

Method for obtaining phenotypic traits of weeds in complex field environments based on instance segmentation

Longzhe Quan*, Bing Wu¹, Shouren Mao¹, Huaiqu Feng¹, Chunjie Yang¹, Wei Jiang¹,

Hengda Li¹

1. College of Engineering, Northeast Agricultural University, Harbin, 150030, China

Abstract

Background:

Weeds are the biggest threat to crop growth, and leaf age and central area of weeds are important phenotypic traits of weeds. They have an important role in understanding the morphological structure of weeds, guiding precision for target weeding and reducing the use of herbicides. However, it is still a substantial challenge to obtain weed types, leaf age and central area in the complex field conditions of light changes, variation in appearance of plants, leaf occlusion. The latest developments in deep learning provide new tools for solving challenging computer vision tasks.

Results:

In this study, we present a weed phenotype segmentation method based on Mask R-CNN that obtains weed types, leaf age and central area in the complex field conditions. By shooting three different angles of the main weeds (*Solanum nigrum*, Barnyard grass, and *Abutilon theophrasti* Medicus) in the field through mobile devices, two datasets (data enhancement and without data enhancement) were produced and used as input to the network, two backbone networks was tested, namely ResNet50 and ResNet101, and the detection results and instance segmentation results of the model were evaluated. The results showed that data enhancement can improve the performance of the model. In the case of data enhancement, the F_1 value

26 with ResNet101 as the backbone network was 0.9214, the mAP scores were 0.6932
27 and 0.5244 (for IOU thresholds of 0.5 and 0.7, respectively), the mIOU reached 0.585,
28 and the best segmentation performance example was obtained. Furthermore, the weed
29 image taken from the top view angle compared to the other two angles achieved the
30 highest detection accuracy.

31 **Conclusion:**

32 Mask R-CNN can achieve accurate segmentation of weeds to obtain the types,
33 leaf age and central area of weeds. Data enhancement and the weed image taken from
34 the top view angle can help to improve the performance of the model. This dataset
35 and research results may provide important resources for the development of
36 precision agriculture in the future.

37 **Keywords:** Leaf age; Mask R-CNN; weed segmentation; Deep learning;
38 Machine vision.

39

40 **Nomenclature**

41 **Symbols**

42	AP	Average precision
43	CNN	Convolutional neural network
44	DCNN	Deep convolutional neural network
45	F_1	Metric function of balance P and R
46	Faster R-CNN	Faster Regions with convolutional neural network features
47	FC	Fully convolutional
48	FCN	Fully convolutional network
49	FN	False Negative
50	FP	False Positive
51	FPN	Feature pyramid network
52	GPU	Graphic processing unit
53	IOU	Intersection over the union
54	mAP	Mean average precision
55	Mask R-CNN	Mask Regions with convolutional neural network features
56	mIOU	Mean intersection over the union
57	P	Precision rate
58	R	Recall rate
59	ROI	Region of interest
60	RPN	Region proposal network
61	Tmax	Maximum thermometer

62 Tmin Minimum thermometer

63 TN True negative

64 TP True positive

65

66 **1. Background**

67 Weeds are the largest threat to crop growth, hindering the growth of crops and
68 promoting the use of herbicides [1], which can lead to a large reduction in crop yield
69 [2], and because weeds compete with major cash crops for space, water, light,
70 nutrients and other resources, the loss of plants is greater in the early stages of plant
71 growth [3]. Therefore, removing weeds as early as possible is critical to ensure good
72 crop yield. The use of herbicides has increased crop yields, and reduced the amount of
73 labour [4]; however, more than 90% of the crops in the United States have had
74 herbicides applied to them [5], and the global use of pesticide is estimated to be 3.5
75 billion kg/year [4]. Excessive use of herbicides has brought about a series of problems
76 such as environmental pollution. Therefore, there is an urgent need to reduce the use
77 of herbicides without affecting the yield of crops, which will also reduce the cost of
78 weed management and is also a valuable goal of precision agriculture [6].

79 The size, shape and growth stage of field weeds are usually uncertain. According
80 to the principle of herbicide and plant physiology [7], the phenotype information of
81 weeds is closely related to the dosage of herbicide. For example, the herbicide dosage
82 required for weeds of different leaf ages is different [8] [9]. The upper limit of
83 herbicide application can kill weeds of different ages, but it can also cause excessive
84 use of herbicides. Leaf age refers to different stages in the development of a plant.
85 Generally speaking, the number of complete leaves grown by the plant is the leaf age
86 [10]. The leaves of the weed are connected to the stem through the petiole. When we
87 look down, it will show the central area formed by the overlapping of the petiole and
88 stem of the top leaves of the weed. This area is also the intersection centre area of
89 weed top leaves [7]. In this study we simply referred to as the central area. Most of
90 the central area is new tissue. Because the epicuticular waxes of plant leaves are

91 closely related to the absorption of herbicides, but the composition of epicuticular
92 waxes in different parts of the same plant is different and will vary with the season,
93 location and age of plants [11] [12], resulting in different organs of the same plant
94 having different sensitivities to herbicides. The new tissue has a large number of
95 stomata and a thin waxy layer, which is beneficial to the absorption of herbicides.
96 Therefore, the central area is more sensitive to herbicides, which is beneficial to
97 herbicide absorption, and the central area has better retention capacity, which is a
98 particularly obvious characteristic especially for weeds over four in leaf age.
99 Therefore, the central area is more sensitive to herbicides and is beneficial to
100 herbicide absorption. Therefore, an effective method is to reduce the dosage of
101 herbicide by guiding the use of herbicide according to the close relationship between
102 the leaf age of the weeds and the central area and the absorption and conduction of the
103 herbicide.

104 The premise of reducing herbicide use is to accurately identify weeds. In recent
105 years, machine vision [13] has been widely used in the agricultural field. Brivot et al.
106 [14] first used machine vision to segment weeds and identify crop rows, proving the
107 feasibility of machine vision in the agricultural field. Researchers have developed a
108 new algorithm that can segment plants according to the soil background under
109 uncontrolled lighting conditions in the wild, separating weeds from crops [15].
110 However, this method is less effective when there is no change in the field image
111 under lighting conditions. Wavelet transform was used to distinguish weeds from
112 crops in the image [16], but when the number of weeds is large, this traditional visual
113 method may not be useful. These methods only distinguish crops from weeds and do
114 not obtain the type and phenotype information of weeds. Although Shirzadifar et al.
115 [17] classified the weeds based on the canopy spectral information of the plants and

116 obtained the species of weeds, achieving good results, in the complex field
117 environment, wind, soil background, and shadows will change the spectral
118 characteristics of plants, affect the performance of the model [18] [19]. A substantial
119 body of related research has also examined the phenotype information of plants.
120 Minervini et al.[20] established finely grained datasets for image-based plant
121 phenotyping, which provided a great contribution to the study of plant phenotypes,
122 and leaf segmentation is also a very important challenge in the field of the plant
123 phenotype analysis. Bell et al. [21] segmented leaves by edge classification and
124 achieved good results for plant overlap. Dobrescu et al. [22] proposed a multi-task
125 deep learning framework for plant phenotypes and achieved good results in leaf
126 counting. But the plant images of these studies were collected under indoor conditions,
127 and images collected indoors tend to have a pure background and light uniformity
128 [23]. These studies are mainly to calculate the number of leaves for leaf segmentation,
129 but an image may have multiple weeds, and our research needs to obtain the leaf age,
130 weed type and the central area of each weed. It can be seen that segmentation of plant
131 phenotypes in a complex farmland environment is still an area that has not been fully
132 studied. Due to the complex environment of farmland, the differences between the
133 plants, and the mutual occlusion of the leaves, it is very difficult to segment the leaf
134 age and the central area, making this study a certain challenge. In this study, the weed
135 species, growth stage and central area were segmented by machine vision in a
136 complex field environment. The purpose is to prepare for guiding the use of
137 herbicides.

138 Deep learning is an emerging field of machine learning that is employed to solve
139 big data analysis problems. The DCNN is a deep learning method that is especially
140 suitable for computer vision problems. In this study, an instance segmentation

141 algorithm based on deep learning is proposed to obtain the weed phenotype in a
142 complex field environment. An agricultural survey has shown that deep learning
143 technology has higher accuracy than traditional image processing technology [24].
144 Because in the complex field environment, the illumination is uncertain, the plants
145 overlap, the climate changes, and the soil background is complex and changeable [25]
146 [26], which makes the DCNN model perform better in this respect. Although the
147 DCNN model can overcome these difficulties, the farmland environment is complex,
148 it needs a large enough dataset to train the deep learning model, which is helpful to
149 overcome the complex field environment and improve the accuracy of the model [27].
150 One method to solve this problem is data enhancement. Data enhancement is also a
151 common method in the field of image recognition. The image is expanded by
152 randomly flipping the image, adding noise, and adjusting the brightness.
153 Geetharamani et al. [28] used a nine-layer deep convolutional neural network to
154 identify plant leaf diseases, and six methods of data enhancement were used to
155 improve the performance of the model, achieving 96.4% classification accuracy.
156 Piedad et al. [29] used the Mask R-CNN model for non-invasive classification of
157 clustered horticultural crops. Due to the limited dataset, the dataset was expanded to
158 improve the accuracy of the model. It can be seen that data enhancement is an
159 important method to enrich the training samples and improve the performance of the
160 model, and can also make the dataset more suitable for the complex farmland
161 environment.

162 When we collect the dataset, the shooting angle in the field [30], the growth
163 stage of the weeds may affect the accuracy of the dataset. Quan and Feng realized the
164 detection of seedlings in different growth cycles and different angles in the field. It
165 was proposed that when the angle between the camera and the vertical direction is 0° ,

166 the detection accuracy is 0.95% lower than other angles [27], It could be seen that
167 data collection from different angles would affect the performance of the model. The
168 position and shape of weeds in the field are complex and changeable, the shape of the
169 same object is different under different shooting angles, which affects the accuracy of
170 the dataset. Therefore, in this study, we collected data from three angles: front view,
171 side view and top view to explore the impact of different angles on the model under
172 study. In this study, an instance segmentation algorithm based on deep learning is
173 proposed to obtain the weed phenotype in a complex field environment. From the
174 previous research, we can know that the DCNN model performs better in dealing with
175 complex environmental problems in the field. Instance segmentation based on deep
176 learning is a new challenge of computer vision [31]. The model used in this study is a
177 new instance segmentation model; its purpose is to detect each object in a weed image
178 and classify each pixel of each instance. The output is the mask and bounding box of
179 the target object [32], which is particularly suitable for solving the problem of leaf
180 adhesion and occlusion [21].

181 Yu et al. [33] proposed an exemplar-based recursive instance segmentation
182 framework to segment plant phenotypes and conducted experiments on a public
183 benchmark to prove the effectiveness of the method. Huang et al. [34] proposed a
184 deep learning model for in-row crop detection in rice fields and constructed a field
185 rice detection dataset with a detection accuracy of 93.22%. It is worth noting that this
186 method identifies a stem-base-centred square region at the plant level, which
187 corresponds to the protected area image of mechanical weeding, and this area is also
188 the central area of the plant. It can be seen that the central area is also of great
189 significance for plant research. The instance segmentation algorithm Mask R-CNN
190 proposed by He et al. [35] can not only identify the bounding box but also mask the

191 target contour, performing better than other models [36]. Jia et al. [37] used an
192 improved Mask R-CNN model to segment overlapping apples with an accuracy rate
193 of 97.31%. Some scholars also used the Mask R-CNN model to complete the fruit
194 detection of the strawberry picking robot in a non-structural environment, which
195 overcame the conditions of overlapping and hidden fruits in a non-structural
196 environment [38].

197 Therefore, the above research provides a feasible basis and reference for the
198 application of DCNN in plant segmentation. It also proves that DCNN can overcome
199 the shortcomings of traditional image segmentation methods. The Mask R-CNN
200 model shows good performance when dealing with complex field environments. The
201 weeds we selected are as follows: *Solanum nigrum*, Barnyard grass, and *Abutilon*
202 *theophrasti* Medicus; these three weeds are commonly found in fields in Northeast
203 China. It can be seen from the above research that we need a sufficient number of
204 well-defined weed datasets to train DCNN models. Weed images should be obtained
205 from real scenes in fields so that they contain the morphological characteristics of the
206 weeds at different growth stages in the complex field environment and cover more
207 variables in the input of the model. We collected three different angle images (front
208 view, side view, top view) of three weeds. The classic DCNN network is modified to
209 improve the accuracy of the model. Therefore, we created two datasets: one
210 containing 4000 weed images without data enhancement, and the other one containing
211 6000 data-enhanced weed images.

212 Based on the above problems, this study proposes a weed phenotype
213 segmentation method based on the improved Mask R-CNN to obtain the weed species,
214 leaf age and central area of weeds. The main objectives of the present study are

215 (1) Evaluate the feasibility of Mask R-CNN to obtain weed species, leaf age and

216 central area in weed phenotype segmentation through seven evaluation indicators.

217 (2) To explore the influence of data collection from different angles (front view,
218 side view, top view) on the phenotypic segmentation of weeds in the complex field
219 environment, and select the most suitable angle.

220 (3) To explore whether data enhancement can improve the performance of the
221 model.

222 (4) To explore whether using Resnet101 combined with the FPN architecture for
223 feature extraction can improve model performance.

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

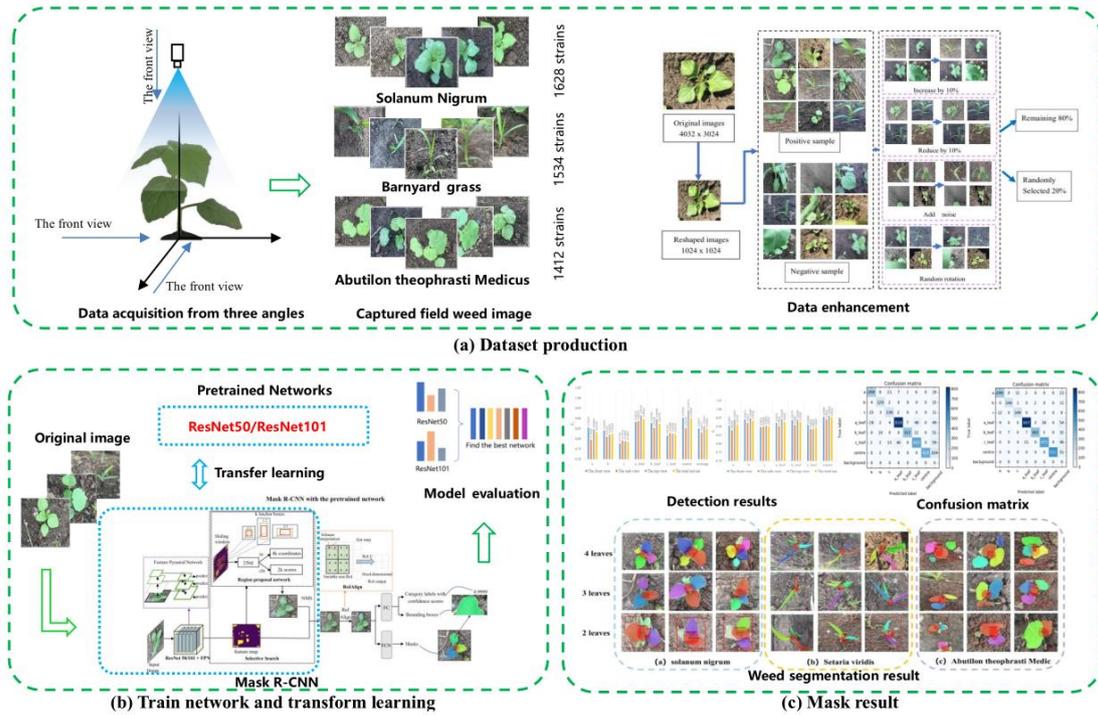
244

245

246 **2. Materials and methods**

247 **2.1. Overview**

248 In this section, we select three typical weeds of *Solanum nigrum*, Barnyard grass,
249 and *Abutilon theophrasti* Medicus from Northeast China. The leaf age and central area
250 of weeds were segmented by the Mask R-CNN model. The dataset needed to train the
251 model was collected in the actual field environment. The datasets were collected at
252 different angles (front view, side view, top view) and were enhanced. The datasets
253 produced were annotated, and the generated file was used as the input of the network
254 to train the network model. The backbone network of the Mask R-CNN initialization
255 model is the residual network combined with the feature pyramid network (FPN).
256 This study uses different backbone networks (ResNet50 and ResNet101) combined
257 with the FPN architecture. The feature extraction performance based on the weed leaf
258 age and central area was evaluated. Figure 1 shows the workflow of the experiment.



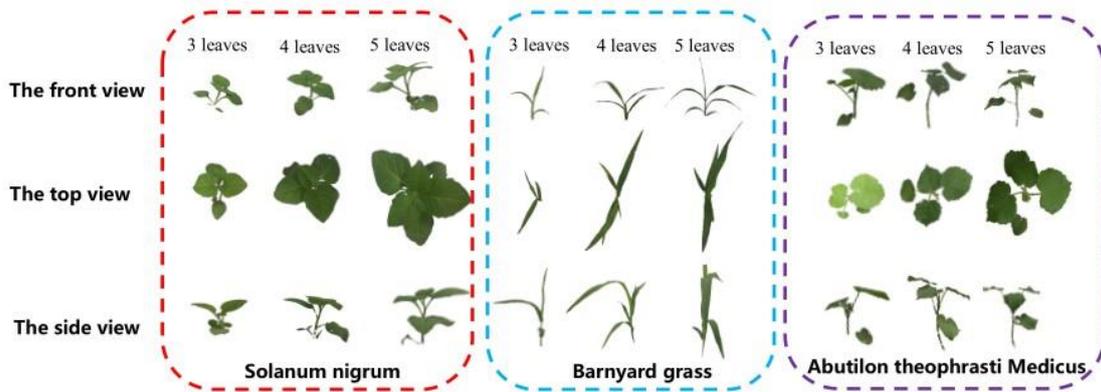
259

260 Fig. 1 Work process of using Mask R-CNN to segment leaf age and central area.

261 **2.2. Image acquisition**

262 The three kinds of weeds selected in this study are: Solanum nigrum, Barnyard
 263 grass, and Abutilon theophrasti Medicus. Solanum nigrum is an annual dicotyledon,
 264 and Barnyard grass is an annual herb, an annual subshrub weed of Abutilon
 265 theophrasti Medicus. These three kinds of weeds are common in fields of Northeast
 266 China. At present, the data image source is the field weed image. Because the
 267 background of greenhouse weeds is single and the image of field weeds is more
 268 complex, the ability of the model to recognize weeds in the natural state can be
 269 verified. The field data images were collected in the Xiangfang District from May to
 270 June 2018. The Xiangfang District is located in the northeast plain and is the main
 271 planting area for maize, soybean and rice. The main weeds in the cornfields of the
 272 Xiangfang District are Solanum nigrum, Barnyard grass, and Abutilon theophrasti
 273 Medicus. The above weed images were collected. The weeds in the field are mainly
 274 weeds of leaf ages of 2-5; thus, we collected only weed images before the 5 leaf stage.

275 Since the weed information obtained from a single shooting angle is not
 276 comprehensive, in order to better show that the weed information obtained from
 277 different angles is different, the field weed background is removed, as shown in
 278 Figure 2. In addition, the shooting weather [39] and acquisition angle [30] had a
 279 greater impact on the dataset, affecting the segmentation precision [40]. Therefore,
 280 from the two-leaf period after crop planting, an iPhone 6s Plus camera with a 4.2 mm
 281 focal length, a maximum aperture of f/2.2 and a maximum resolution of 4032 x 3024
 282 pixels was used to capture images. From May 20, 2019, to June 29, 2019, weed
 283 images of different leaf ages were collected every 2 to 5 days under different weather
 284 conditions, different angles, and different growth stages, to obtain the data of weeds at
 285 each leaf age stage in the growth cycle, as shown in Table 1, and the weed images
 286 were stored in the JPEG file format. The purpose of generating the dataset is to study
 287 the recognition performance of the deep learning model for the leaf age and central
 288 area of individual weeds at different growth stages in the natural state.



290 Fig. 2 Front view, top view, and side view of the three weeds.

291 Table 1 List of images containing environmental information for the experiment

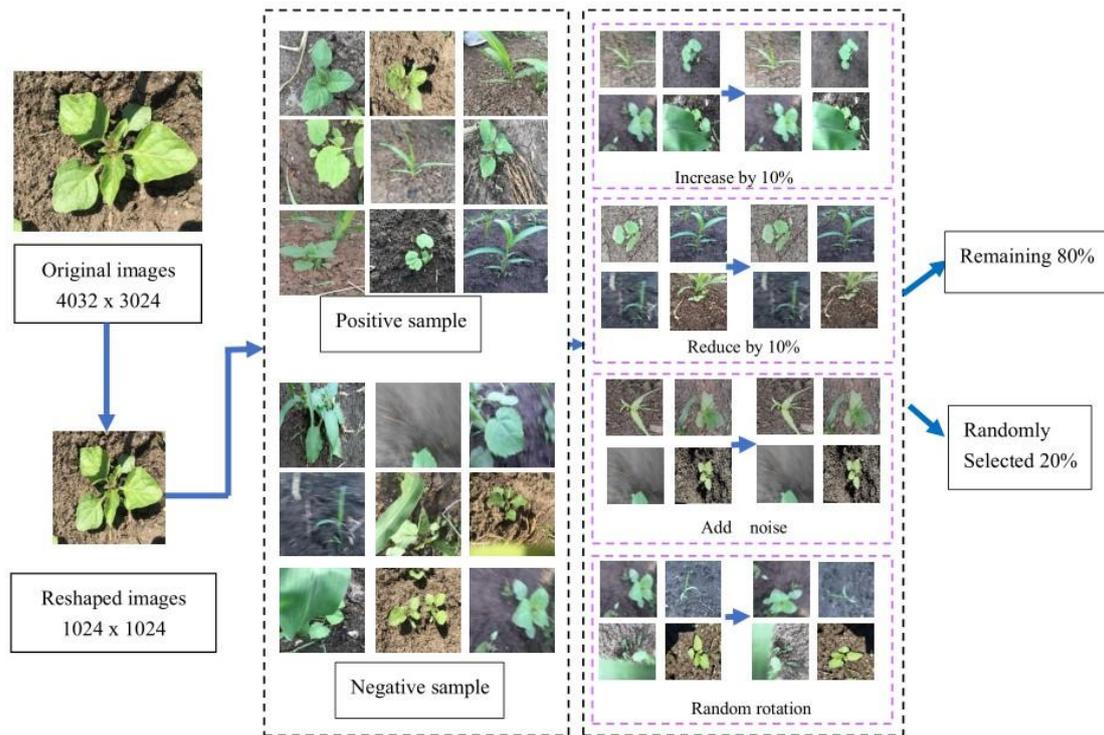
Date	Images	Tmax	Tmin	Weather	The front view	The top view	The side view
25/05/2019	441	30	17	Cloud	201	140	100
30/05/2019	449	22	11	Cloud	112	226	111

04/06/2019	481	24	16	Cloud	132	200	149
08/06/2019	538	24	18	Rain	213	221	104
10/06/2019	422	25	16	Cloud	136	161	125
14/06/2019	426	24	14	Rain	119	203	104
17/06/2019	458	25	13	Cloud	102	254	102
21/06/2019	420	26	14	Cloud	125	107	188
26/06/2019	412	26	17	Cloud	119	136	157
29/06/2019	527	23	15	Cloud	152	214	161
Total	4574	28	16	\	1411	1862	1301

292

293 **2.3. Dataset construction and annotation**

294 When training the network and conducting network testing, the input image size
295 needs to match the input size of the network [41], thus, the images are adjusted to a
296 pixel size of 1024 x 1024 to construct the image dataset of the DCNN. Due to the
297 need to annotate the weeds in the picture, some images that are not suitable for
298 annotation are discarded. Therefore, 4000 images were selected from 4574 images.
299 Due to the limited number of datasets, a data enhancement scheme was adopted to
300 further enrich the images such that they are more representative and reflect the real
301 situation of field data more accurately [27], while additionally improving the training
302 precision of the model [42], expanding the dataset and reducing overfitting [43]
303 (Figure 3).



304

305

Fig. 3 Data enhancement.

306

307

308

309

310

311

312

313

314

315

316

317

318

The images were randomly rotated, and noise was added. Since illumination is an important reason for segmentation, to make the DCNN more robust to the illumination caused by environmental changes, the datasets were further enhanced by simulating the illumination change [44]. The brightness was adjusted to be 10% brighter and 10% darker. At the same time, some blurred, occluded and incomplete images were retained as the dataset of negative samples, and a total of 6000 data-enhanced pictures were obtained. The structure and proportion of the original dataset remain unchanged when data enhancement is carried out. Two datasets were made, one with data enhancement and the other without data enhancement. We reserve 600 of the 4000 images without data enhancement and use them to test the model for model evaluation. Both datasets were randomly divided into a training set and a verification set with a ratio of 8:2. After the training was completed, we used the 600 images previously reserved to test the model to evaluate the model. The VGG

319 Image Annotator labelling tool was used for annotation, as shown in Figure 4, and
320 weed leaves were surrounded by irregular polygons, while the centre areas of the
321 weeds were marked with a circle. Because the number of weeds in a picture is
322 uncertain under actual working conditions and may contain multiple weeds, the
323 number of masked leaves in the picture cannot be used to calculate the leaf age of a
324 single weed, so we used a rectangular frame to mark the outline of the outermost layer
325 of a single weed and calculated the number of leaf masks in the rectangular frame,
326 which is the leaf age of the weed. The rectangular frame is not masked. The labels are
327 divided into 7 categories (Figure 4).



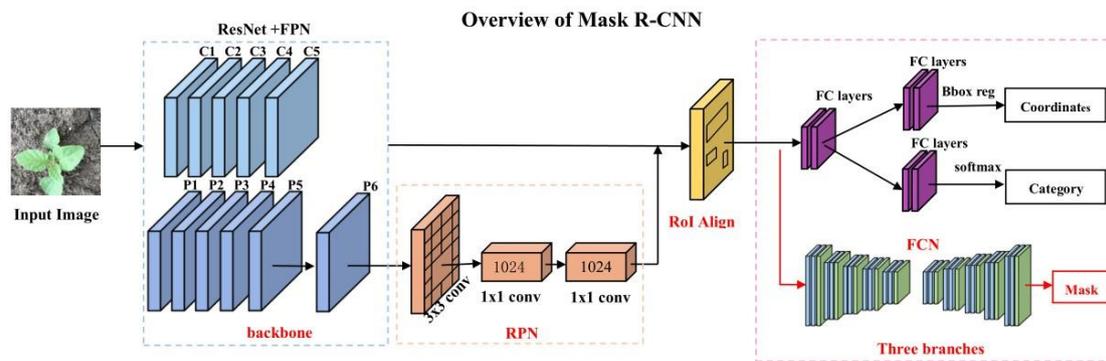
329 Fig. 4 Image annotation example of segmentation.

330 Note: Leaf of Solanum nigrum (a_leaf), leaf of Barnyard grass (b_leaf), leaf of
331 Abutilon theophrasti Medicus (c_leaf), central area (centre), Solanum nigrum (a),
332 Barnyard grass (b), and Abutilon theophrasti Medicus (c).

333 2.4. Structure of Mask R-CNN mode

334 Mask R-CNN extends the target detection framework of Faster R-CNN [45] by
335 adding a masking branch at the end of the model [38]. This process ensures that each
336 output instance segments the proposal box using a fully connected (FC) layer to
337 ensure that the segmentation is parallel to the target detection. To better detect small
338 targets, ROI pooling is changed to ROIAlign. The Mask R-CNN model flow chart is

339 shown in Figure 5. The output consists of three branches: bounding boxes, target
 340 classifications and segmentation masks. The Mask R-CNN backbone networks we
 341 selected are ResNet50 and ResNet101 combined with FPNs. First, the backbone
 342 network extracts the feature map from the input image and then outputs the features
 343 from the backbone network. The map is sent to the region proposal network (RPN)
 344 and ROIAlign to generate the region of interest (ROI). Finally, the ROI predicts the
 345 target category and bounding box through the convolutional layer and the fully
 346 connected layer and segments the target region through the fully convolutional neural
 347 network (FCN). The instance segmentation task of the target is finally completed.

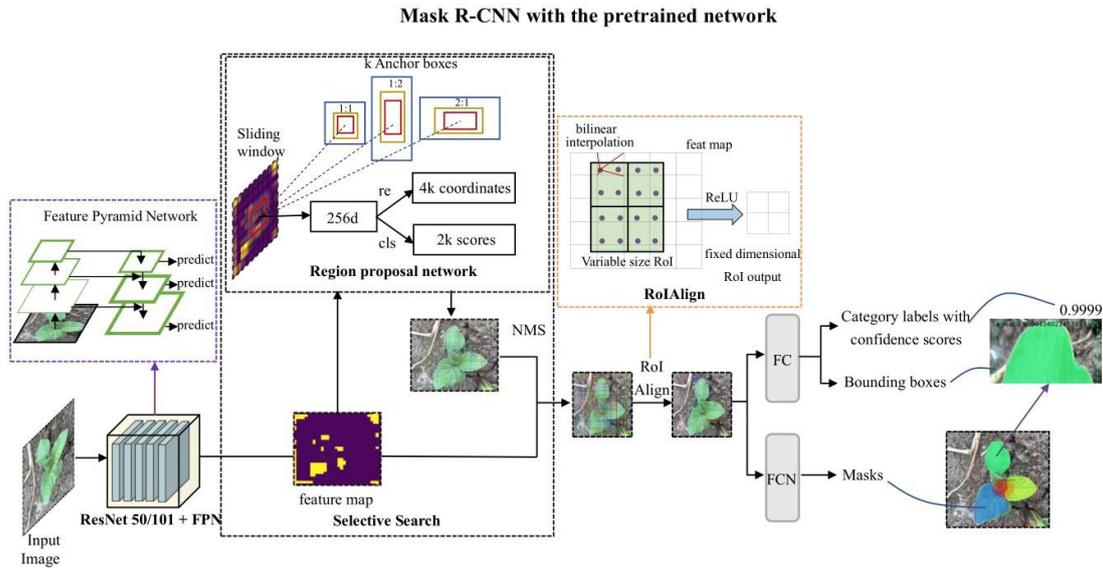


348

349 Fig. 5 The model structure of Mask R-CNN.

350 The FPN is an important module of Mask R-CNN [46] that detects objects of
 351 different scales in object detection. There are top-down and bottom-up paths as well
 352 as lateral connections. As shown in the FPN part of Figure 6, the thicker the contour is,
 353 the stronger the semantics. While the RPN is a sliding window in target detection,
 354 Figure 6 shows that the RPN method is used to generate a region proposal network.
 355 ROIAlign is a very important part of Mask R-CNN. It is an improvement over
 356 ROI Pooling. As shown in the ROIAlign part of Figure 6, it first traverses each
 357 candidate area while keeping the floating point number unchanged and then divides

358 the candidate area into $k \times k$ units. Four fixed coordinate values are calculated in
 359 each unit, and the values of the four positions are calculated by a bilinear interpolation
 360 method. Then, a maximum pooling operation is performed.



361

362

Fig. 6 Mask R-CNN with a pretrained network.

363

2.4.1 Image Feature Extraction (ResNet50 and ResNet101 + FPN)

364

365

366

367

368

369

370

371

372

Mask R-CNN can establish deep neural network models of different depths by designing different weight layers. Currently, the deep learning network models are AlexNet, ZF, GoogLeNet, VGG, and ResNet. Although a deeper number of network layers may lead to higher accuracy, the deeper network layers will result in lower model training and detection speeds. However, since the residual network does not increase the parameters of the model, it can effectively reduce the problem of training degradation and improve the model convergence [38]. Therefore, this paper uses ResNet50 and ResNet101 combined with the FPN as the backbone network to extract the features of weed images.

373

2.4.2 ROI Generation and ROIAlign (RPN+ROIAlign)

374

The RPN uses the convolutional feature map output from the backbone network

375 as its own input, traversing the feature map in a fixed-size sliding-window fashion and
376 finding images containing weed blades and central areas. This area of the RPN scans
377 is called an anchor [47]. The weed target in this study is small, especially in the
378 central area. The size of the weeds changes throughout the growth cycle [38].
379 Therefore, according to the total number of pixels of the target weed in the image, this
380 study designed 32×32 , 64×64 , 128×128 , 256×256 , and 512×512
381 anchors. The length-width ratios were designed to be 1:1, 1:2, and 2:1. A total of
382 fifteen anchor points were selected, and the most likely target area was selected for
383 segmentation and detection.

384 Faster R-CNN [45] is a recognition algorithm proposed for the regional
385 convolutional network R-CNN to perform a large number of repetitive operations in
386 each ROI and then transfer the features extracted by the RPN to the convolution of the
387 last layer. The ROI Pool layer is then added later. However, when performing the
388 quantization operation in ROI Pool, misalignment between the ROI and the extracted
389 features can easily occur, which will have a great negative impact on the results,
390 making this layer unsuitable for this study. Mask R-CNN [35] uses ROI Align instead
391 of the traditional pooling of interest (ROI Pool), which is very suitable for addressing
392 smaller targets and solves the problem of spatial position misalignment caused by
393 ROI Pool. ROI Align uses bilinear interpolation to extract the corresponding features
394 of each ROI on the feature map [38], calculates the exact value of each position,
395 summarizes it by the maximum pooling method, adjusts the dimensions of each ROI
396 to meet the FC requirement, and finally sends it to the FC layer and FCN for target
397 classification, the bounding boxes and the mask.

398 **2.5 Mask R-CNN training model**

399 Before training Mask R-CNN, we introduced a pretraining model based on the

400 COCO dataset [48] using transfer learning. The COCO dataset has 328k images,
 401 including 91 categories. The pretraining model extracted the weights after training on
 402 COCO. On this basis, the datasets established by themselves will be trained again. By
 403 means of this transfer learning, the manpower and cost of training can be reduced, the
 404 training efficiency can be improved, and the parameters of the model can be better
 405 adjusted. The Mask R-CNN model of this experiment was carried out under the deep
 406 learning framework of TensorFlow-gpu 1.14.0 and Keras 2.1.5. This study is based on
 407 the Windows 10 64-bit (DirectX 12) operating system, a six-core Intel Core i7-8700K
 408 @ 3.70 GHz processor, 32 GB of memory, and a GPU built by NVIDIA GeForce
 409 (Santa Clara, CA, USA), the NVIDIA GeForce RTX 2080 Ti graphics card. The
 410 parameters of the pretraining network are shown in Table 2.

411 Table 2 Characteristic parameters of the pretraining network.

Network	Depth	Size (MB)	Parameters (Millions)	Image Input Size	Feature Extraction Layer	ROI Pooling Output Size
ResNet50	50	96	25.6	224 x 224	block_13_expand_relu	[14 14]
ResNet101	101	167	44.6	224 x 224	mixed7	[17 17]

412 Insert ROIAlign after the feature extraction layer, extract the corresponding
 413 features of each ROI, calculate the exact value of each position, and finally transport
 414 it to the fully connected layer and FCN to predict the target class, bounding box and
 415 segmentation mask. We put the marked dataset into the model for training and obtain
 416 a training model based on Mask R-CNN for weed segmentation.

417 2.6 Training and evaluation

418 The weight decay coefficient of Mask R-CNN was set to 0.0005, the momentum
 419 was set to 0.9, the initial learning rate was set to 0.0001, and the training BatchSize
 420 was set to 1. After the parameters were set, training was conducted for 100 rounds,
 421 with training being conducted 1000 times per round. The basic framework of Mask

422 R-CNN involved either ResNet50 or ResNet101.

423 The purpose of the evaluation was to test the ability of the algorithm to identify
424 the weed leaf age and the central area on the image. We used seven key indexes:
425 precision rate (P) (Eq. (1)), recall rate (R) (Eq. (2)), F_1 (Eq. (3)), intersection over
426 union (IOU) (Eq. (4)), average precision (AP), mean average precision (mAP) (Eq.
427 (5)), and mean intersection over the union (mIOU), for evaluation. The Mask R-CNN
428 model completes object detection and object segmentation. For the object detection
429 part, we used the precision rate (P), recall rate (R), F_1 , AP, and mean average precision
430 (mAP) for evaluation. For the object segmentation part, we use the mean intersection
431 over the union (mIOU) error for evaluation.

432 The precision rate and recall rate are defined by the following equation:

$$Precision = \frac{TP}{TP + FP} \quad \#(1)$$

$$Recall = \frac{TP}{TP + FN} \quad \#(2)$$

433 where "true positive (TP)" indicates the number of positive results detected as
434 positive; "false positive (FP)" indicates the number of negative results detected as
435 positive; and "false negative (FN)" indicates the number of positive results detected as
436 negative. The metric function (F_1) [49] of the precision and recall rates is defined by
437 the following equation:

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad \#(3)$$

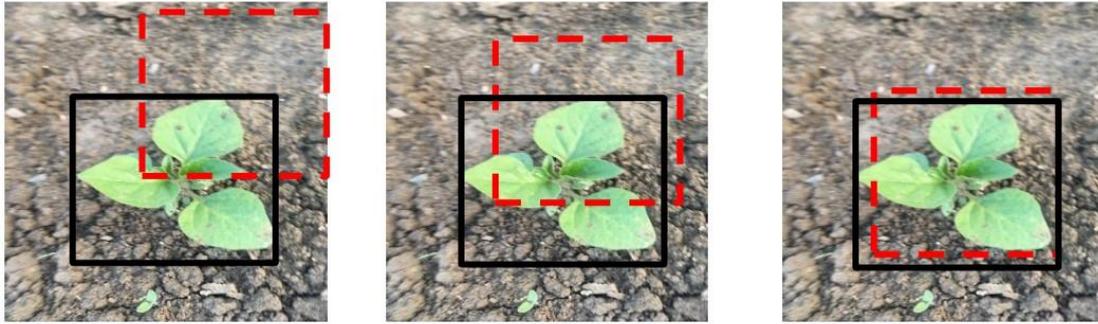
438 To better meet the problem of multiclass imbalance in this paper, we averaged
439 the seven classification indicators [50]. To better evaluate the model algorithm, the
440 IOU is used to conduct measurements [51]. The IOU measures the overlap between
441 two bounding boxes. As shown in Figure 7, the overlap degree between the weed
442 prediction box and the real box on the ground was calculated. The value of the IOU

443 can be divided into three regions. When the IOU threshold is set to 0.7, anchors with
444 IOU values less than or equal to 0.3 are considered to be negative anchors, while
445 anchors with IOU values between 0.3 and 0.7 are neutral anchors, and we do not need
446 to consider these cases. In addition, IOU values greater than or equal to 0.7 are
447 positive anchors. The system will identify the positive anchors and bounding box and
448 match them to the ground-truth boxes to optimize the RPN output of the model. The
449 maximum value of the anchor overlap with the ground-truth boxes is retained by the
450 system. When the IOU threshold is set to 0.5, the IOU values are greater than 0.5, and
451 a positive ROI is observed; when the IOU values are less than 0.5, a negative ROI is
452 observed. The positive ROI is allocated to the mask and ground truth by the system.
453 The detection performance of the model is evaluated by the mean accuracy (mAP)
454 [51]. The mAP has excellent evaluation performance when it is related to the target
455 position information and the category information of the target in the image. The AP
456 can be calculated for each category separately, and then each category can be
457 averaged to calculate the mAP. The larger the mAP value is, the better it is. It can be
458 calculated by the AP. The thresholds were 0.5 and 0.7 in this study. The IOU and mAP
459 are defined by the following equation:

$$IOU = \frac{Area\ Overlap}{Area\ Union} \quad \#(4)$$

$$mAP = \frac{1}{N} \cdot \sum_{i=1}^N AP_i \quad \#(5)$$

460 Note: N represents the number of images.



(a)

(b)

(c)

461

462 Fig.7 Visual example of intersection over union. The red dashed box is the prediction

463 of the detection result. The black solid-line box is the ground truth. The overlap

464 between the two is visible, and different IOU values are depicted from left to right. (c)

465 is the best.

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

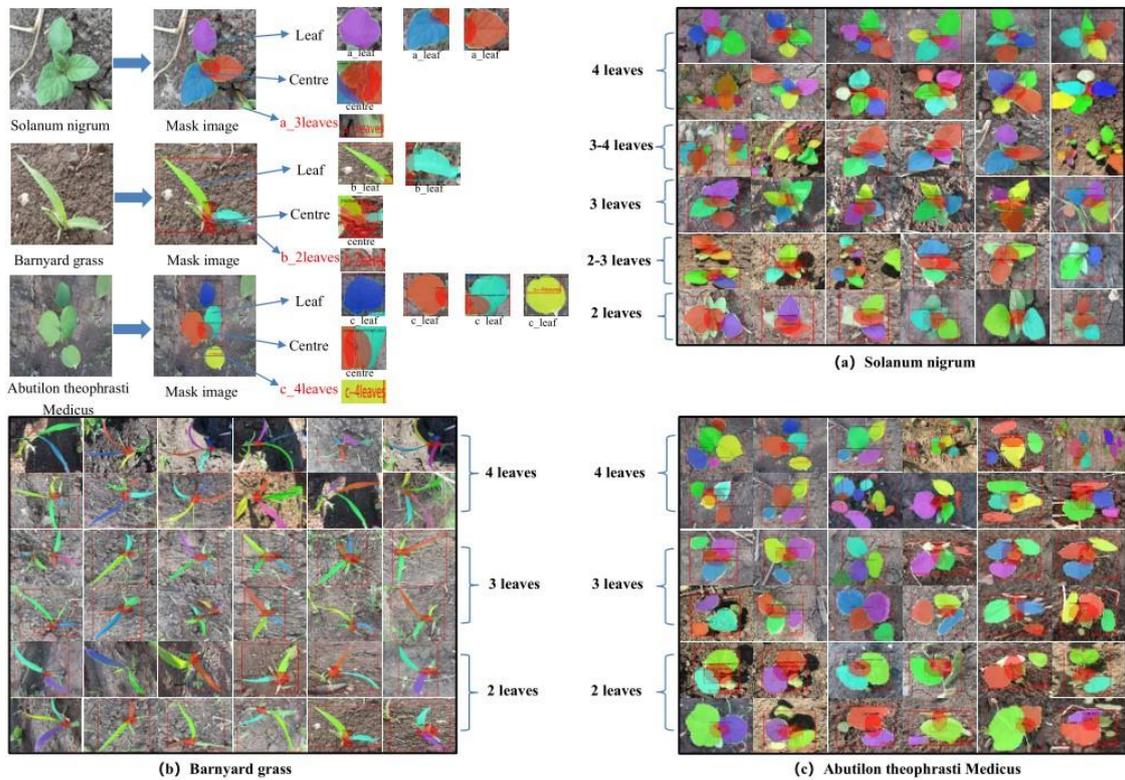
481

482 **3 Results**

483 In the following three stages, the development of a segmentation model based on
484 Mask R-CNN for weed leaf age and central area was completed. First, the image
485 acquisition work was completed in three angles of front view, left view and top view
486 and under different weather conditions, and 4574 weed images were collected. Second,
487 the data pre-processing and image annotation of the acquired image was completed.
488 Two datasets were created: one dataset contains 4000 weed images without data
489 enhancement, and the other dataset contains 6000 data-enhanced weed images. Third,
490 two selected pretrained networks were used to replace the backbone network of the
491 classic Mask R-CNN, network construction and parameter optimization, and an
492 improved Mask R-CNN model processed by the pretraining network was proposed.
493 The training and evaluation of the model were then completed. Data enhancement can
494 improve the performance of the model. The mAP of the Mask R-CNN model using
495 ResNet101 as the backbone network is 0.693, and the mIOU is 0.585, which is better
496 than ResNet50 and can be used for weed segmentation. Moreover, the weed image
497 taken from the top view angle compared to the other two angles has the highest
498 detection accuracy.

499 The result of weed segmentation is shown in Figure 8. A method based on Mask
500 R-CNN for segmentation of weed leaf age and central area is proposed. Two datasets
501 are trained in two different backbones networks (ResNet50 and ResNet101), and the
502 network is chosen that achieves the best balance between mIOU and mAP. The leaf
503 age is determined based on the number of complete leaves of the weed, and the centre
504 area is determined based on the intersection area of the top leaves of the weed [7]. The
505 evaluation of the Mask R-CNN model is carried out from two aspects: first, the

506 detection results of the model are evaluated, and then the segmentation performance
 507 of the model is evaluated.



508

509

Fig. 8 Weed segmentation results.

510 3.1 Weed test results and evaluation

511 The precision rate (P), recall rate (R), F_1 , IOU, AP, and mAP are important
 512 indicators used to evaluate the performance of Mask R-CNN detection. Table 3 lists
 513 the F_1 and mAP values of different datasets (without data enhancement and data
 514 enhancement) on the total test set. The number of the total test set is 600. The total
 515 test set includes 200 front view images, 200 side view images, and 200 top view
 516 images. It can be seen from Table 3 that when data enhancement is used, the F_1 value
 517 of the ResNet101 network in the total test set is 0.9214. Without data enhancement,
 518 the F_1 value of the ResNet101 network on the total test set was 0.8965. The F_1 value
 519 of ResNet50 in the case of data enhancement is 0.0413 higher than the F_1 value
 520 without data enhancement. It can be seen that for the dataset with data enhancement,

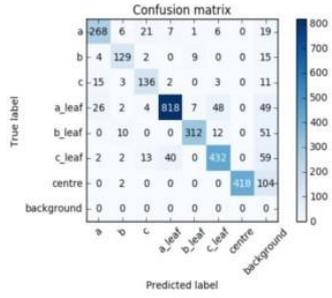
521 the detection results of the two backbone networks ResNet50 and ResNet101 on the
522 total test set are higher than those without data enhancement. It can be seen that the
523 data enhancement effect is better than that without data enhancement. Table 3 also
524 lists the values of mAP under the two thresholds of data enhancement and without
525 data enhancement in the total test set. According to Table 3, for the total test set, when
526 the IOU threshold of ResNet101 is greater than or equal to 0.5, the mAP value is
527 0.6512 without data enhancement, while the mAP value is 0.6932 with data
528 enhancement. It can be seen that no matter which backbone network is used, the map
529 value of the model is greater than that without data enhancement. Although it is only
530 increased by 0.042, it also improves the performance of the network and can reduce
531 the impact of overfitting.

532 Table 3 List of the detection results of different datasets on the total test set

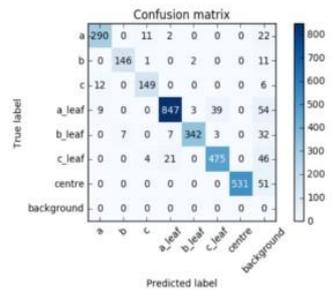
dataset		Without data augmentation	Data augmentation
	F_1	0.8073	0.8486
ResNet50	mAP (IOU \geq 0.5)	0.5421	0.5702
	mAP (IOU \geq 0.7)	0.4534	0.4791
	F_1	0.8965	0.9214
ResNet101	mAP (IOU \geq 0.5)	0.6512	0.6932
	mAP (IOU \geq 0.7)	0.5113	0.5244

533 Note: The total test set consisted of 600 images, including 200 front view images, 200
534 side view images, and 200 top view images.

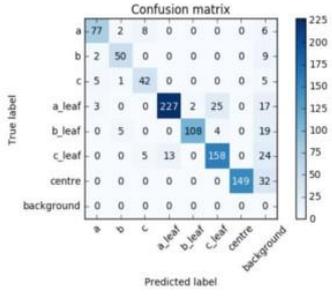
535
536 Figure 9 shows the confusion matrices of the detection results of the model in the
537 case of data enhancement. Figure 10 lists the detection results of the Mask R-CNN
538 model under two backbone networks, three angles, and seven types of labels in the
539 case of data enhancement. The precision rate of the ResNet101 network in the total
540 test set is 0.9523, the recall is 0.8929, and the F_1 value is 0.9214.



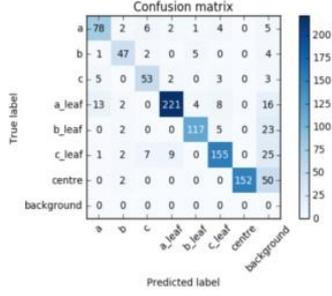
(a) The total test set under ResNet50



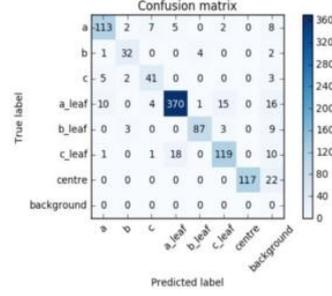
(b) The total test set under ResNet101



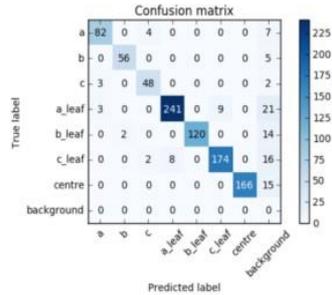
(c) The front view under ResNet50



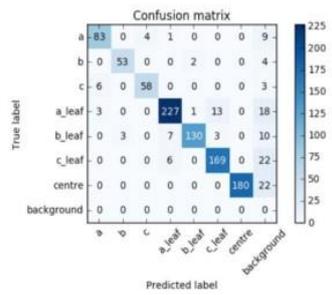
(d) The side view under ResNet50



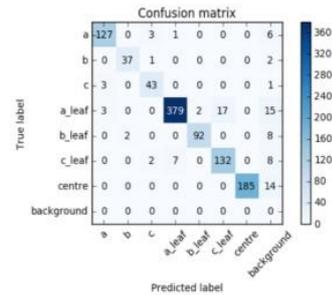
(e) The top view under ResNet50



(f) The side view under ResNet101



(g) The side view under ResNet101



(h) The top view under ResNet101

541

542

Fig. 9 Confusion matrix of the detection results of ResNet50 and ResNet101 in the

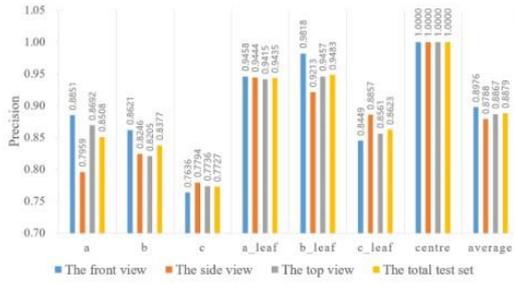
543

case of data enhancement

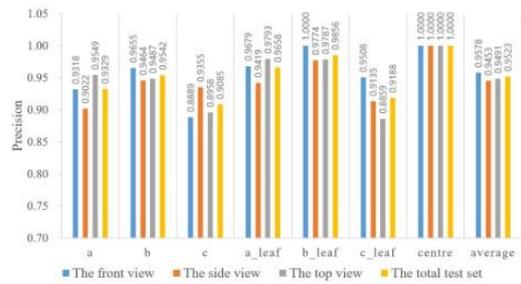
544

Note: The Mask R-CNN with ResNet50 is abbreviated as ResNet50, and Mask R-CNN with ResNet101 is abbreviated as ResNet101. a_leaf represents the leaves of Solanum nigrum, b_leaf represents the leaves of Barnyard grass, c_leaf represents the leaves of Abutilon theophrasti Medicus, a represents Solanum nigrum, b represents Barnyard grass, c represents Abutilon theophrasti Medicus, and centre represents the central area of each weed. The total test set consisted of 600 images, including 200 front view images, 200 side view images, and 200 top view images.

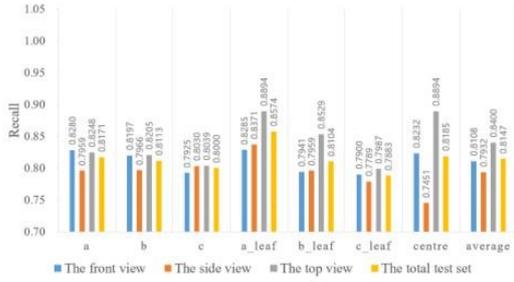
550



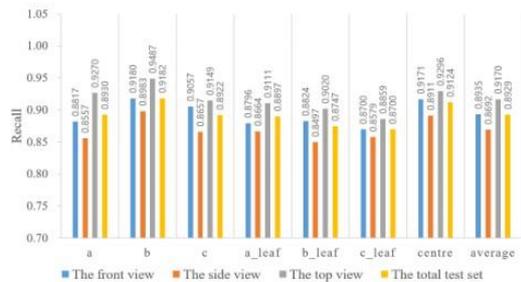
(a) Precision of ResNet50



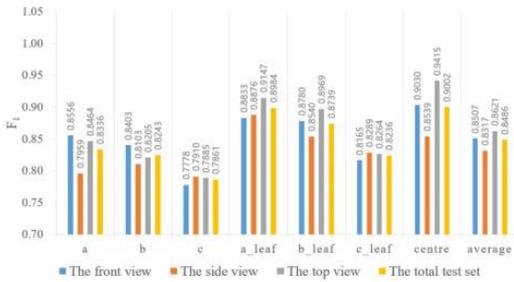
(b) Precision of ResNet101



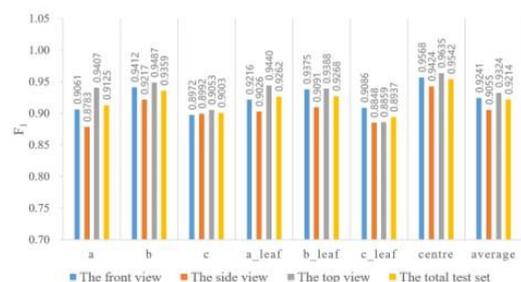
(c) Recall of ResNet50



(d) Recall of ResNet101



(e) F_1 of ResNet50



(f) F_1 of ResNet101

551

552 Fig. 10 Detection results of the Mask R-CNN with pretrained networks in the case of
 553 data enhancement.

554 Note: The Mask R-CNN with ResNet50 is abbreviated as ResNet50, and Mask
 555 R-CNN with ResNet50 is abbreviated as ResNet101. a_leaf represents the leaves of
 556 Solanum nigrum, b_leaf represents the leaves of Barnyard grass, c_leaf represents the
 557 leaves of Abutilon theophrasti Medicus, a represents Solanum nigrum, b represents
 558 Barnyard grass, c represents Abutilon theophrasti Medicus, and centre represents the
 559 central area of each weed. The total test set consisted of 600 images, including 200
 560 front view images, 200 side view images, and 200 top view images.

561

562 It can be seen from Figure 10 that the labels a, b, c, a_leaf, b_leaf, c_leaf, and
 563 centre were not recognized, we surmise that they were recognized as the background.

564 It can be seen from Figure 10 that under the condition of data enhancement, the
 565 precision values of ResNet50 in the front, side, top, and total test sets were greater
 566 than or equal to 0.7636, the recall values were greater than or equal to 0.7451, and the

567 F_1 values were greater than or equal to 0.7778. In comparison, the precision values of
568 ResNet101 were greater than or equal to 0.8859, the recall values were greater than or
569 equal to 0.8497, and the F_1 values were greater than or equal to 0.8783 in the front,
570 side, top, and total test sets. As a result, the precision, recall and F_1 values of
571 ResNet101 in the front, side, top, and total test sets were considerably higher than
572 those of ResNet50, and ResNet101 consistently performed better than ResNet50. The
573 F_1 values of ResNet101 in the front, side, top, and total test sets were 0.9241, 0.9055,
574 0.9324 and 0.9214, respectively. When ResNet101 was used as the backbone network,
575 the recall values of the top view test set were greater than or equal to 0.8859, the
576 recall values of the front view and the total test sets were greater than or equal to
577 0.8700, and the recall values of the side view test set were greater than or equal to
578 0.8497. The top view test set shows the best performance in all classifications. For the
579 total test set, the F_1 values were 0.9359 and 0.9268, and the recall rates were 0.9182
580 and 0.8747 when Barnyard grass and the leaves of Barnyard grass were detected, and
581 on the front, side, and top test sets, the classification performance of Barnyard grass
582 was better than the other two kinds of weeds. For the central area, the precision values
583 of ResNet101 in the front, side, top, and total test sets were 1.0000. Since Figure 9
584 and Figure 10 can show only the classification performance of the model, the
585 recognition accuracy of the model cannot be determined, and the actual environment
586 in the field is complex, which has a certain impact on the identification of weeds.
587 Therefore, the model recognition accuracy is very important for evaluating the model
588 performance. Table 4 lists the mAP values for ResNet50 and ResNet101 for different
589 test sets when IOU values are greater than or equal to 0.5 and greater than or equal to
590 0.7 under data enhancement.

591 Table 4 List of detection results of weeds under different networks and different

592

angles with data augmentation.

Networks		ResNet50	ResNet101
The front view	mAP (IOU \geq 0.5)	0.552	0.701
	mAP (IOU \geq 0.7)	0.456	0.522
The side view	mAP (IOU \geq 0.5)	0.533	0.625
	mAP (IOU \geq 0.7)	0.468	0.458
The top view	mAP (IOU \geq 0.5)	0.625	0.753
	mAP (IOU \geq 0.7)	0.513	0.592
The total test set	mAP (IOU \geq 0.5)	0.570	0.693
	mAP (IOU \geq 0.7)	0.479	0.524

593

The mAP is a commonly used index in target detection. Table 4 shows that the

594

mAP of ResNet101 is higher than that of ResNet50, indicating that ResNet101 has

595

good target detection performance. For the total test set, when the IOU threshold is

596

equal to or greater than 0.5, the mAP value is 0.693, and when the IOU threshold is

597

equal to or greater than 0.7, the mAP value is 0.524. Thus, when the threshold is equal

598

to or greater than 0.5, ResNet101 has good detection performance. When ResNet101

599

is used as the backbone network and the IOU threshold is greater than or equal to 0.5,

600

the mAP value in the top view test set is 0.753, the mAP value in the front view test

601

set is 0.701, the mAP value in the side view test set is 0.625, and the mAP value in the

602

total test set is 0.693. The top view achieved good detection performance.

603

3.2 Instance segmentation results and evaluation

604

For the instance segmentation phase, two backbone networks ResNet50 and

605

ResNet101 were executed. From the above comparison, we can observe that the effect

606

of data enhancement is better than that without data enhancement. Therefore, we

607

tested the data-enhanced dataset, and the results of weed segmentation are shown in

608

Figure 8. The mIOU is an important index for evaluating the segmentation results [31]

609

and is commonly used to evaluate the segmentation performance of the Mask R-CNN

610

model. Table 5 shows the mIOU of Mask R-CNN under different backbone networks

611

and different test sets. The test image of 600 weeds shows that for the total test set, the

612 mIOU of ResNet50 is 0.491 and that the mIOU of ResNet101 is 0.585. The value for
 613 ResNet101 is significantly higher than that of ResNet50. Thus, the ResNet101 model
 614 will result in better network performance and can be applied to the segmentation of
 615 small target objects. Hence, it can meet the needs of weed instance segmentation.
 616 When ResNet101 was used as the backbone network, the mIOU of the top view was
 617 0.603, which is higher than those of other datasets. Therefore, good segmentation
 618 results were achieved.

619 Table 5 Segmentation results of weeds under different networks and different
 620 angles in the case of data enhancement.

Angles	mIOU	
	ResNet50	ResNet101
The front view	0.476	0.585
The side view	0.463	0.566
The top view	0.535	0.603
The total view	0.491	0.585

621

622 4 Discussion

623 In this study, we proposed a weed phenotype segmentation method based on
 624 Mask R-CNN to obtain the growth stage and central area of weeds in a complex
 625 environment. In this study, images of *Solanum nigrum*, Barnyard grass, and *Abutilon*
 626 *theophrasti* Medicus were collected from three angles of front, side, and top view, and
 627 two datasets (data enhancement and without data enhancement) were produced. Two
 628 backbone networks were used, namely, ResNet50 and ResNet101. The results of
 629 detection and instance segmentation are evaluated. Evaluation of this method can
 630 detect weeds at the object level and segment weeds at the pixel level. It is concluded
 631 that data enhancement can improve the performance of the model. The Mask R-CNN
 632 model using ResNet101 as the backbone network has the highest mIOU value, which
 633 is better than ResNet50 and can be used for weed segmentation. Furthermore, the

634 weed image taken from the top view angle compared to the other two angles has the
635 highest detection accuracy.

636 We identified individual plant images of three weeds in the field, but weeds are
637 visual objects with complex structures and rich texture features, and even the same
638 species can have large differences in morphology and colour. However, the deep
639 learning model can automatically learn and extract the features of complex objects,
640 and instance segmentation is a kind of deep learning model that can detect the target
641 pixel by pixel, which solves the problems of blade adhesion and occlusion [52]. Good
642 results have been obtained in the study of plant phenotypes [53]. It can be seen that
643 the instance segmentation model is very suitable for dealing with complex
644 environmental problems in the field. However, the deep learning model relies on a
645 large number of datasets to train the network. Due to the limited number of datasets,
646 we use data enhancement methods to expand the dataset from 4000 to 6000. The
647 experimental results (Table 3) show that the use of data enhancement methods can
648 help improve the performance of the model and help reduce the impact of overfitting.
649 In the case of data enhancement, the leaves of Barnyard grass compared with most
650 broadleaf weeds showed the best performance in all classifications; therefore, because
651 Barnyard grass is an annual herb and its leaves are narrow and long, it is distinct in
652 terms of characteristics from the other two weeds. The less appealing detection result
653 for *Abutilon theophrasti* Medicus occurs because the leaves of *Abutilon theophrasti*
654 *Medicus* are elliptical, similar to those of *Solanum nigrum*. Moreover, a blurred image
655 may result in insufficient image information extraction. Hence, misjudgement of the
656 results can easily occur, which is undesirable because it can lead to errors in leaf age
657 identification. ResNet101 showed a higher F_1 score in the central area, because the
658 characteristics of the centre area are more obvious than other categories when the

659 model is classified. In the total test set, when the IOU threshold of ResNet101 is 0.5
660 or 0.7, the mAP value is greater than that of ResNet50. However, the mIOU of
661 ResNet101 was 0.585 in the total test set, indicating that it can meet the needs of
662 instance segmentation. Thus, the Mask R-CNN model using ResNet101 as the
663 backbone can reliably segment weeds.

664 In the research status at home and abroad, we found that in the research of Quan
665 and Feng, when the angle between the camera and the vertical direction is 0° , the
666 detection accuracy is 0.0095 lower than other angles [27]. But our research found that
667 when ResNet101 is used as the backbone network, the average F_1 value of the top
668 view is 0.0269 and is 0.0083 higher than that of the side view and the front view. I
669 think the reason should be that his research is to obtain maize seedlings and weeds in
670 the image. At the shooting Angle of 0° , the maize in the image contains only one plant,
671 and only contains the top view, but 30° and 75° will produce maize seedlings and
672 weeds from different angles, and also contains the front view and top view. For the
673 detection of maize seedlings and weeds, different angles contain different information,
674 and more comprehensive information is conducive to improving the performance of
675 the model. But in my research, the purpose is to obtain the leaf age and central area of
676 the weed, as shown in Figure 2, we can obtain more comprehensive weed phenotype
677 information from the perspective of the top view, especially given that the central area
678 of weeds will be obtained more, while the side view and the front view will see less
679 from the central area. This may be the cause of the different results.

680 At present, the treatment of weeds in the field mostly involves weed
681 classification and detection. Weed classification can determine only the species of
682 weeds, while the specific position coordinates of weeds cannot be obtained; thus, it is
683 impossible to spray the exact target. Weed detection can facilitate drawing the

684 bounding box of weeds, but weeds are irregular in shape and size, causing the
685 machine to be inaccurate relative to the target, which will result in some herbicide
686 falling to the ground and not being absorbed by the weeds, leading to environmental
687 pollution and wasted herbicide. As a kind of deep learning model, instance
688 segmentation can detect the target pixel by pixel, which solves the problems of blade
689 adhesion and occlusion. The leaf age of weeds and the position of the centre area can
690 be obtained more accurately. In the Northeast Plain of China, the main economic
691 crops are maize, soybeans, and wheat, which are susceptible to annual and perennial
692 weeds. On the one hand, controlling annual and perennial weeds can increase crop
693 yields and reduce the likelihood of damage caused by weeds in the second year [54].
694 On the other hand, studying the interaction between plant phenotype and vision
695 through effective phenotypic analysis can obtain information on plant growth and
696 morphological changes.

697 The limitation of this study is that the identification efficiency of this study is
698 low, and further improvement of efficiency is needed before it can be used in
699 engineering practice. The DCNN model of this study was used to segment only three
700 kinds of weeds. If we expand the kinds of weeds, collect and segment the images of
701 field crops, and increase the number of datasets, the model can achieve a higher
702 segmentation accuracy. According to the obtained leaf age of economic crops, it can
703 also provide an important basis for crop fertilization. For some plants, the central area
704 is the pollination area of flowers, and segmentation of this part will provide an
705 important basis for subsequent studies. Future research will focus on evaluating image
706 datasets covering a wider range of weeds and crop varieties. Most of the images used
707 for model testing contained only single-plant weeds, and only a few contained
708 multiple weeds. Mask R-CNN failed to segment weeds near the edges in a few test

709 images containing multiple weeds, but continuous video input will eliminate edge
710 effects when applied in the field.

711 The results of this study show that the combination of weed phenotype and
712 computer vision is very suitable for dealing with complex field conditions such as
713 light changes, leaf occlusion, and mixed leaf age. The models and methods proposed
714 in this study can be applied to the study of many different types of plants. The data of
715 this study were taken from a complex field environment, and previous studies on plant
716 phenotypes were mostly taken in an indoor environment, in which the image
717 background is often very pure and the illumination is very uniform. Studying the field
718 environment can make the model more suitable for practical applications, and the
719 shooting angle of the dataset determines how much we obtain from the target image
720 information, so it is meaningful to study the segmentation results at different shooting
721 angles. At present, there are few studies on plant phenotypes that are specific to weed
722 phenotypes, but weeds of different leaf ages require different doses of herbicides, so it
723 is of great significance to obtain the information of weed leaf ages to reduce the
724 amount of herbicides. To make this research practical, we can deploy the trained
725 model on the mobile platform of the spray system used for weeding in the future,
726 which can promote the development of precision agriculture and intelligent
727 agriculture.

728 **5 Conclusions**

729 In this paper, we proposed a weed phenotype segmentation method based on the
730 improved Mask R-CNN to obtain the weed species, leaf age and central area of weeds.
731 In the field of plant phenotype research, in the context of complex field environments,
732 obtaining weed phenotypes is still a substantial challenge. According to the research
733 status at home and abroad, we can know that leaf age and central area are important

734 phenotypic traits of weeds. We obtained leaf age and central area, which are of great
735 significance for targeted weeding. Improve the performance of the model through data
736 enhancement. In addition, we found that the weed image taken from the top view
737 angle can help to improve the performance of the model. Weed datasets were
738 constructed by data collection from three angles and data enhancement, the dataset
739 can contain weed information of different growth stages, different angles, and
740 different types. The dataset and research results may provide important resources for
741 future plant phenotype research.

742 Because the DCNN has the ability to extract features from complex
743 environments, it is very suitable for addressing complex image problems in the field.
744 The experimental results show that despite the interference of straw and crop leaves in
745 the background of weeds in the field, Mask R-CNN model using Resnet101 as the
746 backbone network still achieves accurate segmentation of weeds. Good segmentation
747 performance has been achieved. We hypothesize that the Mask R-CNN example
748 segmentation model used here is suitable for all weeds and cash crops in the field,
749 such as *Setaria viridis*, *Cirsium setosum*, maize, and wheat, and that the segmentation
750 of different parts of plants will provide more help in the study of plant phenotypes.
751 Future research will focus on evaluating image datasets that cover a wider range of
752 weeds and crop varieties. The identification efficiency of this study is low, so we need
753 improve model efficiency, and deploy the trained model on the mobile platform of the
754 spray system used for weeding in the future. The model combines artificial
755 intelligence technology with agronomic research and applies it to the development of
756 intelligent agriculture.

757

758

759

760

761

762

763

764

765

766 **6 Declarrations**

767 **Availability of data and materials**

768 Given that the data used in this study were acquired by self-collection, the
769 dataset is being further improved. Thus, the dataset is unavailable for the time being.

770 **Acknowledgments**

771 The authors gratefully appreciate the financial support provided by the Overseas
772 Study and Return Foundation of Heilongjiang LC2018019.

773 **Funding**

774 The authors gratefully appreciate the financial support provided by the Overseas
775 Study and Return Foundation of Heilongjiang LC2018019.

776 **Contributions**

777 QLZ and WB wrote this paper, WB FHQ and MSR collected data, MSR
778 developed the codes, YCJ, LHD and JW supervised the project. All authors read and
779 approved the manuscript.

780 **Conflict of interest**

781 The authors declare no conflict of interest.

782 **Ethics approval and consent to participate**

783 Not applicable.

784 **Consent for publication**

785 Not applicable.

786

787 **7 REFERENCES**

- 788 1. Hashim S, Jan A, Fahad S, Ali HH, Mushtaq MN, Laghari KB, Jabran K, Chauhan BS:
 789 **WEED MANAGEMENT AND HERBICIDE RESISTANT WEEDS: A CASE STUDY**
 790 **FROM WHEAT GROWING AREAS OF PAKISTAN.** *Pakistan Journal of Botany* 2019,
 791 **51(5):1761-1767.**
- 792 2. Slaughter DC, Giles DK, Downey D: **Autonomous robotic weed control systems: A review.**
 793 *Comput Electron Agric* 2008, **61(1):63-78.**
- 794 3. Bakhshipour A, Jafari A, Nassiri SM, Zare D: **Weed segmentation using texture features**
 795 **extracted from wavelet sub-images.** *Biosystems Engineering* 2017, **157:1-12.**
- 796 4. Partel V, Kakarla C, Ampatzidis Y: **Development and evaluation of a low-cost and smart**
 797 **technology for precision weed management utilizing artificial intelligence.** *Computers and*
 798 *Electronics in Agriculture* 2019, **157:339-350.**
- 799 5. Gianessi LP, Reigner NP: **The value of herbicides in US crop production.** *Weed Technology*
 800 2007, **21(2):559-566.**
- 801 6. Pretty J, Bharucha ZP: **Integrated Pest Management for Sustainable Intensification of**
 802 **Agriculture in Asia and Africa.** *Insects* 2015, **6(1):152-182.**
- 803 7. Jeranyama P, Ndlovu F, Morgan J: **Plant Physiology Research.** 2011.
- 804 8. Qinghu L, Guoqi C, Yuhua Z, Zhonghua S, Liyao D: **Sensitivities of Leptochloa**
 805 **chinensis, Echinochloa crusgalli and Digitaria sanguinalis at different leaf stages to**
 806 **cyhalofop-butyl and penoxsulam.** *Journal of Nanjing Agricultural University* 2016,
 807 **39(05):771-776.**
- 808 9. Xiu L, Yu-quan D, Jing-bo L, Yun-yun Z, Chen-zhong J: **Sensitivity of Barnyard Grass at**
 809 **Different Leaf Stage to Bispyribac-Sodium and Cyhalofop-Butyl.** *Journal of Weeds* 2017,
 810 **35(03):22-26.**
- 811 10. YiMing P: **RESEARCH ON ALGORITHM OF RICE DISEASES IDENTIFICATION**
 812 **AND LEAFAGE DETECTION BASED ON MACHINE LEARNING.** *Master of*
 813 *Engineering.* Harbin Institute of Technology; 2019.
- 814 11. Jetter R, Schaffer S: **Chemical composition of the Prunus laurocerasus leaf surface.**
 815 **Dynamic changes of the epicuticular wax film during leaf development.** *Plant physiology*
 816 2001, **126(4):1725-1737.**
- 817 12. Zeisler-Diehl V, Muller Y, Schreiber L: **Epicuticular wax on leaf cuticles does not establish**
 818 **the transpiration barrier, which is essentially formed by intracuticular wax.** *Journal of*
 819 *Plant Physiology* 2018, **227:66-74.**
- 820 13. Garcia-Santillan ID, Pajares G: **On-line crop/weed discrimination through the**
 821 **Mahalanobis distance from images in maize fields.** *Biosyst Eng* 2018, **166:28-43.**
- 822 14. Brivot R, Marchant JA: **Segmentation of plants and weeds for a precision crop protection**
 823 **robot using infrared images.** *IEE Proceedings Part K* 1996, **143(2).**
- 824 15. Jeon HY, Tian LF, Zhu HP: **Robust Crop and Weed Segmentation under Uncontrolled**
 825 **Outdoor Illumination.** *Sensors* 2011, **11(6):6270-6283.**
- 826 16. Bossu J, Gee C, Jones G, Truchetet F: **Wavelet transform to discriminate between crop and**
 827 **weed in perspective agronomic images.** *Comput Electron Agric* 2009, **65(1):133-143.**
- 828 17. Shirzadifar A, Bajwa S, Mireei SA, Howatt K, Nowatzki J: **Weed species discrimination**
 829 **based on SIMCA analysis of plant canopy spectral data.** *Biosyst Eng* 2018, **171:143-154.**

- 830 18. Ozdogan M, Yang Y, Allez G, Cervantes C: **Remote Sensing of Irrigated Agriculture: Opportunities and Challenges**. *Remote Sensing* 2010, **2**(9):2274-2304.
- 831
- 832 19. Huete A, Didan K, Miura T, Rodriguez EP, Gao X, Ferreira LG: **Overview of the radiometric and biophysical performance of the MODIS vegetation indices**. *Remote Sensing of Environment* 2002, **83**(1-2):195-213.
- 833
- 834
- 835 20. Minervini M, Fischbach A, Scharr H, Tsafaris SA: **Finely-grained annotated datasets for image-based plant phenotyping**. *Pattern Recognit Lett* 2016, **81**:80-89.
- 836
- 837 21. Bell J, Dee HM: **Leaf segmentation through the classification of edges**. 2019.
- 838 22. Dobrescu A, Giuffrida MV, Tsafaris SA: **Doing More With Less: A Multitask Deep Learning Approach in Plant Phenotyping**. *Frontiers in Plant Science* 2020, **11**.
- 839
- 840 23. Weng Y, Zeng R, Wu C, Wang M, Wang X, Liu Y: **A survey on deep-learning-based plant phenotype research in agriculture**. *Scientia Sinica Vitae* 2019, **49**(6):698-716.
- 841
- 842 24. Kamilaris A, Prenafeta-Boldu FX: **Deep learning in agriculture: A survey**. *Comput Electron Agric* 2018, **147**:70-90.
- 843
- 844 25. Le VNT, Ahderom S, Apopei B, Alameh K: **A novel method for detecting morphologically similar crops and weeds based on the combination of contour masks and filtered Local Binary Pattern operators**. *Gigascience* 2020, **9**(3).
- 845
- 846
- 847 26. Gao JF, French AP, Pound MP, He Y, Pridmore TP, Pieters JG: **Deep convolutional neural networks for image-based Convolvulus sepium detection in sugar beet fields**. *Plant Methods* 2020, **16**(1).
- 848
- 849
- 850 27. Quan LZ, Feng HQ, Li YJ, Wang Q, Zhang CB, Liu JG, Yuan ZY: **Maize seedling detection under different growth stages and complex field environments based on an improved Faster R-CNN**. *Biosystems Engineering* 2019, **184**:1-23.
- 851
- 852
- 853 28. Geetharamani G, Pandian JA: **Identification of plant leaf diseases using a nine-layer deep convolutional neural network**. *Computers & Electrical Engineering* 2019, **76**:323-338.
- 854
- 855 29. Le TT, Lin CY, Piedad E: **Deep learning for noninvasive classification of clustered horticultural crops - A case for banana fruit tiers**. *Postharvest Biology and Technology* 2019, **156**.
- 856
- 857
- 858 30. Mccool CS, Perez T, Upcroft B: **Mixtures of Lightweight Deep Convolutional Neural Networks: applied to agricultural robotics**. *IEEE Robotics & Automation Letters* 2017, **2**(3):1344-1351.
- 859
- 860
- 861 31. Garcia-Garcia A, Orts-Escolano S, Oprea S, Villena-Martinez V, Garcia-Rodriguez J: **A Review on Deep Learning Techniques Applied to Semantic Segmentation**. 2017.
- 862
- 863 32. Shelhamer E, Long J, Darrell T: **Fully Convolutional Networks for Semantic Segmentation**. *Ieee Transactions on Pattern Analysis and Machine Intelligence* 2017, **39**(4):640-651.
- 864
- 865 33. Yu JG, Li YS, Gao CX, Gao HX, Xia GS, Yu ZL, Li YQ: **Exemplar-Based Recursive Instance Segmentation With Application to Plant Image Analysis**. *Ieee Transactions on Image Processing* 2020, **29**:389-404.
- 866
- 867
- 868 34. Huang SP, Wu SH, Sun C, Ma X, Jiang Y, Qi L: **Deep localization model for intra-row crop detection in paddy field**. *Comput Electron Agric* 2020, **169**.
- 869
- 870 35. He K, Gkioxari G, Dollár P, Girshick R: **Mask R-CNN**. In: *2017 IEEE International Conference on Computer Vision (ICCV): 2017*; 2017.
- 871
- 872 36. He KM, Gkioxari G, Dollár P, Girshick R: **Mask R-CNN**. *IEEE Trans Pattern Anal Mach Intell* 2020, **42**(2):386-397.
- 873

- 874 37. Jia WK, Tian YY, Luo R, Zhang ZH, Lian J, Zheng YJ: **Detection and segmentation of**
875 **overlapped fruits based on optimized mask R-CNN application in apple harvesting robot.**
876 *Comput Electron Agric* 2020, **172**.
- 877 38. Yu Y, Zhang KL, Yang L, Zhang DX: **Fruit detection for strawberry harvesting robot in**
878 **non-structural environment based on Mask-RCNN.** *Computers and Electronics in*
879 *Agriculture* 2019, **163**:9.
- 880 39. Garcia-Santillan ID, Montalvo M, Guerrero JM, Pajares G: **Automatic detection of curved**
881 **and straight crop rows from images in maize fields.** *Biosystems Engineering* 2017,
882 **156**:61-79.
- 883 40. Ho D, Tong M, Ienco D, Gaetano R, Maurel P: **Deep Recurrent Neural Networks for**
884 **mapping winter vegetation quality coverage via multi-temporal SAR Sentinel-1.** *IEEE*
885 *Geoscience & Remote Sensing Letters* 2017, **PP(99)**:1-5.
- 886 41. Ienco D, Gaetano R, Dupaquier C, Maurel P: **Land Cover Classification via Multitemporal**
887 **Spatial Data by Deep Recurrent Neural Networks.** *Ieee Geoscience and Remote Sensing*
888 *Letters* 2017, **14(10)**:1685-1689.
- 889 42. Lee SH, Chan CS, Wilkin P, Remagnino P: **Deep-plant: Plant identification with**
890 **convolutional neural networks.** 2015.
- 891 43. Krizhevsky A, Sutskever I, Hinton GE: **ImageNet Classification with Deep Convolutional**
892 **Neural Networks.** In: *International Conference on Neural Information Processing Systems:*
893 *2012*; 2012.
- 894 44. Ma J, Li Y, Chen Y, Du K, Zheng F, Zhang L, Sun Z: **Estimating above ground biomass of**
895 **winter wheat at early growth stages using digital images and deep convolutional neural**
896 **network.** *European Journal of Agronomy* 2019, **103**:117-129.
- 897 45. Ren SQ, He KM, Girshick R, Sun J: **Faster R-CNN: Towards Real-Time Object Detection**
898 **with Region Proposal Networks.** *Ieee Transactions on Pattern Analysis and Machine*
899 *Intelligence* 2017, **39(6)**:1137-1149.
- 900 46. Wu YT, Tang SM, Zhang SW, Ogai H: **An Enhanced Feature Pyramid Object Detection**
901 **Network for Autonomous Driving.** *Applied Sciences-Basel* 2019, **9(20)**.
- 902 47. Gonzalez S, Arellano C, Tapia JE: **Deepblueberry: Quantification of Blueberries in the**
903 **Wild Using Instance Segmentation.** *Ieee Access* 2019, **7**:105776-105788.
- 904 48. Lin TY, Maire M, Belongie S, Hays J, Perona P, Ramanan D, Dollár P, Zitnick CL: **Microsoft**
905 **COCO: Common Objects in Context.** 2014.
- 906 49. Hripcsak G, Rothschild AS: **Agreement, the F-measure, and reliability in information**
907 **retrieval.** *Journal of the American Medical Informatics Association* 2005, **12(3)**:296-298.
- 908 50. Al-Najjar HAH, Kalantar B, Pradhan B, Saeidi V, Halin AA, Ueda N, Mansor S: **Land Cover**
909 **Classification from fused DSM and UAV Images Using Convolutional Neural Networks.**
910 *Remote Sensing* 2019, **11(12)**.
- 911 51. Zhang E, Yi Z: **Average Precision**; 2016.
- 912 52. Yu J-G, Li Y, Gao C, Gao H, Xia G-S, Yub ZL, Lic Y: **Exemplar-Based Recursive Instance**
913 **Segmentation With Application to Plant Image Analysis.** *IEEE transactions on image*
914 *processing : a publication of the IEEE Signal Processing Society* 2019.
- 915 53. Ubbens JR, Stavness I: **Deep Plant Phenomics: A Deep Learning Platform for Complex**
916 **Plant Phenotyping Tasks.** *Frontiers in Plant Science* 2017, **8**.
- 917 54. Lehoczky E, Nagy P, Lencse T, Toth V, Kismanyoky A: **Investigation of the Damage Caused**

918 **by Weeds Competing with Maize for Nutrients.** *Communications in Soil Science and Plant*
919 *Analysis* 2009, **40**(1-6):879-888.
920
921