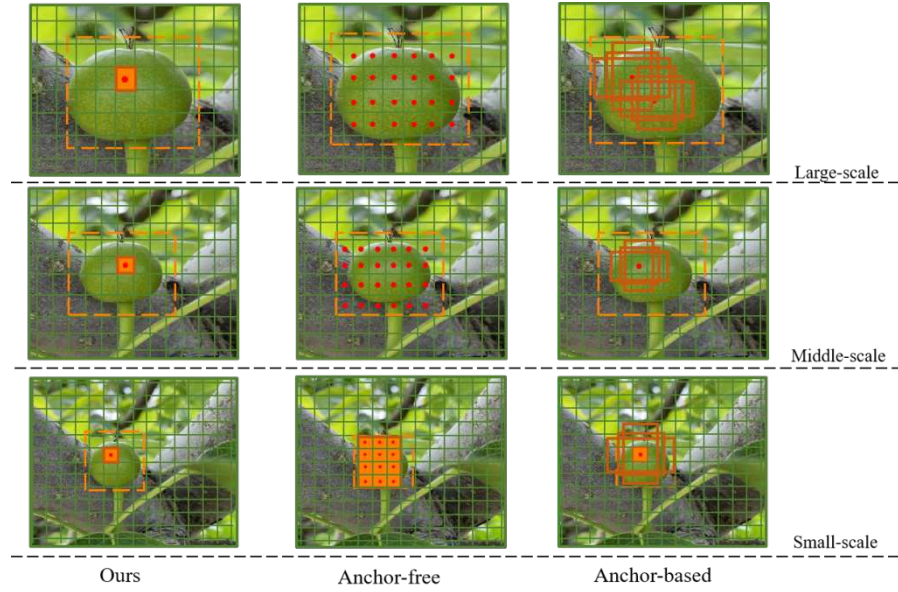
Fig. 10. Comparison of different types of label assignments

383

384

385    The anchor-free algorithms take FCOS as the typical representative and tile the anchor points

386    in the feature map to select positive samples. All anchor points in the ground truth after a feature

387    map is mapped to the original map are selected to calculate the distance from the point to the ground

388    truth: $(l^*, r^*, t^*, b^*)$. FCOS defines the range of $\max(l^*, r^*, t^*, b^*)$ on the multi-scale feature maps to

389    determine the scale on which the object is detected. For example, FCOS stipulates

390    $\max(l^*, r^*, t^*, b^*) \in [128, 256]$ in the top feature map, which means that this feature map is used

391    to detect large objects. The fruit in Figure 10 does not meet this range, so there is no positive sample

392    on the large-scale feature map in this case. Correspondingly, the fruit conforms to the detection

393    range in the small-scale feature map, so the grid of all anchor points in the ground truth is determined

394    as a positive sample.

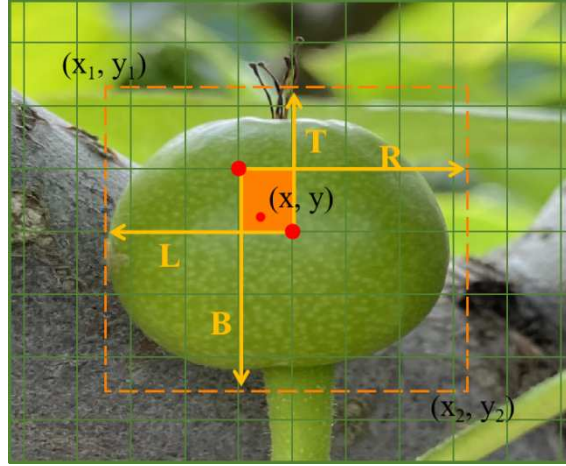395    In the classic anchor-based algorithm, the division of positive and negative samples is

396    completed by calculating the IoU between the ground truth and the anchor boxes. When IoU is

397    greater than the specified threshold, this group of anchors will be determined as positive samples.

398    However, it can be seen from the rightmost column in Figure 10 that the IoU of small objects in the

399     large-scale feature map is zero. This is because the anchor is too different from the ground truth or

400     even completely inside it. Therefore, the small-scale fruit can only produce positive samples in the

401     low-level feature map.

402       These traditional label assignments all define constraints on positive samples, which basically

403     limits the scale range of objects that can be detected at each feature level. The assignments of other

404     algorithms (Kong et al., 2020; Zhu et al., 2019) can also be roughly classified into these two

405     categories. Although these two methods achieve multi-scale detection, larger objects will be

406     allocated more positive samples, and small objects will be easily ignored. This is not conducive to

407     the detection of small objects and makes it difficult to detect fruit around the pear thinning period.

408     **3.4.2 Center-Box**

409       To solve the above problems, ODL Net uses a "fair" label assignment (Zand et al., 2022),

410     Center-Box. As shown on the left of Figure 10, Center-Box cancels the allocation rule of positive

411     and negative samples and directly defines the grid where the object center is located as positive

412     samples (marked in orange) on all levels of feature maps. This strategy prevents the size and shape

413     of objects from dictating the assignment of labels and treats all objects equally at different levels of

414     feature. This means that Center-Box allows the network to learn at all scales of an object, which

415     makes the number of positive samples allocated to small-scale objects and large-scale objects the

416     same. Therefore, the detection will not tend to large-scale objects. To match the positive sample

417     grids, the regression target of Center-Box is defined as the distance from the diagonal vertex of the

418     grid to the ground truth, which is shown in Figure 11 of (L, T, B, R). The coordinates of the upper

419     left corner and the lower right corner of the ground truth are represented as $(x_1, y_1)$ and $(x_2, y_2)$,

420     respectively. The coordinates of the center point are represented as $(x, y)$.

Fig. 11. Diagram of the Center-Box regression

The regression target of Center-Box is the distance between the upper left corner of the grid

where the center point is and the right and upper boundaries of the ground truth and the distance

between the lower right corner of the grid and the left and lower boundaries of the ground truth.

They are represented by $(L^*, T^*, B^*, R^*)$, which is shown in Equation 10:

$$\begin{cases} L^{(i)*} = (x/s_i + 1) - x_1^{(i)}/s_i \\ T^{(i)*} = (y/s_i + 1) - y_1^{(i)}/s_i \\ R^{(i)*} = x_2^{(i)}/s_i - x/s_i \\ B^{(i)*} = y_2^{(i)}/s_i - y/s_i \end{cases} \tag{10}$$

Where i represents the feature scale of [1, 2, 4]. $(\frac{x}{s_i}, \frac{y}{s_i})$ and $(\frac{x}{s_i} + 1, \frac{y}{s_i} + 1)$ in Equation 10

represent the coordinates in the upper left and lower right corners of the grid, respectively. Further

explanation is that $L^{(i)*} + R^{(i)*} = (x_2^{(i)} - x_1^{(i)}) + 1$, $T^{(i)*} + B^{(i)*} = (y_2^{(i)} - y_1^{(i)}) + 1$, where $x_2^{(i)} -$

$x_1^{(i)} = w^i = \frac{w}{s_i}$, and $y_2^{(i)} - y_1^{(i)} = h^i = \frac{h}{s_i}$. w and h represent the size of the ground truth in the

original image, and $w^i$ and $h^i$ represent the width and height of the ground truth on scale i,

respectively.

The learning process of the regression target is shown in Equation 11:

$$\begin{cases} L^{(i)} = (\alpha \times \text{Sigmoid}(l))^2 * 2^i \\ T^{(i)} = (\alpha \times \text{Sigmoid}(t))^2 * 2^i \\ R^{(i)} = (\alpha \times \text{Sigmoid}(r))^2 * 2^i \\ B^{(i)} = (\alpha \times \text{Sigmoid}(b))^2 * 2^i \end{cases} \tag{11}$$

436 Where (l, t, r, b) represent the predicted values in network for the distance in four directions, and

437 their values are controlled between 0 and 1 by the sigmoid function. $i \in \{1, 2, 4\}$, represents the

438 scale of different feature maps, and $2^i$ is used to distinguish different scales in the learning

439 process. $\alpha$ is a range constant used to expand the detection coverage. It is set as 1.0 in the experiment

440 because of the small size of most objects around the pear thinning period, and it can be adjusted

441 according to the size of the object in other studies. $(L^{(i)}, T^{(i)}, R^{(i)}, B^{(i)})$ is the predicted result on

442 the i-th layer of feature map. This predicted distance is used to compare with the real distance and

443 adjust the network parameters according to loss function for learning.

444 The Center-Box approach does not specifically aim to detect fruits of corresponding sizes on

445 feature maps of different scales, but rather ensures that fruits can be learned on feature maps of all

446 scales where they exist. It directly assigns the grid cell containing the center of the fruit as a positive

447 sample and regresses the distances from the top-left to the bottom-right corners of the grid cell to

448 the true box. This strategy allows for an equal number of positive samples to be assigned to both

449 large-scale and small-scale fruits, enabling the ODL Net to treat the detection of fruits at different

450 scales equally. Consequently, this implicitly enhances the detection capability of the ODL Net for

451 small-scale fruits. In the real working environment of the ODL Net, as supported by the statistical

452 information provided in Section 2, small-scale fruits constitute nearly half of the overall quantity.

453 Hence, the scale-agnostic nature of the Center-Box approach empowers the ODL Net to deliver

454 satisfactory performance in detection tasks before and after thinning in pear orchards.

455 **3.5 Loss Function**

456      The loss function of ODL Net consists of three parts: classification loss, confidence loss and

457      bounding box loss. The network loss is the weighted sum of the above three, which is shown in

458      Equation 12. The impact of each loss can be adjusted by weight $\lambda$.

$$\text{Loss} = \lambda_1 L_{cls} + \lambda_2 L_{conf} + \lambda_3 L_{obj} \tag{12}$$

460      **3.5.1 Classification and Confidence Loss**

461      For detection tasks in the pear orchard, only "pear" is the category of prediction tag output

462      from the network. At this point, the common binary cross entropy loss BCE with logit loss is used

463      as the classification loss:

$$y_i = \text{Sigmoid}(x_i) = \frac{1}{1+e^{-x_i}} \tag{13}$$

$$L_{cls} = -\sum_{n=1}^{N} y_i^* \log(y_i) + (1 - y_i^*)\log(1 - y_i) \tag{14}$$

466      Where $x_i$ represents the predicted value of the current category. $y_i$ represents the probability of the

467      current category obtained after activating the function. $y_i^*$ is the true value of the class, expressed

468      as 0 or 1.

469      The confidence level of the prediction box indicates its reliability. The higher the value, the

470      more reliable the prediction box is, and the closer it is to the ground truth. The confidence loss is

471      the same type as the classification loss, using the binary cross entropy loss. It should be noted that

472      the total confidence loss is obtained by weighted addition of the confidence losses on the three

473      prediction branches:

$$L_{conf} = \beta_1 L_{conf}^L + \beta_2 L_{conf}^M + \beta_3 L_{conf}^S \tag{15}$$

475      Where $\beta_1, \beta_2, \beta_3$ represents the influence of the confidence loss of the feature map with the

476      resolution from high to low. The weight is set to (5.0, 1.0, 0.5) in the experiment to improve the

477      detection accuracy of small fruits. Because small-scale objects are detected on the high-resolution

478   feature map, $\beta_1$ is adjusted higher to facilitate small fruit detection.

479   **3.5.2 Bounding Box Loss**

480       The goal in the process of boundary box regression is to minimize the distance between the

481   prediction box and the ground truth. For the relative position of the bounding boxes, the classic

482   method is to calculate the IoU value of the two boxes. IoU is usually used to express the coincident

483   area of two object positions. On this basis, many more advanced methods have been proposed

484   (Rezatofighi et al., 2019; Zheng et al., 2020). For example, DIoU takes into account the distance

485   between the ground truth and the prediction box, the overlap rate and the scale:

486   $$\mathrm{DIoU} = \mathrm{IoU} - \frac{\rho^2(b, b^{gt})}{c^2} \tag{16}$$

487   Where $b, b^{gt}$ represents the center of the prediction box and the ground truth, respectively.

488   $\rho(b, b^{gt})$ represents the Euclidean distance between two central points. The symbol c represents

489   the diagonal distance of the smallest area that can contain both boxes.

490       In this study, to cooperate with the Center-Box, a loss method with scale invariance is proposed.

491   The goal in the regression process is to minimize the distance between the prediction box and the

492   ground truth. As explained in Section 3.4.2, each box is represented by four distances. Therefore, it

493   is the   hope of this study that the distance in four directions can be taken into account in the loss

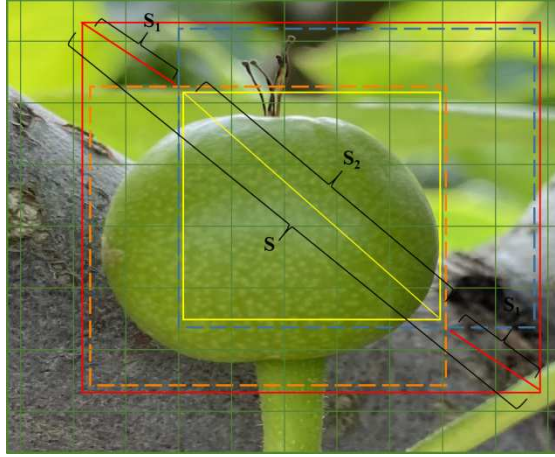494   of the bounding box, which is shown in Figure 12.

Fig. 12. Calculation diagram of bounding box loss

497    In the loss of bounding boxes, the overlapping area (yellow box), non-overlapping area and

498    minimum inclusion area (red box) are all considered. They are expressed in square Euclidean

499    Distance as Equation 17:

500
$$\begin{cases} S_1 = (L^* - L)^2 + (T^* - T)^2 + (R^* - R)^2 + (B^* - B)^2 \\ S_2 = (\min(L^*, L) + \min(R^*, R) - 1)^2 + (\min(T^*, T) + \min(B^*, B) - 1)^2 \quad (17) \\ S = (\max(L^*, L) + \max(R^*, R) - 1)^2 + (\max(T^*, T) + \max(B^*, B) - 1)^2 \end{cases}$$

501    Where $(L^*, T^*, B^*, R^*)$ and $(L, T, B, R)$ represent the predicted and true values, respectively. The

502    expression of bounding box loss is as Equation 18:

503
$$L_{obj}(L^*, T^*, R^*, B^*) = 1 - \frac{(S_2 - S_1)}{S} \qquad (18)$$

504    ## 4. Experiments

505    The experiment is conducted on a server equipped with the Ubuntu 16.04 operating system,

506    which is equipped with four GTX 3090 graphics cards and V11.4 CUDA. During the training, two

507    graphics cards are used, and 16 images are set for each batch. The initial learning rate of the

508    experiment is set to 0.005, and 0.0001 is used as the weight attenuation to prevent over-fitting. To

509    update and calculate the network parameters and minimize the loss function, a random gradient

510    descent (SGD) optimizer with a momentum of 0.9 is used to assist the training. Image enhancement

511    methods are used to enrich the dataset before training to reduce over-fitting. Finally, 300 epochs are

512     trained for the ODL Net.

## 4.1 Evaluation Index

514     The average precision (AP) is selected as the evaluation index of algorithm performance in the

515     experiment. It is the area under the PR curve with recall as the horizontal axis and precision as the

516     vertical axis. The calculate method is shown in Equation 19. Other evaluation indicators used in the

517     experiment also belong to the same type: $AP_{50}$ is the measured value of AP when the IOU threshold

518     is 0.5; $AP_{75}$ is the AP measurement value when IOU is 0.75; $AP_s$, $AP_m$ and $AP_l$ represent AP

519     measurement values of small, medium and large objects, respectively. Fruit with the number of

520     pixels less than 174×174 are defined as small-scale objects, fruit with the number of pixels greater

521     than 523×523 are defined as large-scale objects, and fruit with the number of pixels between them

522     are defined as medium-scale objects. The formulation of the scale range is explained in section 2.2.

$$AP = \int_0^1 P(R)dR \tag{19}$$

524     In this definition of AP, P is represents the proportion of the number of predicted positive

525     samples to the number of real positive samples; R represents the proportion of positive samples

526     correctly predicted by the algorithm in the real positive samples. The calculation equations are

527     shown in Equation 20 and Equation 21, where TP represents the number of detection frames whose

528     intersection to parallel ratio is greater than the set threshold; FP represents the number of detection

529     frames whose intersection ratio is less than the set threshold, or the number of redundant detection

530     frames generated under the same target; FN indicates the number of targets not detected.

$$\text{Precision} = \frac{TP}{TP+FP} \times 100\% \tag{20}$$

$$\text{Recall} = \frac{TP}{TP+FN} \times 100\% \tag{21}$$

533     In addition, the average recall (AR) is also a supplementary evaluation index, although AP is

534     more authoritative. AR refers to the maximum recall in a given number of detection results on each

535     image.

536     **4.2 Comparative Experiments**

537        In this section, a comparison is made between ODL Net and classical CNN-based detection

538     algorithms since 2020. The detection accuracy on the dataset prior to thinning the pears is presented

539     in Table 2. Two key results are emphasized in the experiment: the overall detection accuracy of the

540     algorithm (referred to as AP) and the detection accuracy specifically for small-scale fruits (referred

541     to as APs). Table 2 reveals that ODL Net achieves the highest detection accuracy, reaching 56.2%,

542     surpassing other algorithms by a margin of at least 0.5 percentage points. Among the detection

543     algorithms developed in the past two years, AutoAssign (Zhu et al., 2020) demonstrates the closest

544     accuracy to ODL Net for small fruits, with a mere 0.3 percentage point difference. However, its

545     overall detection accuracy is unsatisfactory. Similarly, the AP of TOOD (Feng et al., 2021) reaches

546     55.7%, but its performance on small fruits falls significantly behind our algorithm.

547        From the aforementioned results, it is evident that one of the key advantages of ODL Net lies

548     in its capability to enhance the detection accuracy of small-scale fruits without compromising its

549     overall AP. Building upon the YOLOV5 baseline, ODL Net exhibits an increase of 1.4 percentage

550     points in AP and a 2.1 percentage point increase in APs. In contrast, the other algorithms in Table 2,

551     such as NAS-FCOS (Wang et al., 2020), an enhanced version of FCOS, yield considerably lower

552     accuracy compared to ODL Net, differing by more than 3.0 percentage points. Additionally, ODL

553     Net achieves the highest AR of 61.6% and ARs of 40.0%.

554            Table 2. Comparative experiments on the pear dataset before thinning

| | AP% | $AP_{50}$% | $AP_{75}$% | $AP_s$% | $AP_m$% | $AP_l$% | AR% | $AR_s$% | $AR_m$% | $AR_l$% |
|---|---|---|---|---|---|---|---|---|---|---|

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| ATSS (Zhang et al., 2020) | 50.3 | 79.0 | 50.8 | 22.6 | 72.3 | 85.4 | 56.3 | 33.4 | 77.2 | 89.1 |
| AutoAssign (Zhu et al., 2020) | 52.8 | 84.2 | 53.3 | 29.0 | 75.6 | 87.6 | 60.8 | 39.3 | 80.5 | 91.3 |
| Double-Head RCNN (Wu et al., 2020) | 51.8 | 83.8 | 54.1 | 26.6 | 71.9 | 79.9 | 56.7 | 36.5 | 75.6 | 83.3 |
| NAS-FCOS (Wang et al., 2020) | 53.1 | 82.9 | 54.4 | 28.3 | 73.2 | 86.8 | 60.0 | 38.8 | 79.3 | 91.0 |
| TOOD (Feng et al., 2021) | 55.7 | 84.7 | 56.4 | 28.7 | 76.8 | 90.2 | 61.6 | 40.0 | 81.3 | 93.3 |
| YOLOV5 | 54.8 | 81.6 | 56.3 | 27.2 | 75.8 | 88.9 | 61.2 | 39.8 | 80.5 | 92.4 |
| ODL Net | 56.2 | 83.0 | 57.5 | 29.3 | 77.0 | 91.2 | 61.3 | 39.0 | 81.8 | 93.8 |

555     The experimental results on the pear dataset after fruit thinning are presented in Table 3. The

556     dataset exhibits a significantly lower fruit density compared to the pre-thinning dataset, with

557     minimal instances of fruit overlap. In this scenario, ODL Net demonstrates a notable improvement

558     over YOLOV5, with an increase of 2.4 percentage points in both AP and APs. Similarly, for ATSS,

559     which exhibits relatively good performance, there is a 2.0 percentage point increase in both AP and

560     APs. But in terms of detecting small-scale fruits after pear thinning, Double-Head RCNN and NAS-

561     FCOS achieve comparable or slightly higher detection accuracy compared to ODL Net. However,

562     these algorithms tend to focus excessively on small objects and lack sensitivity towards objects of

563     other scales, resulting in an AP that is 3.6-5.0 percentage points lower than ODL Net. They prioritize

564     small object detection at the expense of AP.

565     The experiments demonstrate that ODL Net enhances the detection accuracy of small objects

566     while also considering the overall detection accuracy (AP) of the algorithm. This is because SEM

567     enlarges the receptive field of the network, enabling the network to perceive fruits at all scales. In

568     addition, effective feature fusion also enables ODL Net to capture more and richer features for

569     accuracy detection.

570                  Table 3. Comparative experiments on the pear dataset after thinning

| | AP% | $AP_{50}$% | $AP_{75}$% | $AP_s$% | $AP_m$% | $AP_l$% | AR% | $AR_s$% | $AR_m$% | $AR_l$% |
|---|---|---|---|---|---|---|---|---|---|---|
| ATSS (Zhang et al., 2020) | 63.1 | 79.5 | 70.0 | 36.9 | 81.7 | 90.1 | 69.3 | 51.6 | 85.1 | 92.2 |
| AutoAssign (Zhu et al., 2020) | 60.0 | 76.9 | 68.0 | 37.8 | 77.7 | 86.3 | 70.6 | 55.1 | 84.6 | 90.0 |
| Double-Head RCNN (Wu et al., 2020) | 61.5 | 80.5 | 72.0 | 40.4 | 75.7 | 83.5 | 66.4 | 52.9 | 77.6 | 85.4 |
| NAS-FCOS (Wang et al., 2020) | 60.1 | 78.1 | 68.8 | 38.9 | 75.5 | 86.0 | 67.3 | 51.6 | 80.9 | 88.6 |
| TOOD (Feng et al., 2021) | 62.3 | 79.4 | 67.6 | 37.7 | 61.7 | 88.3 | 68.2 | 48.4 | 69.3 | 91.9 |
| YOLOV5 | 62.7 | 77.4 | 68.5 | 36.5 | 80.6 | 88.7 | 69.3 | 51.2 | 86.2 | 91.0 |
| ODL Net | 65.1 | 78.6 | 70.4 | 38.9 | 81.7 | 91.2 | 70.8 | 53.9 | 86.1 | 92.5 |

571 **4.3 Ablation Experiments**

572     Considering the experimental nature of the dataset, the ablation experiments on the pear dataset

573     after thinning can fully and clearly show the role of each module. The experimental data are shown

574     in Table 4. The addition of the Center-Box focuses on improving the detection accuracy of small-

575     scale objects. At this time, the accuracy of medium-scale and large-scale objects is almost

576 unchanged. It is also the semantic enhancement module for small objects, which further improves

577 the detection accuracy of small fruits. They can be used separately in the detection of other small

578 objects. In addition, the feature enhancement module and the position enhancement module are very

579 helpful in improving the overall detection accuracy. Compared with the SEM, although they are not

580 aimed at small objects, they improve the overall detection accuracy of ODL Net. Finally, the AP and

581 APs of ODL Net are 2.4 percentage points higher than those of the baseline algorithm. The detection

582 accuracy of the algorithm for small-scale objects in the dataset is up to 39.3%, although the overall

583 accuracy is not the highest at this time.

584 <div align="center">Table 4. Ablation experiments on the pear dataset after the thinning period</div>

| Structure | | | | AP% | $AP_s$% | $AP_m$% | $AP_l$% |
|---|---|---|---|---|---|---|---|
| Center-Box | SEM | PEM | FEM | | | | |
| × | × | × | × | 62.7 | 36.5 | 80.6 | 88.7 |
| √ | × | × | × | 63.8 | 38.9 | 80.6 | 89.3 |
| √ | √ | × | × | 64.1 | 39.3 | 80.1 | 89.7 |
| √ | × | √ | × | 64.8 | 38.2 | 81.2 | 88.1 |
| √ | √ | √ | × | 64.9 | 38.5 | 82.8 | 91.0 |
| √ | √ | √ | √ | 65.1 | 38.9 | 81.7 | 91.2 |

585 What needs to be specially explained is the location of the semantic enhancement module in

586 ODL Net. In fact, it is obvious that the top-level feature map brings the richest semantic feature.

587 However, we still confirm the location of SEM through experiments to eliminate the impact of the

588 pear dataset. The experimental data are shown in Table 5, which shows the accuracy of ODL Net

589 when SEM is added at different feature map layers. At this time, ODL Net is only constructed with

590 Object-Box and SEM, but no other modules. Table 5 shows that the accuracy of ODL Net reaches

591 64.1% when SEM is added to the top level of the feature fusion network. The detection of small

592 objects reaches 39.3%. Therefore, in this study, SEM is finally added to the top-level feature map

593 $C_2$ to enhance the detection accuracy of small-scale fruit.

594

Table 5. Ablation experiments of the layer with SEM

| Layer | AP | $AP_{50}$ | $AP_{75}$ | $AP_s$ | $AP_m$ | $AP_l$ |
|-------|------|------|------|------|------|------|
| $C_2$ | 64.1 | 77.9 | 69.7 | 39.3 | 80.1 | 89.7 |
| $C_3$ | 62.9 | 77.6 | 68.5 | 36.9 | 79.4 | 89.7 |
| $C_2, C_3, C_4$ | 63.6 | 77.4 | 69.3 | 38.1 | 81.0 | 90.1 |

## 595 4.4 Sample Results

596 The detection effect of the algorithms on the pear dataset is shown in Figure 13 and Figure 14

597 before and after fruit thinning, respectively. A representative image is selected to show the effect

598 before the thinning period. The area with objects (marked with orange rectangle in the original

599 image) is enlarged in all of the result images to display the detection results more intuitively. In the

600 image before fruit thinning, there are ten fruits to be detected in the selected image of Figure 13.

601 Box redundancy occurs in the Yolov5, AutoAssign and Double-Head RCNN. That is, there are

602 multiple prediction boxes on a fruit, or nonexistent objects are detected. ATSS and NAS-FCOS miss

603 approximately three numbers of the target fruit. While, TOOD successfully detected all fruits. ODL
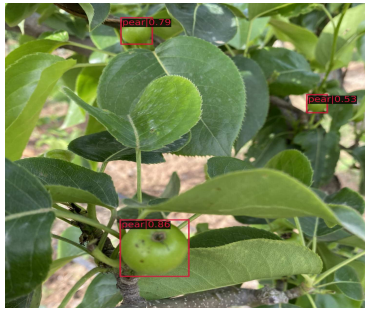
604 Net recognizes nine fruits with the highest scores.

original image

ATSS

AutoAssign

Double-Head
RCNN

NAS-FCOS

TOOD

Fig. 13. Comparison images of algorithms on the pear dataset before thinning

605

606　　　The fruit density after thinning is relatively sparse, and there is basically no problem of box

607　　redundancy when detected. The decrease in fruit density also makes the detection easier. However,

608　　for the incomplete and fuzzy fruit in the lower right corner in the first image of Figure 14, most of

609　　the algorithms fail to detect it. The other two detected images after thinning are also shown in Figure

610　　14.

AutoAssign

Double-Head
RCNN

NAS-FCOS

TOOD

YOLOV5

ODL Net



611                    Fig. 14. Comparison images of algorithms on the pear dataset after thinning

612           The other detected images before thinning are shown in Figure 15 as examples.
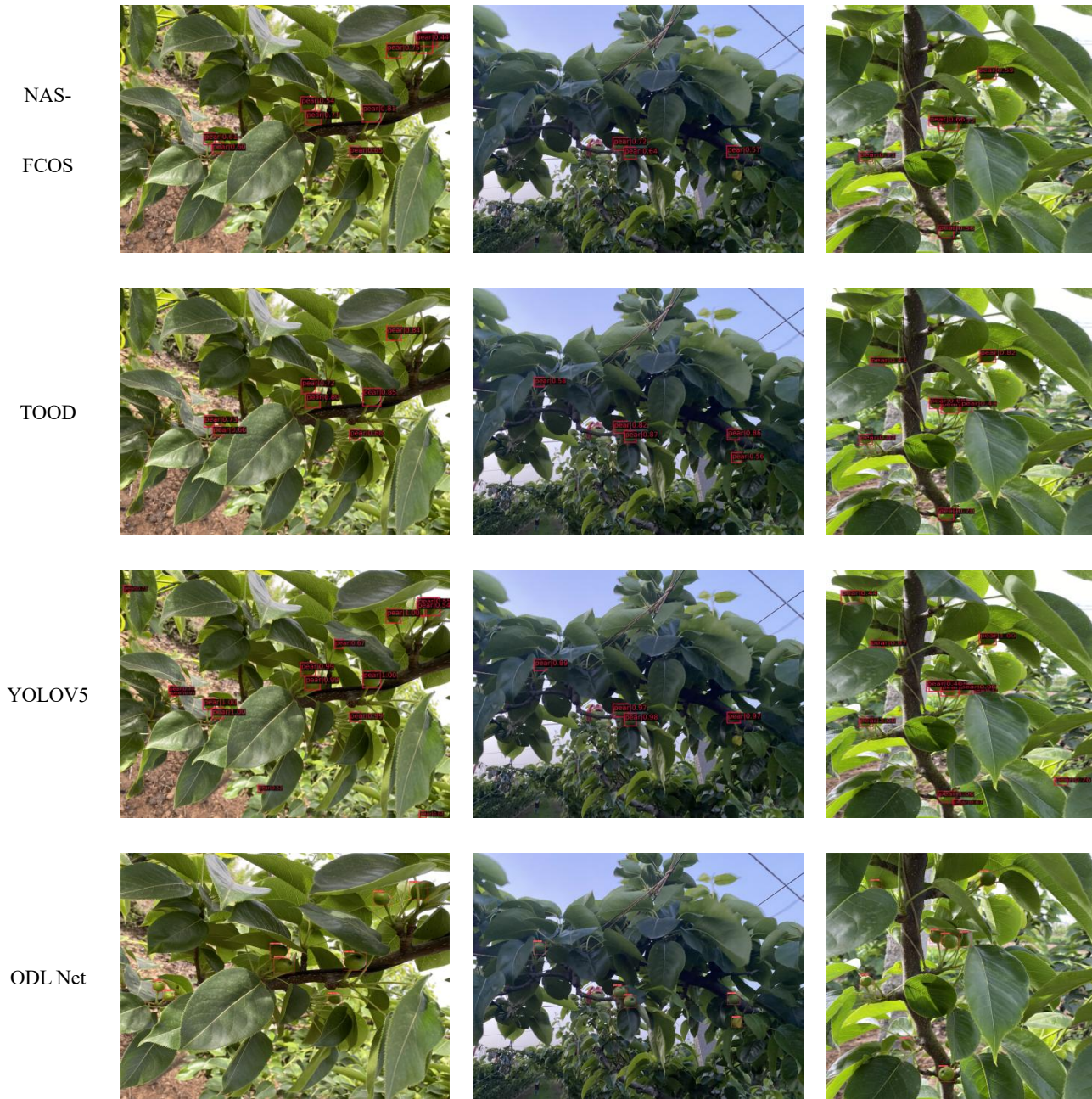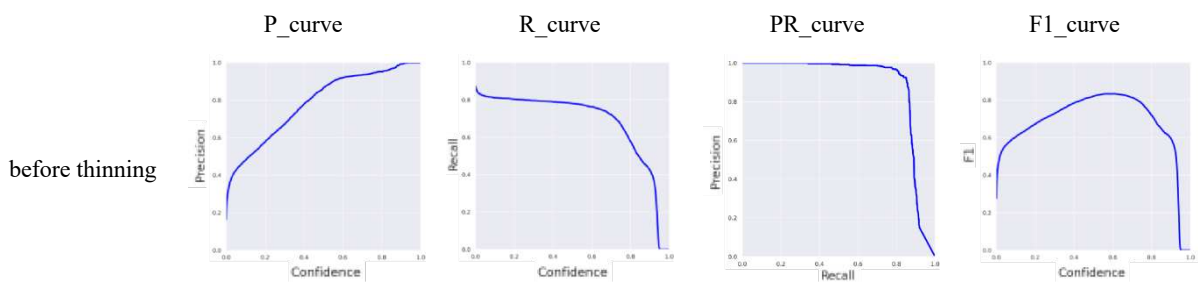
original image

ATSS

AutoAssign
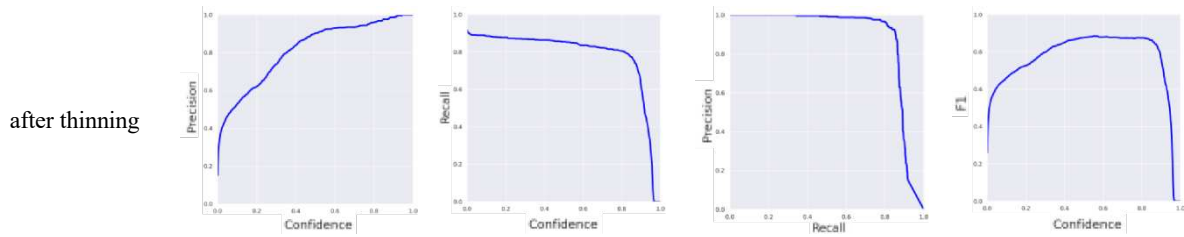
Double-Head RCNN

613                    Fig. 15. Comparison images of algorithms on the pear dataset around the thinning period

614          The curves of ODL Net during training are shown in Figure 16.

after thinning

Fig. 16. Curves of ODL Net during training

## 5. Conclusion and Feature Work

This study proposes a detection algorithm called ODL Net, specifically designed for detecting small-scale fruits before and after thinning in pear orchards. It enhances the detection of small objects through the SEM and the label assignment strategy called Center-Box. Additionally, the modules of FEM and PEM are constructed to further improve the overall detection performance of ODL Net. These enhancement modules can also be used individually.

For fruit varieties that require thinning, the detection of fruits by ODL Net before the thinning stage can guide the thinning process. Moreover, the detection of fruits by ODL Net after thinning enables calculations for irrigation and fertilizer requirements, facilitates scientific yield measurement, and supports intelligent management of orchards. In the case of fruit varieties that do not require thinning, ODL Net provides continuous monitoring throughout the fruit growth period. Particularly in the early stages of fruit growth, where the fruit size is small and detection poses significant challenges, a high-performance detection algorithm is crucial. ODL Net, designed specifically for small-scale fruits, can partially address this issue. And It is precisely because of this characteristic that ODL Net offers significant assistance in intelligent orchard management and fills the research gap in various small-scale fruits detection during the thinning period in orchards.

Although ODL Net has completed the improvement for small-scale objects, the difficulties have not been completely overcome. The complete green appearance of pear during fruit thinning

still caused some difficulties in the detection work. In future studies, we hope to propose a detection

algorithm for both indistinguishable green color and small size.

## Acknowledgments

## References

Bargoti S, Underwood J. Deep fruit detection in orchards. 2017 IEEE international conference on robotics and automation (ICRA). IEEE, 2017: 3626-3633.

Chen Q, Wang Y, Yang T, et al. You only look one-level feature. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2021: 13039-13048.

Ebrahimi M, Khoshtaghaza M, Minaei S, et al. Vision-based pest detection based on SVM classification method. Computers and Electronics in Agriculture, 2017, 137: 52-58.

Elfwing S, Uchibe E, Doya K. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. Neural Networks, 2018, 107: 3-11.

Feng C, Zhong Y, Gao Y, et al. Tood: Task-aligned one-stage object detection. IEEE/CVF International Conference on Computer Vision (ICCV). IEEE Computer Society, 2021: 3490-3499.

Fu L, Gao F, Wu J, et al. Application of consumer RGB-D cameras for fruit detection and localization in field: A critical review. Computers and Electronics in Agriculture, 2020, 177: 105687.

Gongal A, Amatya S, Karkee M, et al. Sensors and systems for fruit detection and localization: A review. Computers

656      and Electronics in Agriculture, 2015, 116: 8-19.

657      Hussain D, Hussain I, Ismail M, et al. A simple and efficient deep learning-based framework for automatic fruit

658      recognition. Computational Intelligence and Neuroscience, 2022, ID 6538117.

659      Jia W, Meng u, Ma X, et al.. Efficient detection model of green target fruit based on optimized Transformer

660      network.Transactions of the Chinese Society of Agricultural Engineering, 2021, 37(14): 163-170.

661      Rabbi J, Ray N, Schubert M, et al. Small-object detection in remote sensing images with end-to-end edge-enhanced

662      GAN and object detector network[J]. Remote Sensing, 2020, 12(9): 1432.

663      Kong T, Sun F, Liu H, et al. Foveabox: Beyound anchor-based object detection. IEEE Transactions on Image

664      Processing, 2020, 29: 7389-7398.

665      Lin T Y, Dollár P, Girshick R, et al. Feature pyramid networks for object detection. Proceedings of the IEEE

666      conference on computer vision and pattern recognition. 2017: 2117-2125.

667      Liu S, Qi L, Qin H, et al. Path aggregation network for instance segmentation. Proceedings of the IEEE conference

668      on computer vision and pattern recognition. 2018: 8759-8768.

669      Liu Z, Hu H, Lin Y, et al. Swin transformer v2: Scaling up capacity and resolution. Proceedings of the IEEE/CVF

670      Conference on Computer Vision and Pattern Recognition. 2022: 12009-12019.

671      Liu Z, Lin Y, Cao Y, et al. Swin transformer: Hierarchical vision transformer using shifted windows. Proceedings of

672      the IEEE/CVF International Conference on Computer Vision. 2021: 10012-10022.

673      Maheswari P, Raja P, Apolo-Apolo O E, et al. Intelligent fruit yield estimation for orchards using deep learning based

674      semantic segmentation techniques—a review. Frontiers in Plant Science, 2021, 12: 684328.

675      Mai X, Zhang H, Meng M. Faster R-CNN with classifier fusion for small fruit detection. IEEE International

676      Conference on Robotics and Automation (ICRA). IEEE, 2018: 7166-7172.

677      Ngugi LC, Abelwahab M, Abo-Zahhad M. Recent advances in image processing techniques for automated leaf pest

and disease recognition–A review. Information processing in agriculture, 2021, 8(1): 27-51.

Patrício D I, Rieder R. Computer vision and artificial intelligence in precision agriculture for grain crops: A systematic review. Computers and electronics in agriculture, 2018, 153: 69-81.

Rezatofighi H, Tsoi N, Gwak J Y, et al. Generalized intersection over union: A metric and a loss for bounding box regression. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2019: 658-666.

Sa I, Ge Z, Dayoub F, et al. Deepfruits: A fruit detection system using deep neural networks. sensors, 2016, 16(8): 1222.

Sun M, Xu L, Chen X, et al. Bfp net: balanced feature pyramid network for small apple detection in complex orchard environment. Plant Phenomics, 2022, 2022.

Tang, Y., Zhou, H., Wang, H., et L.. Fruit detection and positioning technology for a Camellia oleifera C. Abel orchard based on improved YOLOv4-tiny model and binocular stereo vision, Expert Systems with Applications 2023, 211:118573.

Tey Y S, Brindal M. A meta-analysis of factors driving the adoption of precision agriculture. Precision Agriculture, 2022, 23(2): 353-372.

Tu S, Pang J, Liu H, et al. Passion fruit detection and counting based on multiple scale faster R-CNN using RGB-D images. Precision Agriculture, 2020, 21(5): 1072-1091.

Wang N, Gao Y, Chen H, et al. Nas-fcos: Fast neural architecture search for object detection. proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2020: 11943-11951.

Wu H T, Tsai C W. An intelligent agriculture network security system based on private blockchains. Journal of Communications and Networks, 2019, 21(5): 503-508.

Wu Y, Chen Y, Yuan L, et al. Rethinking classification and localization for object detection. Proceedings of the

IEEE/CVF conference on computer vision and pattern recognition. 2020: 10186-10195.

Xu B, Cui X, Ji W, et al. Apple grading method design and implementation for automatic grader based on Improved

   YOLOv5. Agriculture, 2023,13,124.

Yang L, Chen Y, Tian Z, et al. Field road segmentation method based on improved UNet. Transactions of the Chinese

   Society of Agricultural Engineering, 2021, 37(09): 185-191. (in Chinese)

Yang X, Yang J, Yan J, et al. Scrdet: Towards more robust detection for small, cluttered and rotated objects.

   Proceedings of the IEEE/CVF International Conference on Computer Vision. 2019: 8232-8241.

Yu F, Wang D, Shelhamer E, et al. Deep layer aggregation. Proceedings of the IEEE conference on computer vision

   and pattern recognition. 2018: 2403-2412.

Zand M, Etemad A, Greenspan M. Objectbox: From centers to boxes for anchor-free object detection. European

   Conference on Computer Vision. Springer, Cham, 2022: 390-406.

Zhang S, Chi C, Yao Y, et al. Bridging the gap between anchor-based and anchor-free detection via adaptive training

   sample selection. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2020:

   9759-9768.

Zhang W, Wang S, Thachan S, et al. Deconv R-CNN for small object detection on remote sensing images. IEEE

   International Geoscience and Remote Sensing Symposium. IEEE, 2018: 2483-2486.

Zhao K, Yan W Q. Fruit detection from digital images using CenterNet. International Symposium on Geometry and

   Vision. Springer, Cham, 2021: 313-326.

Zheng Z, Wang P, Liu W, et al. Distance-IoU loss: Faster and better learning for bounding box regression.

   /Proceedings of the AAAI conference on artificial intelligence. 2020, 34(07): 12993-13000.

Zhu B, Wang J, Jiang Z, et al. Autoassign: Differentiable label assignment for dense object detection. arXiv preprint

   arXiv:2007.03496, 2020.

722    Zhu C, He Y, Savvides M. Feature selective anchor-free module for single-shot object detection. Proceedings of the

723    IEEE/CVF conference on computer vision and pattern recognition. 2019: 840-849.

724    Zhu X, Su W, Lu L, et al. Deformable detr: Deformable transformers for end-to-end object detection. arXiv preprint

725    arXiv:2010.04159, 2020.