

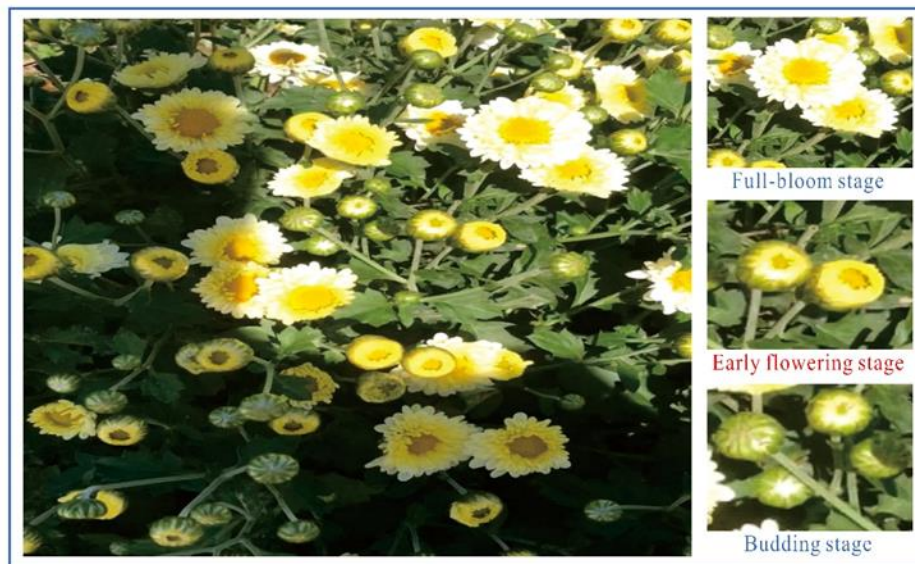
22 **Abstract:** Tea chrysanthemum detection at its flowering stage is one of the key components for
23 selective chrysanthemum harvesting robot development. However, it is a challenge to detect
24 flowering chrysanthemums under unstructured field environments given variations on illumination,
25 occlusion and object scale. In this context, we propose a highly fused and lightweight deep learning
26 architecture based on YOLO for tea chrysanthemum detection (TC-YOLO). First, in the backbone
27 component and neck component, the method uses the Cross-Stage Partially Dense network
28 (CSPDenseNet) and the Cross-Stage Partial ResNeXt network (CSPResNeXt) as the main networks,
29 respectively, and embeds custom feature fusion modules to guide the gradient flow. In the final head
30 component, the method combines the recursive feature pyramid (RFP) multiscale fusion reflow
31 structure and the Atrous Spatial Pyramid Pool (ASPP) module with cavity convolution to achieve
32 the detection task. The resulting model was tested on 300 field images using a data enhancement
33 strategy combining flipping and rotation, showing that under the NVIDIA Tesla P100 GPU
34 environment, if the inference speed is 47.23 FPS for each image (416×416), TC-YOLO can achieve
35 the average precision (AP) of 92.49% on our own tea chrysanthemum dataset. Through further
36 validation, it was found that overlap had the least effect on tea chrysanthemum detection, and
37 illumination had the greatest effect on tea chrysanthemum detection. In addition, this method
38 (13.6M) can be deployed on a single mobile GPU, and it could be further developed as a perception
39 system for a selective chrysanthemum harvesting robot in the future.

40 **Keywords:** Tea chrysanthemum; Flowering stage detection; Deep convolutional neural network;
41 Agricultural robotics

42 1. Introduction

43 Current studies show that tea chrysanthemums have significant commercial value (Liu et al.,

44 2020; Liu et al., 2019). Not only that, but tea chrysanthemums can offer a range of health benefits
45 (Hou et al., 2017; Yue et al., 2018). For example, they can significantly inhibit the activity of
46 carcinogens and have distinct anti-aging, cholagogic and antihypertensive effects (Zheng et al.,
47 2021). In addition, apigenin and naringenin in tea chrysanthemums can effectively reduce or even
48 reverse the prevalence of obesity (Thaiss et al., 2016). In natural environments, a tea
49 chrysanthemum plant could present multiple flower heads varying in different stages and sizes. The
50 examples of different flowering stages of chrysanthemums are shown in Fig. 1.



51
52
53
54
55
56
57
58 **Fig. 1.** A typical tea chrysanthemum plant with three different flowering stages.

59 At present, tea chrysanthemum is mainly manually harvested at the early flowering stage, and
60 this process is labor-intensive and time-consuming. The development of artificial intelligence and
61 selective harvesting robots will help to prevent crop wastage due to the current shortages in skilled
62 labourers (Li et al., 2021). Therefore, an automated selective harvesting robot is required by industry.
63 Typically, a robotic harvesting process is divided into two steps. First, a computer vision system is
64 used to detect the objects of interest. Then the manipulator moves towards the detected objects and
65 the gripper harvests the objects guided by the detection results. Accurately detecting the

66 chrysanthemum at the early flowering stage is critical for guiding the manipulator to the correct
67 position. Although methods based on machine learning technology have made remarkable
68 achievements in agricultural applications (Gao et al., 2020), it is still difficult to develop a
69 lightweight network for a selective harvesting robot under unstructured environments.

70 The study on chrysanthemum detection based on machine vision began as early as 1996. Kondo
71 et al. (1996) developed a system for cutting chrysanthemums. The chrysanthemums were detected
72 using the direction, position, size and shape features. Unfortunately, this study failed to provide
73 definite detection results. Subsequently, Warren (2000) developed a system for identifying
74 chrysanthemum leaves through extracting image features, including the length and width of the
75 leaves, and the shape of the leaf apices. For the system, the precision of detection was up to 75%,
76 and generally, its running speed was roughly the same as the manual running speed. Tarry et al.
77 (2014) developed an integrated system for chrysanthemum bud detection by extracting color
78 characteristics. The system can be utilized to automate labor-intensive work in flower greenhouses
79 with a detection precision of 78.2%. Tete&Kamlu (2017) extracted the characteristics of
80 chrysanthemum leaves using two algorithms, i.e., threshold segmentation and K-means clustering,
81 and detected the leaves infected with pests and diseases. However, this study merely provided the
82 images of segmented chrysanthemums without definite qualitative results. Yuan et al. (2018)
83 extracted the characteristics of chrysanthemum petals through the random forest algorithm and
84 tested the chrysanthemums for variety. The experimental results showed that the detection precision
85 increased with the number of extracted characters, and AP was up to 95% at maximum. Yang et al.
86 (2018) segmented the White Chrysanthemums based on HSV color space, extracted target image
87 features from an irrelevant background, and developed a robot plucking system for Hangzhou White

88 Chrysanthemum, of which the precision for identification reached 85%, and the average plucking
89 time was 0.4s. However, the robustness for the extraction method based on color features is poor.
90 Considering the shortcomings, [Yang et al. \(2019\)](#) further extracted the color and texture features of
91 the images by the RGB value and the gray level co-occurrence matrix. The experimental results
92 have shown that the system can effectively detect the image of Hangzhou White Chrysanthemum
93 under different light conditions with average plucking time of 0.7s and plucking success rate of 90%.
94 This study added the chrysanthemum texture features based on the original study. The detection
95 precision can be effectively improved by the extraction and reasonable utilization of image features,
96 but the speed for inference is still to be further improved.

97 There are advantages in the detection precision for the chrysanthemum detection method based
98 on traditional machine learning. However, in a complicated unstructured environment, it is difficult
99 to detect the chrysanthemum images, and the actual speed for detection is too slow to meet the
100 requirements for large-scale operations. Deep learning, a subset of machine learning, enables
101 learning of hierarchical representations and the discovery of potentially complex patterns from large
102 datasets ([Sezer&Altan, 2021b](#)). It has shown impressive advancements on various problems in
103 natural language processing and computer vision, and the performance of deep convolutional neural
104 networks (CNNs) on image classification, segmentation and detection are of particular note
105 ([Altan&Karasu, 2020](#); [Sezer&Altan, 2021a](#); [Togacar et al., 2020](#)). Deep learning in the agriculture
106 domain is also a promising technique with growing popularity ([Gao et al., 2021](#)). [Liu et al. \(2019\)](#)
107 proposed a deep learning model based on VGG16 and ResNet50 networks for identifying big
108 chrysanthemums. Relevant experimental results have shown that the precision of the model may be
109 up to 78%, and the inference speed for detecting a chrysanthemum image can reach 10ms. [Liu et al.](#)

110 (2020) tested the flowering periods of 7 commercial chrysanthemum species with the DNN
111 algorithm, and the detection precision of the model reached 96%. Van Nam (2020) detected the
112 chrysanthemum image with noise with Faster R-CNN algorithm, by which the detection precision
113 of the model was up to 76%, and the inference speed for detecting a chrysanthemum image was
114 0.3s. The above studies have shown that deep learning approaches outperform the traditional
115 machine learning methods regarding to inference speed, with no advantage over detection accuracy.
116 Furthermore, for these studies, the chrysanthemum detection was carried out under controlled
117 environments, and there is still a gap to reach a reliable detection in fields. A CNN architecture was
118 designed to understand common features of these environments, and to take advantage of the
119 synergy between them. In many works, fusing features of different scales is an important approach
120 to improve detection performance. Low-level features have a high resolution and contain more
121 details, but due to few convolution operations, these features have less semantic information and
122 more noise. In contrast, high-level features have more semantic information, but have a low
123 resolution and poor ability to perceive details (Liu et al., 2021). Apparently, the key to improving
124 the performance of detection models is to efficiently fuse the features of different convolutional
125 layers.

126 With the significant improvements in computer performance and the rapid development of
127 deep learning, the advantages of feature fusion have become increasingly prominent. Feature fusion
128 algorithms can be classified into three categories: algorithms based on Bayesian decision theory,
129 algorithms based on sparse representation theory, and algorithms based on deep learning theory. In
130 object detection, algorithms based on deep learning theory are the mainstream — that is, multi-class
131 features from several neural networks are fused to obtain fused features, examples of which include

132 spatial pyramid pooling (SPP) structure, pyramid pooling module (PPM), and atrous spatial pyramid
133 pooling (ASPP). The information of a large number of small objects can be obtained by feature
134 fusion, thus improving the detection accuracy (Bae et al., 2020). However, extending the
135 architecture usually requires more computation. This decreases the model inference speed, resulting
136 in low detection efficiency.

137 The task of harvesting chrysanthemum at a specific maturity stage usually requires less
138 inference time on small devices, which poses a complex challenge for computer vision algorithms.
139 In response to this, a lightweight network based on feature fusion is proposed in this paper, which
140 has a potential to be deployed on an embedded GPU chip without performance sacrifice. The main
141 purpose of the network design is to enable the architecture to achieve more abundant gradient
142 combinations and simultaneously reduce the amount of computation. The gradient flow can
143 propagate through different network paths by dividing the feature map at the base layer into two
144 parts, and then merging them through the proposed cross-stage hierarchy. By switching series and
145 transition steps, a large correlation difference is produced in the propagated gradient information.
146 The contributions of this paper are as follows:

147 1. A fusion detection model was designed that alternates between cross-stage partial DenseNet
148 (CSPDenseNet) and cross-stage partial ResNeXt (CSPResNeXt) as the main network, and is
149 equipped with several combination modules. The model can obtain abundant gradient combinations
150 and effectively truncate the redundant gradient flow.

151 2. The impact of dataset sizes, different data enhancement methods, and complex unstructured
152 scenarios were quantified with the TC-YOLO model, and the superiority of the TC-YOLO model
153 was validated by comparing it with a number of state-of-the-art detection models.

154 3. A lightweight detection architecture was designed for detecting chrysanthemum at the early
155 flowering stage, which can adapt to complex unstructured environments (illumination variation,
156 occlusion, overlapping, etc.).

157 The organization of this paper is as follows. Section 2 describes the materials and methods,
158 especially the settings of the dataset and the design details of the proposed TC-YOLO network.
159 Section 3 introduces the experimental results of TC-YOLO. Section 5 discusses the contribution of
160 this paper, the limitations of the study and unresolved problems. We reflect on these problems and
161 indicate the future possible solutions. Finally, a brief summary and conclusions of this paper are
162 provided in Section 6.

163 **2. Materials and methods**

164 *2.1. Dataset*

165 The chrysanthemum datasets used in this paper were collected from chrysanthemum breeding
166 bases in Sheyang County in Jiangsu Province, Dongzhi County in Anhui Province, and Nanjing
167 Agricultural University in Jiangsu Province from October 2019 to October 2020. An Apple X phone
168 camera was used to capture images with a resolution of 1080×1920 . All images were taken under
169 natural light, including variations in background, illumination, growth stages, occlusion, and
170 overlapping.

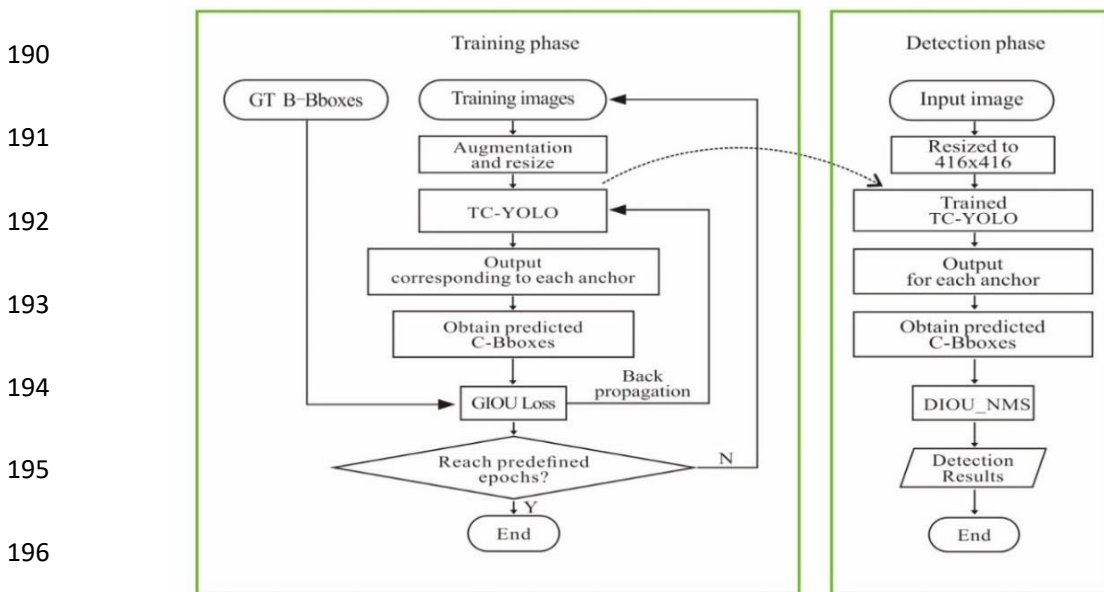
171 Chrysanthemums were captured in these images at three stages: the budding stage, the early
172 flowering stage, and the full-bloom stage. The budding stage refers to the stage when buds on the
173 main stem or branch of the chrysanthemum plant can be identified by the naked eye and the petals
174 are not opened. The early flowering stage refers to the stage when the petals are not fully opened,
175 and the full bloom stage refers to the stage when the petals are fully opened. A total of 1000

176 chrysanthemum image samples were randomly divided into training (60%), validation (30%), and
 177 test datasets (10%). Each prediction module had three prior anchors of different scales. Through k -
 178 means clustering [39], nine prior anchors — (10,13), (16,30), (33,23), (30,61), (62,45), (59,119),
 179 (116,90), (156,198), and (373,326) — of the chrysanthemum dataset were obtained. The pre-
 180 processed chrysanthemum image $416 \times 416 \times 3$ was transformed into a $208 \times 208 \times 12$ feature map
 181 via the focus slicing operation, and then one convolution operation of 32 convolution kernels was
 182 performed to finally form a $208 \times 208 \times 32$ feature map. The examples of preprocessed
 183 chrysanthemum images are shown in Fig. 2.



188 Fig. 2. Preprocessed chrysanthemum images. The preprocessed methods include flip and rotation.

189 2.2. TC-YOLO



197 Fig. 3. The flow chart of TC-YOLO training and detection. In the training phase, the original chrysanthemum image

198 is enhanced and the image size is adjusted to 416*416. The coordinates of the predicted boxes are obtained after the
199 image passes through the proposed TC-YOLO framework. The GIOU Loss function and back-propagation method
200 are used to train the network. The training stops when the preset epoch is reached. In the detection phase, the
201 DIOU_NMS is used to evaluate the distance between the predicted box and the ground true box, so as to output
202 accurate detection results.

203 TC-YOLO is a lightweight CNN (13.6M) that can adapt to complex unstructured environments
204 (illumination variation, occlusion, overlap, etc.). Fig. 3 shows the flow chart of TC-YOLO training
205 and detection. The network architecture is deep both vertically and horizontally, that is, it has both
206 top-down connection and horizontal connection, as shown in Fig. 4. TC-YOLO has three main
207 components: the Backbone, the Neck and the Head components. For the Backbone component, the
208 main network is CSPDenseNet, and the CBL addon module and the SPP addon module are added
209 to the component. For the Neck component, the feature extraction network is CSPResNeXt, and the
210 CBL addon module is added to the component. For the Head component, the multi-scale fusion
211 network is the Recursive Feature Pyramid (RFP)+ Atrous Spatial Pyramid Pooling (ASPP) structure.
212 The detection component is a post-processing step, involving iterative region proposal and non-
213 maximum suppression DIOU_NMS based on anchor frame. Mish activation function, drop block
214 regularization method and GIOU loss function are used in the whole network.

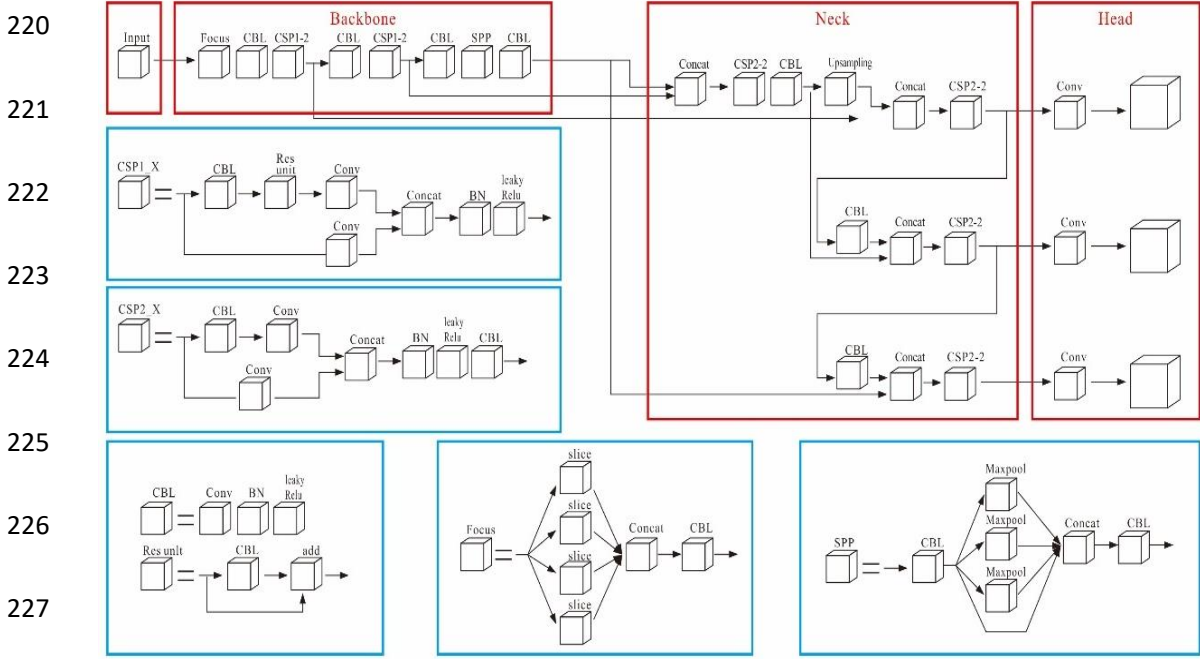
215

216

217

218

219



228 **Fig. 4.** Structure of the proposed TC-YOLO network. Convolution + Batch normalization + Leaky ReLU (CBL),
 229 Cross-Stage Partial (CSP), Spatial Pyramid Pooling (SPP), Batch Normalization (BN).

230 *2.2.1. Backbone*

231 *2.2.1.1. CSPDenseNet*

232 The architecture of one stage of the proposed CSPDenseNet network is shown in Fig. 4. The
 233 stage of the CSPDenseNet network is composed of a partial dense block and a partial transition
 234 layer. In the partial dense block, the feature maps of the base layer in the stage are split into two
 235 parts through the channel: $x_0 = [x_0', x_0'']$. Between x_0' and x_0'' , the former is directly linked to
 236 the end of the stage, and the latter will go through a dense block. All steps involved in the partial
 237 transition layer are as follows. First, the output of dense layers, $[x_0', \dots, x_k]$, will go through a
 238 transition layer. Second, the output of this transition layer, x_T , will be concatenated with x_0'' and
 239 go through another transition layer, and then generate an output x_U . The feed-forward pass and
 240 weight updating of CSPDenseNet are shown in Eq. (1) and Eq. (2), respectively.

$$x_k = w_k * [x_0, x_1, \dots, x_{k-1}] \quad (1)$$

$$\begin{aligned}
x_T &= w_T * [x_0, x_1, \dots, x_k] \\
x_U &= w_U * [x_0, x_T] \\
w_k' &= f(w_k, g_0, g_1, g_2, \dots, g_{k-1}) \\
w_T' &= f(w_T, g_0, g_1, g_2, \dots, g_k) \\
w_U' &= f(w_U, g_0, g_T)
\end{aligned} \tag{2}$$

241 where, $*$ is the convolution operator; $[x_0, x_1, \dots]$ is the connection x_0, x_1, \dots, w_i ; x_i represents the
242 weight and output of the i -th dense layer; f is the function (LeakyReLU) of weight updating; and
243 g_i is the gradient propagating to the i -th dense layer. It was found that a large amount of gradient
244 information is repeatedly used to update the weights of different dense layers, due to which different
245 dense layers may repeatedly learn the copied gradient information.

246 From Eq. (1) and Eq. (2), we can see that, gradients from dense layers are integrated separately
247 and the feature mapping x_0' , which does not pass through the dense layer is also integrated
248 separately. For gradient information updated by weights, both sides do not contain repeated gradient
249 information of the other sides. The CSPDenseNet network proposed in this section retains the
250 advantages of DenseNet feature reuse characteristics, but at the same time prevents excessive
251 repeated gradient information by truncating the gradient flow. This idea was featured by designing
252 a hierarchical feature fusion strategy and applying it to partial transition layers.

253 2.2.1.2. Partial Dense Block

254 The purpose of the partial dense block design is to (1) increase gradient paths — the number
255 of gradient paths can be doubled by splitting and merging, and thanks to the cross-stage strategy,
256 the shortcomings of using explicit feature map copy for splicing can be alleviated; (2) balance the
257 computation of each layer — generally, the number of channels at the base layer of DenseNet is

258 much greater than the growth rate, and since the number of channels at the base layer involved in
 259 dense layer operations in some partial dense blocks only accounts for half of the original number,
 260 nearly half of the computational bottleneck can be effectively removed; and (3) reduce the memory
 261 flow — assuming that the size of the basic feature map of a dense block in the DenseNet is $w * h * c$,
 262 the growth rate is d , and there are m dense layers in total. Then, the CIO of the dense module is
 263 calculated as follows:

$$264 \quad (c * m) \frac{(m^2 + m) * d}{2} \quad (3)$$

265 The CIO of partial dense blocks is calculated as follows:

$$266 \quad \frac{(c * m) + (m^2 + m) * d}{2} \quad (4)$$

267 Although m and d are usually much smaller than c , the a partial dense module can still save up
 268 to half of the memory flow in the network.

269 2.2.1.3. CBL + SPP

270 CBL is the smallest component in the whole network structure, which is closely connected by
 271 one convolutional layer, batch normalization (BN), and Leaky ReLU activation functions. In order
 272 to maximize the difference in gradient flow passing through CBL, an SPP component is added to
 273 the top layer of the CBL component, and kernel sizes and strides of varied sizes are used to output
 274 different receptive field features; then, a concatenation operation is performed. Specifically, the
 275 equations of the feature mapping matrix of the CBL + SPP component are as follows:

$$K_h = \text{ceil} \left(\frac{h_{in}}{n} \right) \quad (5)$$

$$S_h = \text{ceil} \left(\frac{h_{in}}{n} \right)$$

$$P_h = \text{floor} \left(\frac{K_h * n - h_{in} + 1}{2} \right)$$

$$H = 2 * P_h + h_{in}$$

$$K_w = \text{ceil} \left(\frac{w_{in}}{n} \right) \quad (6)$$

$$S_w = \text{ceil} \left(\frac{w_{in}}{n} \right)$$

$$P_w = \text{floor} \left(\frac{k_w * n - w_{in} + 1}{2} \right)$$

$$w = 2 * P_w + w$$

276 where K_h , S_h , P_h , and h are the height of the kernel, the step size in the height direction of the
 277 feature mapping matrix, the number of fillings in the height direction of the feature mapping matrix,
 278 and the height of the feature mapping matrix, respectively; $\text{ceil}()$ is the rounding up symbol; floor
 279 $()$ is the rounding down symbol; and h_{in} is the height of the input data.

280 By combining Eq. (5) and Eq. (6), the equation of the feature mapping matrix is obtained, as
 281 follows:

$$\left[\frac{h + 2p - f}{s} + 1 \right] * \left[\frac{w + 2p - f}{s} + 1 \right] \quad (7)$$

283 where p stands for padding, s stands for stride, and f is the input data size.

284 2.2.2. Neck

285 The neck component uses the same structure as the backbone. It should be noted that the
 286 adopted CSPResNeXt structure can greatly promote the extraction of low-level chrysanthemum
 287 features and high fusion at different feature scales. However, when redundant gradient flow enters
 288 the neck, in order to truncate the redundant gradient flow and avoid excessive GPU calculation and
 289 over-fitting, the main network of the neck component is changed from CSPDenseNet to
 290 CSPResNeXt. The CSPResNeXt structure is composed of five groups of blocks of different sizes.
 291 The number of ResNeXt networks of each group of blocks is designed to be 1. ResNeXt, in essence,
 292 is a group convolution, which controls the number of groups by variable cardinalities. Group
 293 convolution is a compromise between normal convolution and depth-wise separable convolution —

294 that is, the number of channels of the feature map generated by each branch is n ($n > 1$). The equation
 295 is as follows:

$$296 \quad y = x + \sum_{i=1}^C T_i(x) \quad (8)$$

297 where $+$ is a shortcut; C is the variable cardinality of the group convolution; T_i is a series of
 298 convolutions; and T is composed of continuous convolutions ($1*1 \rightarrow 3*3 \rightarrow 1*1$).

299 2.2.3. Head

300 The Head component is the prediction part of the network, and the scale of the final predicted
 301 feature map is 76×76 , 38×38 , and 19×19 . The structure of the prediction part is Recursive Feature
 302 Pyramid (RFP). RFP adds feedback connections to FPN as highlighted in Fig. 5. R_i denotes the
 303 feature transformations before connecting them back to the bottom-up backbone. Then, $\forall i = 1, \dots$,
 304 S , the output feature f_i of RFP is defined by

$$305 \quad f_i = F_i(f_{i+1}, X_i), X_i = B_i(X_{i-1}, R_i(f_i)) \quad (9)$$

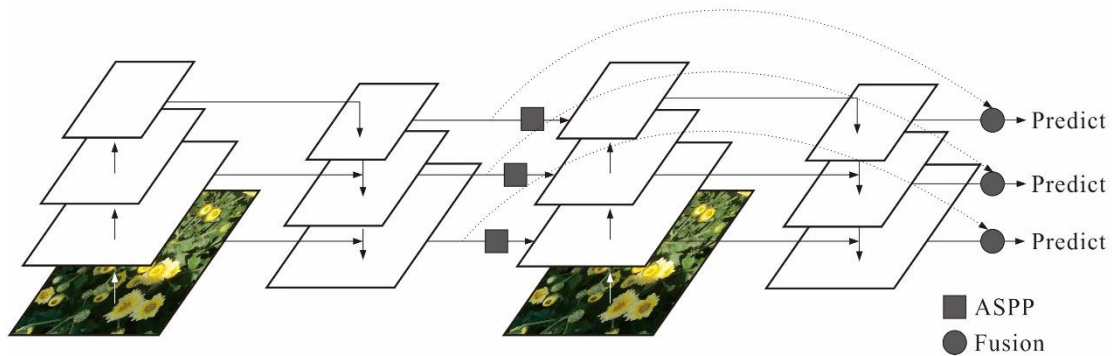
306 Which makes RFP a recursive operation. We unroll it to a sequential network, *i.e.*, $\forall i = 1, \dots, S, t =$
 307 $1, \dots, T$,

$$308 \quad f_i^t = F_i^t(f_{i+1}^t, X_i^t), X_i^t = B_i^t(X_{i+1}^t, R_i^t(f_i^{t-1})) \quad (10)$$

309 where T is the number of unrolled iterations, and we use superscript t to denote operations and
 310 features at the unrolled step t . f_i^0 is set to 0. In our implementation, F_i^t and R_i^t are shared across
 311 different steps.

312 We used Atrous Spatial Pyramid Pooling (ASPP) to implement the connecting module R,
 313 which takes a feature f_i^t as its input and transforms it to the RFP feature used in Fig. 5. In this
 314 module, there are four parallel branches that take f_i^t as their inputs, the outputs of which are then
 315 concatenated together along the channel dimension to form the final output of R. Three branches of

316 them use a convolutional layer followed by a ReLU layer, the number of the output channels is 1/4
 317 the number of the input channels. The last branch uses a global average pooling layer to compress
 318 the feature, followed by a 1x1 convolutional layer and a ReLU layer to transform the compressed
 319 feature to a 1/4-size (channel-wise) feature. Finally, it is resized and concatenated with the features
 320 from the other three branches. The convolutional layers in those three branches are of the following
 321 configurations: kernel size = [1, 3, 3], atrous rate = [1, 3, 6], padding = [0, 3, 6]. Unlike the original
 322 ASPP, we do not have a convolutional layer following the concatenated features as in here R does
 323 not generate the final output used in dense prediction tasks. Note that each of the four branches
 324 yields a feature with channels 1/4 that of the input feature, and concatenating them generates a
 325 feature that has the same size as the input feature of R.



330 **Fig. 5.** The architecture of Recursive Feature Pyramid (RFP).

331 2.3. Evaluation metrics

332 For object detection applications, average precision (AP) is a standard metric for evaluation of
 333 model performance. In our case, we calculated the average accuracy for one category, tea
 334 chrysanthemums, to check the overall performance of the model. Precision is a ratio of true object
 335 detections to the total number of objects that a model predicted. Recall is a ratio of true object
 336 detections to the total number of objects in the dataset. The equation is as follows:

$$337 \quad AP = \sum_{k=1}^N P(k) \Delta recall(k) \quad (11)$$

338 where N is the total number of images in the test dataset, $P(k)$ is the precision value at k images and
339 *recall* (k) is the change of the recall between k and $k - 1$ images.

340 2.4. Experimental Setup

341 Four experiments were designed to test the performance of the proposed detection model. In
342 the first experiment (Section 3.1), we randomly selected and established eight datasets in varied
343 sizes to investigate the influence of dataset size for TC-YOLO modeling. In the second experiment
344 (Section 3.2), in order to understand the impact of data augmentation on TC-YOLO, based on the
345 results of the first experiment, we selected 14 data augmentation methods and performed the test
346 using different configuration schemes. In the third experiment (Section 3.3), in order to investigate
347 the robustness of TC-YOLO in complex unstructured scenarios, we set up nine unstructured
348 scenarios based on our previous chrysanthemum harvesting experience. In the fourth experiment
349 (Section 3.4), in order to thoroughly study the detection performance of TC-YOLO, we tested nine
350 latest target detection frameworks on the chrysanthemum dataset.

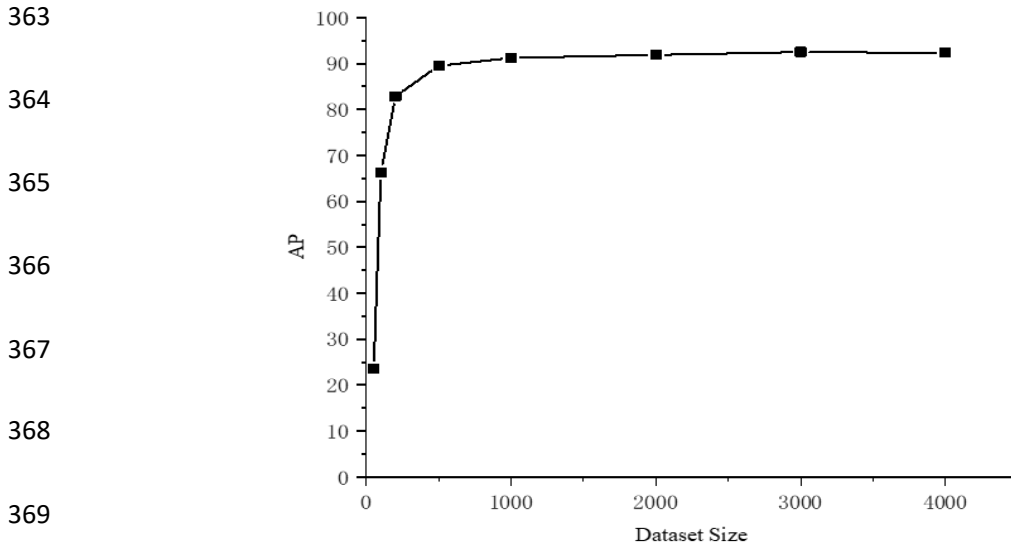
351 The experiments were carried out on a server with NVIDIA Tesla P100, CUDA 11.2. The
352 operating system was Ubuntu 16. The basic detection frameworks were CSPDenseNet and
353 CSPResNeXt. In the training process, the key hyperparameter settings were as follows: epoch = 300;
354 learning rate = 0.01; momentum = 0.937; gamma = 0.1; weight decay = 5×10^{-4} ;
355 warmup_momentum = 0.8; warmup_bias_lr = 0.1. The optimizer used was a stochastic gradient
356 descent (SGD).

357 3. Results

358 3.1. Impact of Dataset Size on the TC-YOLO

359 To verify the impact of the dataset size on the chrysanthemum detection task, eight datasets of

360 varied sizes were randomly selected from the chrysanthemum training set, which contain 50, 100,
361 200, 500, 1000, 2000, 3000, 4000 images, respectively. The AP of the chrysanthemum is shown in
362 [Fig. 6](#).



370 **Fig. 6.** Impact of dataset size on the detection task. Eight different dataset sizes were adopted to validate the optimal
371 dataset size in this paper.

372 It can be seen that the performance of the detection model improves with the increase in dataset
373 size in [Fig. 6](#). When the number of images was less than 500, the AP value increased rapidly (23.62%
374 to 66.29%, improved by 180.65%) with the increase in the number of images. When the number of
375 samples reached at 1000, only small improvements can be obtained (from 91.22% to 92.49%) by
376 adding more samples. The performance of the detection model converged with AP value at 93%
377 after 3000 images. The optimal dataset size was set as 3,000 images for the experiment in the
378 following sections due to its high AP value achieved under the minimal sample number required.

379 3.2. The impact of data augmentation on the TC-YOLO

380 To explore the influence of data augmentation on TC-YOLO, we selected 14 data augmentation
381 methods for the test, and the test results are shown in [Table 1](#). Our experimental design is to

382 configure these data augmentation methods in turn, and select the best-performing data
 383 augmentation methods according to the performance of corresponding models. Surprisingly, the
 384 most advanced data enhancement methods, such as Cutout, Mixup, Cutmix and Mosaic, only
 385 achieved performance APs of 86.11%, 85.16%, 86.51% and 88.62% respectively. This may be
 386 because massive redundant gradient information will significantly reduce the learning ability of the
 387 network. We found that Flip and Rotation performed best, with 91.02% and 91.83% AP values,
 388 respectively. Moreover, when both Flip and Rotation were configured, the performance of the model
 389 was slightly improved (AP=92.49%). Therefore, this combination was adopted for data
 390 augmentation. Note that other than Flip and Rotation, the improvement of model performance
 391 gained by using other data augmentation methods was rather limited.

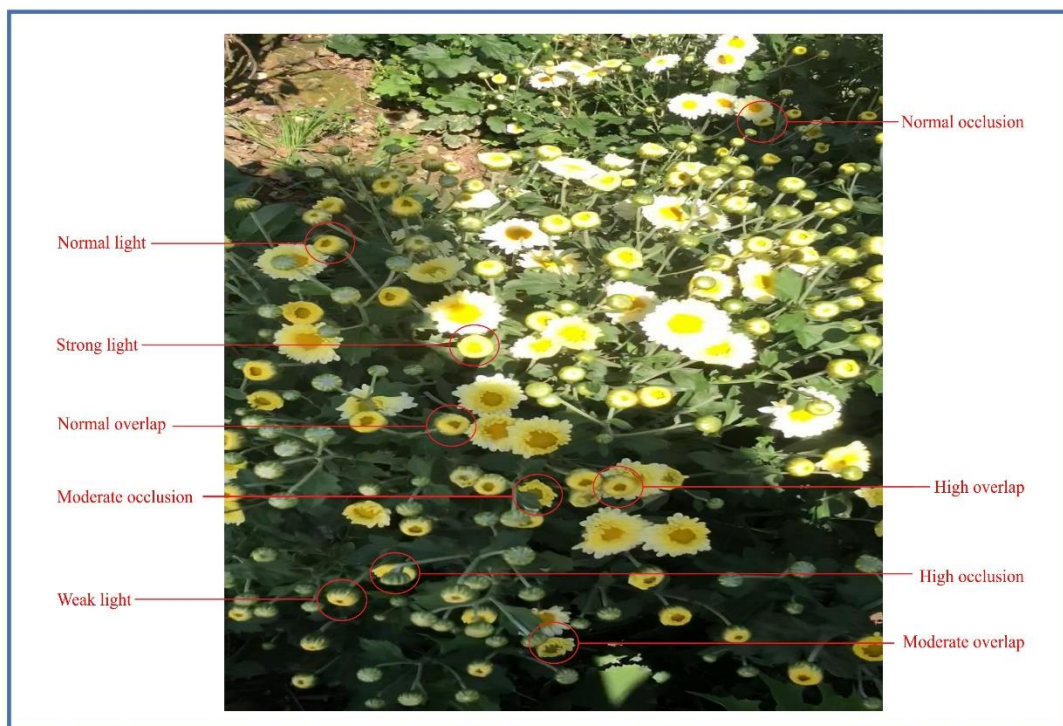
392 **Table 1**

393 The impact of data augmentation on the TC-YOLO. Fourteen different data enhancement methods are employed to
 394 identify the best data enhancement strategy.

| Flip | Shear | Crop | Rotation | Grayscale | Hue | Saturation | Exposure | Blur | Noise | Cutout | Mixup | Cutmix | Mosaic | AP |
|------|-------|------|----------|-----------|-----|------------|----------|------|-------|--------|-------|--------|--------|-------|
| √ | | | | | | | | | | | | | | 91.02 |
| | √ | | | | | | | | | | | | | 87.96 |
| | | √ | | | | | | | | | | | | 87.63 |
| | | | √ | | | | | | | | | | | 91.83 |
| | | | | √ | | | | | | | | | | 85.11 |
| | | | | | √ | | | | | | | | | 88.67 |
| | | | | | | √ | | | | | | | | 88.23 |
| | | | | | | | √ | | | | | | | 84.52 |
| | | | | | | | | √ | | | | | | 83.26 |
| | | | | | | | | | √ | | | | | 88.89 |
| | | | | | | | | | | √ | | | | 86.11 |
| | | | | | | | | | | | √ | | | 85.16 |
| | | | | | | | | | | | | √ | | 86.51 |
| | | | | | | | | | | | | | √ | 88.62 |
| √ | | | √ | | | | | | | | | | | 92.49 |

395 *3.3. Impact of Different Unstructured Scenarios on the TC-YOLO*

396 This study examined the robustness of the proposed model in different unstructured
397 environments, including strong light, weak light, normal light, high overlap, moderate overlap,
398 normal overlap, high occlusion, moderate occlusion, and normal occlusion. There were 30236
399 chrysanthemums in the early flowering stage in nine unstructured environments. Note that there is
400 no strict criterion for reference as defining various scenarios. Some criteria are set artificially for
401 the convenience of testing. First of all, most scenarios may be differentiated empirically, as shown
402 in Fig. 7, which has provided nine scenarios. Additionally, according to our definition, tea
403 chrysanthemums are counted separately in different scenarios, e.g., for a tea chrysanthemum in
404 normal light, normal overlap, and normal occlusion scenarios. In counting, the number of tea
405 chrysanthemums is increased by 1 respectively in the three scenarios.



415 Fig. 7. Example of 9 unstructured scenarios.

416 Table 2 shows that, under normal environmental conditions, the AP of chrysanthemums at the
417 early flowering stage reached impressive values of 96.52%, 97.11%, and 96.88% for normal light,

418 normal overlap and normal occlusion, respectively. When the unstructured environment became
 419 complex, the AP of chrysanthemums at the early flowering stage decreased significantly, especially
 420 under the scenario of strong light, where the AP was only 80.54%. Interestingly, in all unstructured
 421 environments, the error rate under strong light was the highest, reaching 7.54%, while the miss rate
 422 under high overlap was the highest, reaching 14.22%. In addition, in general, overlapping had the
 423 least influence on the detection of chrysanthemums at the early flowering stage. Under high overlap,
 424 the AP, error rate, and miss rate were 82.51%, 3.27%, and 14.22%, respectively. Illumination had
 425 the biggest impact on the detection of chrysanthemums at the early flowering stage. With strong
 426 light, the accuracy, error rate and miss rate were 80.54%, 7.54%, and 11.92%, respectively.

427 **Table 2**

428 Impact of Different Unstructured Scenarios on the TC-YOLO. Nine different scenarios were artificially set up to
 429 verify the generalisability of the proposed model in different scenarios.

| Environment | Count | Correctly Identified | | Falsely Identified | | Missed | |
|--------------------|-------|----------------------|----------|--------------------|----------|--------|----------|
| | | Amount | Rate (%) | Amount | Rate (%) | Amount | Rate (%) |
| Strong light | 6853 | 5519 | 80.54 | 517 | 7.54 | 817 | 11.92 |
| Weak light | 16511 | 14639 | 88.66 | 1174 | 7.11 | 698 | 4.23 |
| Normal light | 23542 | 22723 | 96.52 | 438 | 1.86 | 381 | 1.62 |
| High overlap | 5699 | 4702 | 82.51 | 186 | 3.27 | 811 | 14.22 |
| Moderate overlap | 12998 | 11814 | 90.89 | 476 | 3.66 | 708 | 5.56 |
| Normal overlap | 21883 | 21251 | 97.11 | 182 | 0.83 | 450 | 2.06 |
| High occlusion | 8352 | 6806 | 81.49 | 601 | 7.2 | 945 | 11.31 |
| Moderate occlusion | 14552 | 13076 | 89.86 | 586 | 4.03 | 890 | 6.11 |

| | | | | | | | |
|------------------|-------|-------|-------|-----|------|-----|------|
| Normal occlusion | 23523 | 22789 | 96.88 | 296 | 1.26 | 438 | 1.86 |
|------------------|-------|-------|-------|-----|------|-----|------|

430 To intuitively show the effect of illumination variation on the performance of TC-YOLO, we
431 collected two images of the same tea chrysanthemum from different angles for testing. The image
432 on the left was collected under weak light, while that on the right was collected under Normal light.

433 The results are shown in Fig. 8.



440 Fig. 8. The detection performance of the same tea chrysanthemum under different light intensities. The undetected
441 red box on the left image is manually marked. It can be clearly seen that light intensity plays a major role on the
442 performance of TC-YOLO.

443 3.4. Comparisons with State-of-the-Art Detection Methods

444 To comprehensively verify the performance of the proposed method, We have compared the
445 proposed TC-YOLO with the latest object detection technologies based on feature fusion, including
446 FPN-based RetinaNet models; EfficientDet (Cao et al., 2020) models, based on EfficientB + BiFPN;
447 multi-level feature pyramid network (MLFPN)-based M2Det (Zhang&Li, 2020) models; YOLOv3
448 models based on DarkNet53 + FPN; modified sine-cosine algorithm (MSCA)-based VGG16 +
449 PFPNet (Wang et al., 2019) models; RefineDet (Zhang et al., 2018) models, based on Advanced
450 RISC Machine (ARM) + object detection module (ODM); YOLOv4 models based on

451 CSPDarknet53; and YOLOv5 models based on CSPDenseNet.

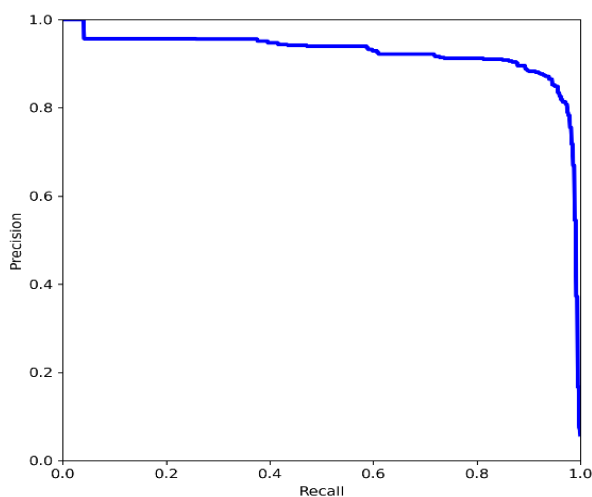
452 **Table 3**

453 Comparisons with state-of-the-art detection methods. We used the same dataset and the same computing device to

454 compare the proposed model with thirty-four different models for validating the superiority of the proposed model.

| Method | Backbone | Size | FPS | AP |
|----------------|--------------|-------------|-------|-------|
| RetinaNet | ResNet101 | 800 × 800 | 4.53 | 71.51 |
| RetinaNet | ResNet50 | 800 × 800 | 5.33 | 78.66 |
| RetinaNet | ResNet101 | 500 × 500 | 7.26 | 83.82 |
| RetinaNet | ResNet50 | 500 × 500 | 7.88 | 83.26 |
| EfficientDetD6 | EfficientB6 | 1280 × 1280 | 5.23 | 85.67 |
| EfficientDetD5 | EfficientB5 | 1280 × 1280 | 6.22 | 85.54 |
| EfficientDetD4 | EfficientB4 | 1024 × 1024 | 8.02 | 85.06 |
| EfficientDetD3 | EfficientB3 | 896 × 896 | 9.16 | 86.21 |
| EfficientDetD2 | EfficientB2 | 768 × 768 | 11.63 | 86.45 |
| EfficientDetD1 | EfficientB1 | 640 × 640 | 15.29 | 82.09 |
| EfficientDetD0 | EfficientB0 | 512 × 512 | 37.63 | 80.53 |
| M2Det | VGG16 | 800 × 800 | 7.03 | 80.67 |
| M2Det | ResNet101 | 320 × 320 | 16.88 | 75.88 |
| M2Det | VGG16 | 512 × 512 | 21.28 | 73.29 |
| M2Det | VGG16 | 300 × 300 | 42.56 | 70.22 |
| YOLOv3 | DarkNet53 | 608 × 608 | 12.21 | 87.44 |
| YOLOv3(SPP) | DarkNet53 | 608 × 608 | 15.67 | 84.62 |
| YOLOv3 | DarkNet53 | 416 × 416 | 43.23 | 81.23 |
| YOLOv3 | DarkNet53 | 320 × 320 | 46.92 | 75.52 |
| PFPNet (R) | VGG16 | 512 × 512 | 24.43 | 74.98 |
| PFPNet (R) | VGG16 | 320 × 320 | 33.42 | 80.33 |
| PFPNet (s) | VGG16 | 300 × 300 | 41.33 | 81.58 |
| RFBNetE | VGG16 | 512 × 512 | 21.46 | 78.83 |
| RFBNet | VGG16 | 512 × 512 | 36.11 | 78.36 |
| RFBNet | VGG16 | 512 × 512 | 45.49 | 83.58 |
| RefineDet | VGG16 | 512 × 512 | 31.24 | 78.66 |
| RefineDet | VGG16 | 448 × 448 | 43.09 | 76.85 |
| YOLOv4 | CSPDarknet53 | 608 × 608 | 19.28 | 87.62 |
| YOLOv4 | CSPDarknet53 | 512 × 512 | 24.69 | 87.43 |
| YOLOv4 | CSPDarknet53 | 300 × 300 | 46.61 | 84.92 |
| YOLOv5s | CSPDenseNet | 416 × 416 | 47.62 | 89.86 |
| YOLOv5l | CSPDenseNet | 416 × 416 | 42.22 | 88.96 |
| YOLOv5m | CSPDenseNet | 416 × 416 | 36.89 | 87.23 |
| YOLOv5x | CSPDenseNet | 416 × 416 | 32.23 | 85.26 |
| Ours | CSPD+R | 416 × 416 | 47.23 | 92.49 |

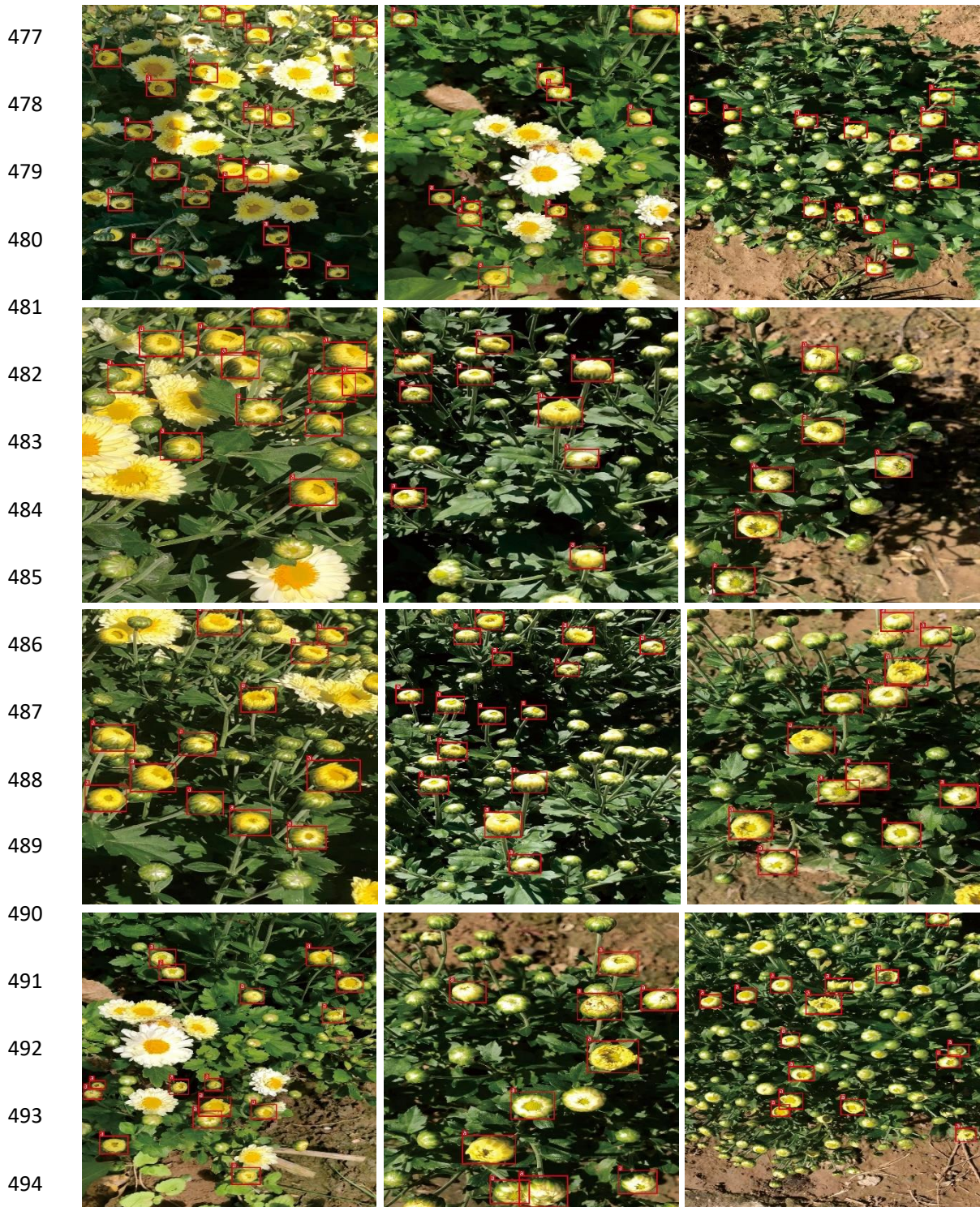
455 First, the experimental results in Table 3 were analyzed from the perspective of AP. The
456 proposed method achieved the highest AP of 92.49%, which is 2.63% higher than the best-
457 performing feature fusion-based region detector YOLOv5s. There may be two reasons for this. One
458 is that the proposed method based on some fusion components and CSPDenseNet+CSResNeXt
459 alternate fusion can better guide the gradient flow of tea chrysanthemum features, and the other
460 reason is that in the complex unstructured chrysanthemum dataset, adding RFP structure can
461 increase the gradient backflow in model training, which may activate more detailed features in the
462 chrysanthemum image. It is worth mentioning that the detection speed of the proposed method is
463 not the fastest; based on the 416x416 input scale, YOLOv5s (the fastest) is 0.39 FPS quicker than
464 our approach. This may be because the addition of the gradient backflow mechanism slightly
465 increases the reasoning scale of the model. However, compared with the significant improvement
466 of AP, it is totally acceptable to lose some reasoning speed. Fig. 9 shows the PR curves during the
467 training process. Fig. 10 displays the qualitative results of TC-YOLO on the chrysanthemum dataset.



474 **Fig. 9.** The PR curves during the training process.

475

476

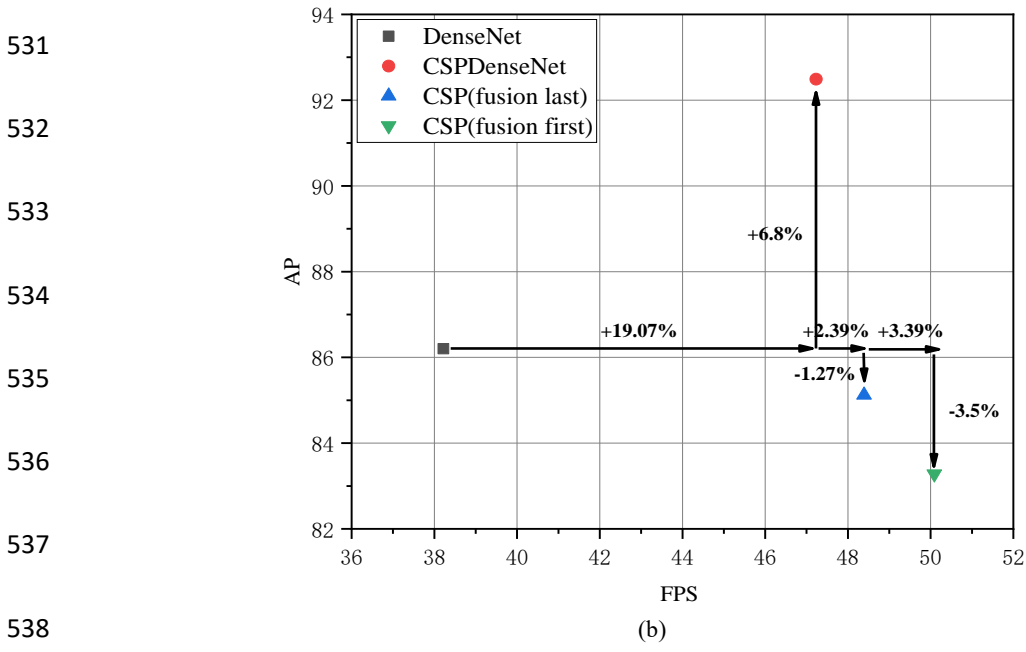
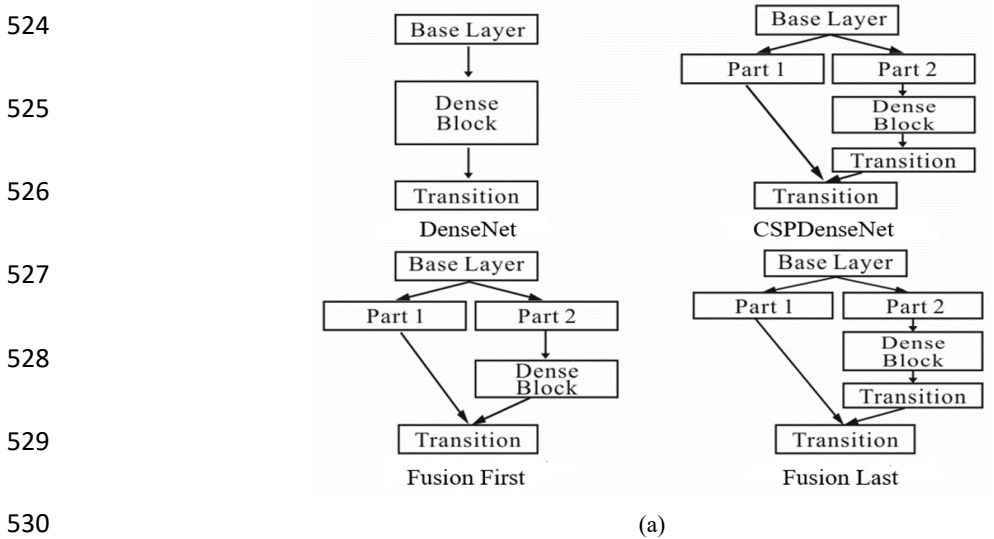


495 **Fig. 10.** Qualitative results of our method. As can be seen from the figure, the proposed method can satisfactorily
496 detect chrysanthemums in the early flowering stage in complex unstructured scenarios including illumination
497 variation, occlusion and overlap in real time.

498 **4. Discussion**

499 In this paper, a lightweight object detection model called tea chrysanthemum - YOLO (TC-
500 YOLO) is proposed, which can adapt to illumination variation, occlusion and overlapping scenarios.
501 In order to realize the detection task, a separate fusion CNN was constructed to fuse different
502 function modules, the purpose of which is to transform chrysanthemum features into common
503 subspace, in which features can be combined linearly or nonlinearly. Moreover, data augmentation
504 methods were applied, and customized specific loss function was used to train these modules. In
505 this way, these function modules can better understand the feature of chrysanthemums at different
506 scales, thus improving the end-to-end performance of the lightweight network in complex
507 unstructured environment. Local transition layer CSPDenseNet and CSPResNeXt were added into
508 different function modules as the main network. The main purpose of the design is to enrich the
509 architecture with gradient combinations and to reduce the amount of calculation. By dividing the
510 feature map of the base layer into two parts and combining them through the proposed cross-stage
511 hierarchy, gradient flow can travel through different network paths. Through the transformation of
512 series and step size, a large correlation difference is generated in the gradient information. We also
513 visualized the features learned by the tea chrysanthemum – YOLO. Also, to demonstrate the role of
514 the local transition layer, Fig. 11 (a) shows four different fusion strategies. CSP (fusion first) refers
515 to connecting the feature maps generated by two parts and then performing the conversion. If this
516 strategy is adopted, a lot of gradient information will be reused. As for the CSP (fusion last) strategy,
517 the output of dense blocks will pass through the transition layer, and then be connected to the feature
518 map from the first part. If the CSP (fusion last) strategy is adopted, gradient information will not be
519 reused, because the gradient flow is truncated. With DenseNet as the baseline, the results are shown
520 in Fig. 11 (b). In the process of information integration, the cross-stage splitting and merging

521 strategy can effectively reduce the possibility of repetition. It can be seen from the results in Fig. 11
 522 that the learning ability of the network will be significantly improved if repeated gradient
 523 information is effectively reduced.



539 **Fig. 11.** Four architectures as shown in Fig. 11 (a) were used for image detection, with DenseNet as the baseline,
 540 and the detection results were as shown in Fig. 11 (b). The CSPDenseNet strategy for image detection reduced the
 541 computational cost by 19.07%, and increased AP by 6.8%, while the CSP (Fusion Last) strategy reduced the
 542 computational cost by 21.46% though decreased AP by 1.27%. Similarly, the CSP (Fusion First) reduced the

543 computational cost by 24.85% though decreased AP by 3.5%.

544 Also, the features learned from TC-YOLO are visualized. Fig. 12 shows the network activation
545 of the four input images in the third, tenth, and twentieth convolutional layers. The red, yellow, and
546 blue represent the activated part, the potentially activated part, and the non-activated part
547 respectively. Generally, the regions are different in terms of the degree of activation at different
548 layers for the feature map. Taking the first row in Fig. 12 for instance, starting from the left, the first
549 image is a preprocessed input image. The second image shows the activation of the input image in
550 the third convolutional layer. The upper left region of the feature map has been activated to a high
551 degree, but the red areas are sparse, indicating a low degree of activation. The third image shows
552 the activation of the input image in the tenth convolutional layer. The convergence of red areas is
553 more significant than that in the second image, suggesting that the input image has been activated
554 deeply at the tenth convolutional layer. The fourth image shows the activation of the input image in
555 the twentieth convolutional layer. Obviously, peripheral areas that have not been activated are
556 noticed by the feature fusion mechanism, and features of a large number of peripheral areas are
557 activated.

558

559

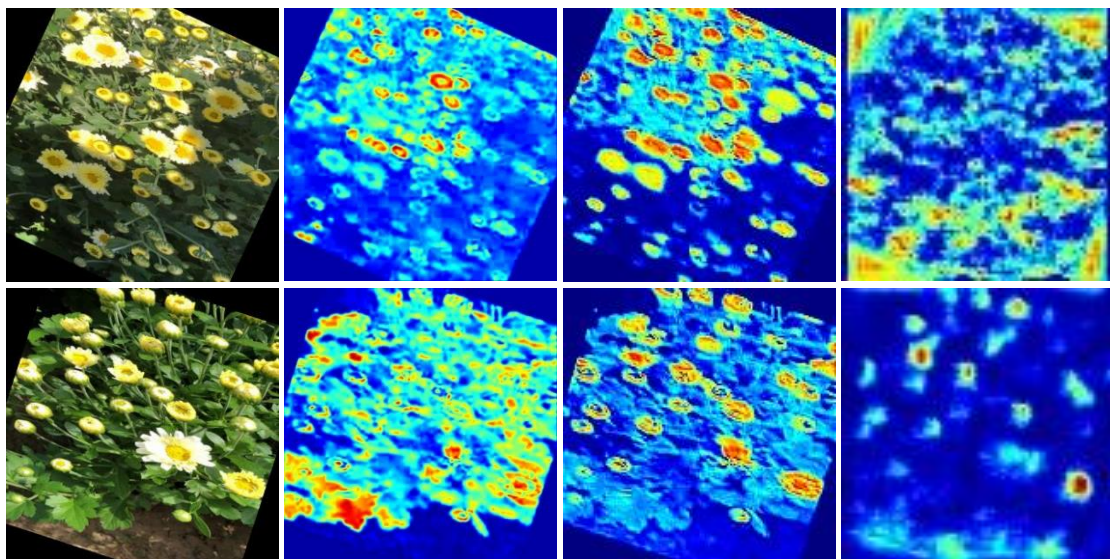
560

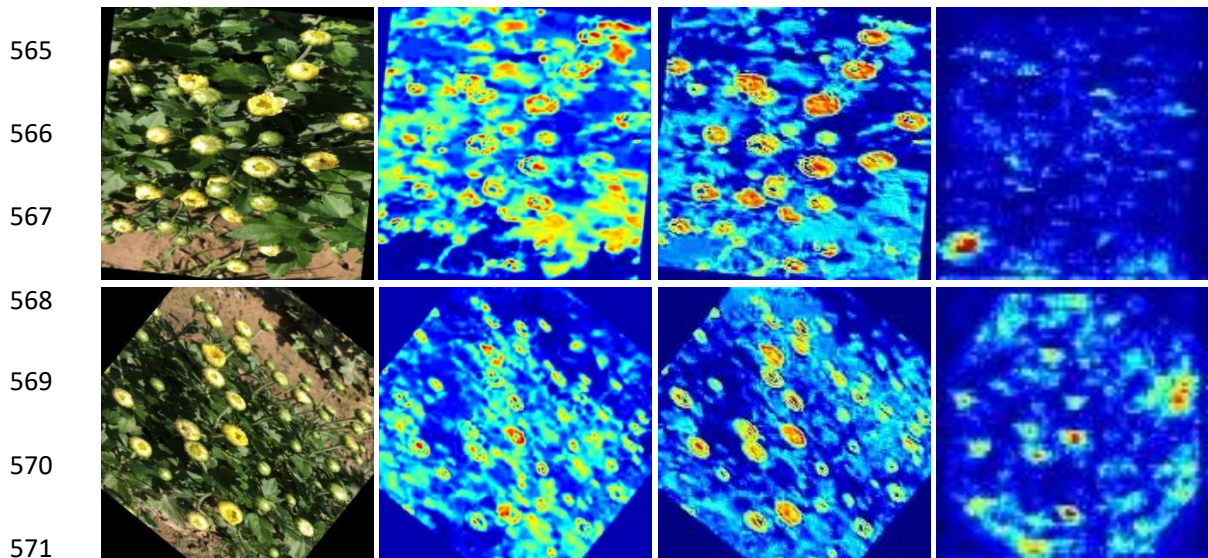
561

562

563

564





572 **Fig. 12.** Visualization results of some input images. The first, second, third and fourth columns are for the
 573 preprocessed input image, the input image activated at the third convolutional layer, the input image activated at the
 574 tenth convolutional layer, and the input image activated at the twentieth convolutional layer, respectively.

575 Research on machine learning-based chrysanthemum detection is currently still not fully
 576 explored. Traditional machine learning-based chrysanthemum detection has non-negligible
 577 drawbacks in terms of reliability in fields and inference speed (Kondo et al., 1996; Tarry et al., 2014;
 578 Tete&Kamlu, 2017; Warren, 2000; Yang et al., 2018; Yang et al., 2019; Yuan et al., 2018). Among
 579 them, the fastest inference speed is the robotic chrysanthemum picking system for Hangzhou white
 580 chrysanthemums developed by Yang et al. (2018) in 2018, recognizing an image with an inference
 581 speed of 0.4 s (2.5 FPS). Our TC-YOLO detection of an image can reach an inference speed of
 582 approximately 21.17 ms (47.23 FPS), which is about 19 times faster than the proposed method in
 583 Yang et al. (2018). The inference speed of chrysanthemum detection based on deep convolutional
 584 networks has improved significantly (Liu et al., 2020; Liu et al., 2019; Van Nam, 2020). The
 585 inference speed of the VGG16 and ResNet50-based chrysanthemum detection model can reach an
 586 impressive 10ms, about twice as fast as TC-YOLO (Liu et al., 2020), but with only 78% accuracy,

587 which is much lower than the detection accuracy of TC-YOLO at 92.49%. This might be because
588 VGG16 and ResNet50 extract feature information from a single layer for detection, while TC-
589 YOLO not only improves the feature extraction method, but also incorporates a more advanced
590 feature fusion mechanism. It is worth mentioning that most of the current research results for
591 chrysanthemum detection have been tested in relatively ideal environments, while TC-YOLO was
592 tested in nine unstructured environments.

593 Although the TC-YOLO has shown its satisfactory results, there is still space to improve further.
594 First, from the perspective of performance, the detection performance of TC-YOLO under
595 illumination variation needs to be further investigated as this variation limits the overall detection
596 accuracy. One of the possible solutions to solve this is to add more training images collected under
597 different illuminations to fine-tune the model. Currently, as far as the authors' knowledge, there is
598 no publicly available tea chrysanthemum dataset. Tea chrysanthemum usually matures once a year
599 and typically needs to be picked in the early flowering stage. However, the early flowering period
600 of tea chrysanthemum is very short, which leads to the difficulty to obtain a large amount of data
601 for training networks. Generative Adversarial Networks (GANs) ([Goodfellow et al., 2014](#)) are
602 possible to be employed to expand the dataset for training the state-of-the-art deep learning models.
603 These methods are capable of generating realistic images that could increase and diversify the
604 original training datasets. This might result in an improved generalization capability of the learned
605 visual classifiers. TC-YOLO is a lightweight (13.6M) fusion network, but whether it can maintain
606 satisfactory speed and accuracy in a mobile embedded device is still not fully investigated. In our
607 future work, we will explore the possibility of GANs to generate additional tea chrysanthemum
608 datasets. Moreover, the pretrained TC-YOLO will be deployed in an embedded platform such as

609 NVIDIA Jetson and its performance will be evaluated to develop the tea chrysanthemum harvesting
610 robot.

611 **5. Conclusions**

612 This paper presents a lightweight CNN architecture called TC-YOLO to detect tea
613 chrysanthemums in the early flowering stage under complex and unstructured environments. We
614 collected 1000 original images (1080×1920) as the original dataset. Under the NVIDIA Tesla P100
615 GPU environment, the AP reached 92.49% in the test dataset, and the inference speed was 47 FPS.
616 The optimal dataset size is 3,000 images for training TC-YOLO, and the best data augmentation
617 strategy is the combination of flip and rotation. We also learn that overlapping has the least impact
618 on the detection of tea chrysanthemums, while lighting has the greatest impact. The proposed
619 method not only achieves the highest AP (92.49%), which is 2.63% higher than the optimal region
620 detector YOLOv5s based on feature fusion, but also keeps a high inference speed (47.23FPS). The
621 proposed lightweight model TC-YOLO has the potential to be integrated into a selective harvesting
622 robot for automatic early flowering tea chrysanthemums harvesting in the future.

623 **CRedit authorship contribution statement**

624 **Chao Qi:** Conceptualization, Methodology, Software, Writing - original draft, Writing - review &
625 editing. **Junfeng Gao:** Conceptualization, Writing - review & editing. **Simon Pearson:** Writing -
626 review & editing. **Helen Harman:** Writing - review & editing. **Kunjie Chen:** Supervision, Writing -
627 - review & editing. **Lei Shu:** Supervision, Project administration, Funding acquisition, Writing -
628 review & editing.

629 **Declaration of Competing Interest**

630 The authors declare that they do not have any financial or non-financial conflict of interests.

631 **Acknowledgments**

632 This work was supported by Lincoln Agri-Robotics as part of the Expanding Excellence in
633 England (E3) Programme.

634 **References**

635 Altan, A., & Karasu, S. (2020). Recognition of COVID-19 disease from X-ray images by hybrid
636 model consisting of 2D curvelet transform, chaotic salp swarm algorithm and deep learning
637 technique. *Chaos Solitons Fractals*, 140, 110071.

638 Bae, K. I., Park, J., Lee, J., Lee, Y., & Lim, C. (2020). Flower classification with modified
639 multimodal convolutional neural networks. *Expert Systems with Applications*, 159, 11.

640 Cao, L., Zhang, X., Pu, J., Xu, S., Cai, X., & Li, Z. (2020). The Field Wheat Count Based on the
641 Efficientdet Algorithm. *In Proceedings of the IEEE 3rd International Conference on
642 Information Systems and Computer Aided Education* (pp. 557-561).

643 Gao, J. F., French, A. P., Pound, M. P., He, Y., Pridmore, T. P., & Pieters, J. G. (2020). Deep
644 convolutional neural networks for image-based *Convolvulus sepium* detection in sugar beet
645 fields. *Plant Methods*, 16, 12.

646 Gao, J. F., Westergaard, J. C., Sundmark, E. H. R., Bagge, M., Liljeroth, E., & Alexandersson, E.
647 (2021). Automatic late blight lesion recognition and severity quantification based on field
648 imagery of diverse potato genotypes by deep learning. *Knowledge-Based Systems*, 214, 15.

649 Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.,
650 & Bengio, Y. J. (2014). Generative Adversarial Networks, Accessed June 01.
651 <https://ui.adsabs.harvard.edu/abs/2014arXiv1406.2661G>

652 Hou, Y., Li, G., Wang, J., Pan, Y., Jiao, K., Du, J., Chen, R., Wang, B., & Li, N. (2017). Okanin,

653 effective constituent of the flower tea *Coreopsis tinctoria*, attenuates LPS-induced microglial
654 activation through inhibition of the TLR4/NF- κ B signaling pathways. *Scientific Reports*, 7,
655 45705.

656 Kondo, N., Ogawa, Y., Monta, M., & Shibano, Y. Year. Visual sensing algorithm for chrysanthemum
657 cutting sticking robot system. In *Proceedings of the International Society for Horticultural*
658 *Science* (pp. 383-388).

659 Li, Y. T., He, L. Y., Jia, J. M., Lv, J., Chen, J. N., Qiao, X., & Wu, C. Y. (2021). In-field tea shoot
660 detection and 3D localization using an RGB-D camera. *Computers and Electronics in*
661 *Agriculture*, 185, 12.

662 Liu, C. L., Lu, W. Y., Gao, B. Y., Kimura, H., Li, Y. F., & Wang, J. (2020). Rapid identification of
663 chrysanthemum teas by computer vision and deep learning. *Food Science & Nutrition*, 8, 1968-
664 1977.

665 Liu, Y., Sun, P., Wergeles, N., & Shang, Y. (2021). A survey and performance evaluation of deep
666 learning methods for small object detection. *Expert Systems with Applications*, 172, 14.

667 Liu, Z. L., Wang, J., Tian, Y., & Dai, S. L. (2019). Deep learning for image-based large-flowered
668 chrysanthemum cultivar recognition. *Plant Methods*, 15.

669 Sezer, A., & Altan, A. (2021a). Detection of solder paste defects with an optimization-based deep
670 learning model using image processing techniques. *Soldering & Surface Mount Technology*,
671 33, 291-298.

672 Sezer, A., & Altan, A. (2021). Optimization of Deep Learning Model Parameters in Classification
673 of Solder Paste Defects. In *Proceedings of the 3rd International Congress on Human-Computer*
674 *Interaction, Optimization and Robotic Applications* (pp. 1-6).

675 Tarry, C., Wspanialy, P., Veres, M., & Moussa, M. (2014). An Integrated Bud Detection and
676 Localization System for Application in Greenhouse Automation. In *Proceedings of the*
677 *Canadian Conference on Computer and Robot Vision* (pp. 344-348).

678 Tete, T. N., & Kamlu, S. (2017). Detection of plant disease using threshold, k-mean cluster and ann
679 algorithm. In *Proceedings of the 2nd International Conference for Convergence in Technology*
680 (pp. 523-526).

681 Thaïss, C. A., Itav, S., Rothschild, D., Meijer, M. T., Levy, M., Moresi, C., Dohnalová, L.,
682 Braverman, S., Rozin, S., Malitsky, S., Dori-Bachash, M., Kuperman, Y., Biton, I., Gertler, A.,
683 Harmelin, A., Shapiro, H., Halpern, Z., Aharoni, A., Segal, E., & Elinav, E. (2016). Persistent
684 microbiome alterations modulate the rate of post-dieting weight regain. *Nature*, 540, 544-551.

685 Togacar, M., Comert, Z., Ergen, B., & Ozyurt, F. (2020). Classification of brain MRI using hyper
686 column technique with convolutional neural network and feature selection method" *Expert*
687 *Systems with Applications* (vol 149, 113274, 2020). *Expert Systems With Applications*, 150.

688 Van Nam, N. (2020). Application of the Faster R-CNN algorithm to identify objects with both noisy
689 and noiseless images. *International Journal of Advanced Research in Computer Engineering*
690 *& Technology*, 9.

691 Wang, T., Anwer, R. M., Cholakkal, H., Khan, F. S., Pang, Y., & Shao, L. (2019). Learning Rich
692 Features at High-Speed for Single-Shot Object Detection. In *Proceedings of the IEEE/CVF*
693 *International Conference on Computer Vision* (pp. 1971-1980).

694 Warren, D. (2000). Image analysis in chrysanthemum DUS testing. *Computers and Electronics in*
695 *Agriculture*, 25, 213-220.

696 Yang, Q. H., Chang, C., Bao, G. J., Fan, J., & Xun, Y. (2018). Recognition and localization system

697 of the robot for harvesting Hangzhou White Chrysanthemums. *International Journal of*
698 *Agricultural and Biological Engineering*, 11, 88-95.

699 Yang, Q. H., Luo, S. L., Chang, C., Xun, Y., & Bao, G. J. (2019). Segmentation algorithm for
700 Hangzhou white chrysanthemums based on least squares support vector machine. *International*
701 *Journal of Agricultural and Biological Engineering*, 12, 127-134.

702 Yuan, P. S., Ren, S. G., Xu, H. L., & Chen, J. 2018. Chrysanthemum Abnormal Petal Type
703 Classification using Random Forest and Over-sampling. In *Proceedings of the IEEE*
704 *International Conference on Bioinformatics and Biomedicine* (pp. 275-278).

705 Yue, J., Zhu, C., Zhou, Y., Niu, X., Miao, M., Tang, X., Chen, F., Zhao, W., & Liu, Y. (2018).
706 Transcriptome analysis of differentially expressed unigenes involved in flavonoid biosynthesis
707 during flower development of Chrysanthemum morifolium ‘Chuju’. *Scientific Reports*, 8,
708 13414.

709 Zhang, T., & Li, L. (2020). An Improved Object Detection Algorithm Based on M2Det. In
710 *Proceedings of the IEEE International Conference on Artificial Intelligence and Computer*
711 *Applications* (pp. 582-585).

712 Zhang, Z., Qiao, S., Xie, C., Shen, W., Wang, B., & Yuille, A. L. (2018). Single-Shot Object
713 Detection with Enriched Semantics. In *Proceedings of the IEEE/CVF Conference on Computer*
714 *Vision and Pattern Recognition* (pp. 5813-5821).

715 Zheng, J. Y., Lu, B. Y., & Xu, B. J. (2021). An update on the health benefits promoted by edible
716 flowers and involved mechanisms. *Food Chemistry*, 340, 17.

717