

Beet Seedling and Weed Recognition Based on Convolutional Neural Network and Multi-modality Images

Jun Sun (✉ sun2000jun@sina.com)

Jiangsu university <https://orcid.org/0000-0003-3019-6086>

Yu Yang

Jiangsu University

Xiaofei He

JSU: Jiangsu University

Long Wang

Jiangsu University

Xiaohong Wu

Jiangsu University

Jifeng Shen

Jiangsu University

Research

Keywords: beets and weeds, deep learning, deformable convolution, multi-modality images, object detection

Posted Date: May 25th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-520181/v1>

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

[Read Full License](#)

Version of Record: A version of this preprint was published at Multimedia Tools and Applications on December 14th, 2021. See the published version at <https://doi.org/10.1007/s11042-021-11764-5>.

1 Beet seedling and weed recognition based on convolutional neural network and multi-
2 modality images

3 Jun Sun¹, Yu Yang¹, Xiaofei He², Long Wang³, Xiaohong Wu⁴, Jifeng Shen⁵

4 School of electrical and Information Engineering, Jiangsu University, Zhenjiang, China

5 Correspondence

6 Jun Sun, School of electrical and Information Engineering, Jiangsu University, Zhenjiang
7 212000, China.

8 Email: sun2000jun@sina.com

9 Tel: 13775544650

10 **Abstract**

11 **Background:** Difficulties in the recognition of beet seedlings and weeds can arise from
12 a complex background in the natural environment and a lack of light at night. In the
13 current study, a novel depth fusion algorithm was proposed based on visible and near-
14 infrared imagery.

15 **Results:** Visible (RGB) and near-infrared images were superimposed at the pixel-level
16 via a depth fusion algorithm and were subsequently fused into three-channel multi-
17 modality images in order to characterize the edge details of beets and weeds. Moreover,
18 an improved region-based fully convolutional network (R-FCN) model was applied in
19 order to overcome the geometric modeling restriction of traditional convolutional kernels.
20 More specifically, for the convolutional feature extraction layers, deformable convolution
21 was adopted to replace the traditional convolutional kernel, allowing for the entire
22 network to extract more precise features. In addition, online hard example mining was

23 introduced to excavate the hard negative samples in the detection process for the
24 retraining of misidentified samples. A total of four models were established via the
25 aforementioned improved methods. Results demonstrate that the average precision of the
26 improved optimal model for beets and weeds were 84.8% and 93.2%, respectively, while
27 the mean average precision was improved to 89.0%.

28 **Conclusion:** Compared with the classical R-FCN model, the performance of the
29 optimal model was not only greatly improved, but the parameters were also not
30 significantly expanded. Our study can provide a theoretical basis for the subsequent
31 development of intelligent weed control robots.

32 **Keywords**

33 beets and weeds, deep learning, deformable convolution, multi-modality images, object
34 detection

35

36 **Background**

37 The presence of weeds in the field can cause great damage to crop seedlings. More
38 specifically, weeds compete with crops for sunlight and nutrients, thus seriously affecting
39 the photosynthesis of seedlings and increasing the spread of diseases and insect pests.

40 Therefore, the removal of weeds is of great significance in order to maintain crop yield¹.

41 Generally, the weeding methods based on machine vision are performed under good
42 daylight conditions. However, in order to improve production efficiency and reduce the
43 damage of weeds to crop seedlings, continuous operation at night is required to identify

44 weeds. Traditional weed elimination methods are time-consuming and laborious, and
45 mainly depends on artificial excavation or pesticide spraying. Moreover, it is difficult to
46 identify weeds that are similar to crops². Pesticide residues produced by spraying not only
47 poses a great threat to human health, but can also damage the ecological environment³.
48 With the increasing application of precision agriculture, it is particularly important to
49 establish the real-time accurate identification of weeds and crops for the rational
50 application of pesticides, crop yield increases, the reduction of environmental pollution,
51 and the implementation of intelligent weeding⁴⁻⁶.

52 At present, most weed detection methods are developed based on machine vision. A
53 classification method for weed classification based on a back propagation (BP) neural
54 network was proposed⁷. Following the fuzzy classification of the features, a genetic
55 algorithm was used to optimize the network for the identification of weeds. A method to
56 identify weeds based on machine vision during the maize seedling stage was designed⁸.
57 After distortion correction, HSI (Hue, Saturation and Intensity) color space conversion
58 and threshold segmentation, the collected images of maize plants and weeds were
59 identified according to the shape and color features. A weed segmentation network based
60 on an artificial neural network was designed⁹. The single-stage wavelet transform was
61 used to extract weed texture features with 14 texture features selected to optimize the
62 algorithm by using the principal component analysis method. Finally, the features were
63 sent into the neural network in order to identify the weeds. The aforementioned studies
64 combine shallow feature extraction and pattern recognition to identify weeds. However,

65 the feature extraction of such methods is time consuming and the applicability is weak.
66 In addition, due to the influence of the complex field background, the weed characteristics
67 extracted by humans can be ambiguous and uncertain, which consequently results in weed
68 identification based on traditional machine vision at low accuracies.

69 Traditional machine vision methods can only generate low-level image features, which
70 limits the expression ability. However, the convolutional neural network (CNN) has a
71 strong capability to represent image feature without manual feature selection. It can
72 further extract high-level abstract features in images with a high accuracy and efficiency
73 based on the gradient descent and back propagation algorithm¹⁰. Therefore, CNNs are
74 widely used in image classification¹¹, object detection¹² and semantic segmentation¹³. The
75 neurons in convolutional layers can extract the primary visual features of the image by
76 local receptive fields, and reduce the number of parameters through weight sharing. The
77 pooling layer can realize the invariance of displacement, scaling and distortion, while
78 simultaneously performing feature dimensionality reduction.

79 Recently, many networks based on pre-trained CNNs have achieved promising results
80 in weed and crop seedling detection. A weed identification method based on a deep CNN
81 and hash code was proposed¹⁴, which was able to effectively compress the high-
82 dimensional features of the weed through a binary hash layer to detect weeds. Andrea
83 used a classification method based on a CNN that identified maize seedlings and weeds
84 by optimizing the number of convolutional kernels on the basis of the original
85 classification network¹⁵. Joseph Proposed a UAV for weed detection and fixed-point

86 removal, which limited the use of pesticides and improved the efficiency of weed
87 removal¹⁶. Results from the aforementioned literature demonstrate that CNNs can not
88 only automatically extract the shallow features (texture, color, etc.) of weeds and crops,
89 but can also learn deeper abstract features. Moreover, CNNs are able to reduce the cost
90 of feature extraction and are more robust for weed detection in a complex environment.
91 Therefore, CNNs have the potential to be applied to detect beet seedlings and weeds.
92 However, the models in the current literature are not sensitive to the feature information
93 of weeds due to the weak light at night and the use of traditional convolutional kernels.
94 This has resulted in feature extraction difficulties and low recognition accuracies.

95 At present, the computer vision-based detection of crops and weeds in the field is
96 limited. This is a result of the noise commonly present in images of certain scenes (e.g.
97 at night with insufficient light). Thus, the detection process is easily affected changes in
98 light, thus misinterpreting the shape, color, texture and additional information of the
99 detected object, resulting in a poor visual performance. Moreover, traditional CNNs only
100 use the original regular convolutional kernel to extract the features of crops and weeds,
101 with a limited ability in geometric modeling. This may weaken the feature extraction
102 ability of the model. Furthermore, during the training of the model, the uneven
103 distribution of positive and negative samples will also lead to the poor generalization of
104 the model. In view of the above limitations of the current methods, three improved
105 methods were proposed in this study:

106 (1) A deep fusion algorithm was adopted to fuse the RGB (Red, Green and Blue) and

107 near-infrared (NIR) images of beets and weeds into three-channel multi-modality
108 images, then the fusion images were sent into CNN for training.

109 (2) In the feature extraction layer, the traditional convolutional kernel was replaced by
110 the deformable convolution.

111 (3) The hard negative samples were fully excavated by using online hard example mining
112 (OHEM) in the detection process, and then they were sent into the network for re-
113 identification.

114 **Results**

115 In the current study, an improved R-FCN model was proposed to detect and identify beet
116 seedlings and weeds under poor light (at night) and complex field backgrounds. Based on
117 the classic R-FCN network, the visible and near-infrared images of beets and weeds were
118 fused into three-channel multi-modality images at pixel-level using a deep fusion
119 algorithm. The fusion images were then sent to the convolutional neural network for
120 training, so as to improve the mean average precision of beets and weeds. Furthermore,
121 considering that the traditional convolutional kernel would restrict the geometric
122 modeling ability of the model, deformable convolution was adopted in the feature
123 extraction layer. Moreover, online hard example mining was introduced to excavate the
124 hard negative samples in the detection process to retrain the misidentified samples.
125 Through the aforementioned improvement measures, the average precision of beets and
126 weeds were 84.8% and 93.2% respectively, and the mean average precision was increased
127 from 82.3% to 89.0%. Compared with the original model, the detection accuracy was

128 improved by approximately 7 percentage points. Our results demonstrate that the
129 performance of the optimized model is not only greatly improved, but the quantity of the
130 model parameters is also maintained at a reasonable amount, provide a theoretical basis
131 for the subsequent development of intelligent weed control robots.

132

133 *Performance evaluation of the model*

134 Precision and recall are widely used in the field of information retrieval. As with all
135 machine learning problems, in order to calculate precision and recall, it is necessary to
136 explain the following: True Positive (TP), the number of positive classes predicted as
137 positive classes; False Positive (FP), the number of negative classes predicted as positive
138 classes (error rate); True Negative (TN), the number of negative classes predicted as
139 negative classes; and False Negative (FN), the number of positive classes predicted as
140 negative classes (missing rate). Based on this, we define the following:

141

$$Precision =$$

$$142 \quad TP/(TP + FP) \quad (3)$$

143

$$Recall = TP/(TP + FN)$$

$$144 \quad (4)$$

145 The average precision (AP) is then calculated as follows:

146

$$AP = \int_0^1 p(r) d(r)$$

147

$$(5)$$

148 where p represents *Precision*, r represents *Recall*, and p is a function taking r

170 detection accuracy of the model greatly improved with the use of the fusion data set, and
171 the detection system was also more robust.

172

173 **Table 1. The APs and mAP for different types of images**

174

175 The detection results were visualized to further demonstrate the performance of multi-
176 modality images. As can be seen from figure 7, image triplet (a) demonstrates the
177 detection result of the improved R-FCN model on the RGB data set, while image triplet
178 (b) depicts the result of the same model using the multi-modality data set. Due to the
179 fusion of the near-infrared and visible images for the multi-modality image, the feature
180 information of beet seedlings and weeds were able to be better characterized under the
181 terrible light and complex field backgrounds. Thus, the detection performance of the
182 multi-modality fusion images was better than that of the improved R-FCN model. Hence,
183 the subsequent ablation experiments were conducted using on the fusion data set.

184

185 **Figure 1. Testing results using different types of the images**

186

187 *The impact of deformable convolution on mAP*

188 Models 2 and 3 used deformable convolutions, while models 0 and 1 implemented
189 traditional convolutions (Table 2). Compared with the traditional convolution model, the
190 deformable convolution model was able to improve the mAP of beets and weeds by 3-4
191 percentage points. This can be attributed to the offset variables of the deformable
192 convolution kernel, allowing for the feature expression of the CNN to automatically adapt

193 to changes in the morphological of the target object.

194

195 **Table 2. Model parameter settings and mAP**

196

197 Figure 8 depicts the convolution results for the multi-modality data set. The deformable
198 convolution results (image triplet (b)) exhibited a higher detection accuracy and a lower
199 missing rate than the traditional convolutional model (image triplet (a)). This is a result
200 of the ability of the deformable convolution model to adaptively detect irregular
201 geometric edges of seedlings when detecting smaller target objects.

202

203 **Figure 2. Testing results using deformable convolution and traditional convolution**

204

205

206 *The impact of OHEM on mAP*

207 OHEM retrains the hard samples with large loss values, as such samples may lead to the
208 misclassification of weeds and beets. Compared to other models without OHEM, the
209 detection accuracy of models 1 and 3 with the OHEM algorithm were improved (Table
210 2). This indicates that the OHEM method can suppress the simple samples and the small
211 number of samples, making the training process more efficient. In addition, the OHEM
212 algorithm also eliminated several heuristics and hyper-parameters by automatically
213 selecting hard examples, thus simplifying the training process²⁷.

214

215 *Detection results of the optimal model*

216 In order to verify the prediction results of the optimal model under the actual field
217 environment, six images of beets and weeds were selected from the test set. As
218 demonstrated in figure 9, with the use of deformable convolution and multi-modality
219 fusion images, our improved model was able to maintain a high detection accuracy. This
220 indicated that the optimal proposed model exhibits strong generalization and robustness
221 abilities for the detection of beet seedlings and weeds under the poor light and complex
222 field backgrounds.

223

224 **Figure 3. Detection results of the optimal model**

225

226 **Discussion**

227 In modern agriculture, efficient detection and identification of weeds and crops is an
228 important step to improve agricultural efficiency. As an cutting-edge research method,
229 deep learning is more and more widely used in precision agriculture. In the research of
230 automatic weeding machine, effective image preprocessing is the first step for the
231 machine to complete the recognition of images under different lighting conditions. The
232 high detection accuracy means that it can realize the accurate distinction between weeds
233 and crops, reduce the artificial input and improve the agricultural efficiency.

234 In this study, we successfully achieved the fusion of RGB images and near-infrared
235 images, which confirmed that the accurate detection of beet plants and weeds can still be
236 achieved in the absence of light conditions. The model was further improved by
237 deformable convolution and ohem.

238 **Conclusion**

239 This study proved that, through the fusion of RGB image and near-infrared image, the
240 improved deep learning model R-FCN can achieve better detection results. Aim to get a
241 better detection result, online hard example mining was introduced to excavate the hard
242 negative samples in the detection process for the retraining of misidentified samples.
243 Deformable convolution can also improve the effect of feature extraction. However, there
244 are still a lot of work to be improved, such as how to further improve the accuracy of the
245 model, whether the model can be compressed to the mobile portable robot. Therefore, it
246 is necessary to carry out further research

247 **Materials and methods**

248 *Data source*

249 In order to investigate the detection and identification of beets and weeds in complex
250 backgrounds, images of beets and weeds were collected at the University of Bonn,
251 Germany, in 2016. For more details about the dataset, see the link and study below
252 (<http://www.ipb.unibonn.d-e/data/sugarbeets2016/>). The images were collected via a
253 multi-modality camera (JAI AD-130GE), equipped with two high-sensitivity CCD
254 multispectral sensors of 1.3 million pixels. The camera can simultaneously collect visible
255 (400nm~650nm) and near infrared (NIR) (760nm~1000nm), with an output image size
256 of 1296×1296 pixels¹⁷. The dataset contains a total of 2,093 images of beets and weeds
257 at different growth stages (Figure.1). In the process of data acquisition, beet seedlings and
258 weeds with different levels of maturity under varying angle transformations were consider

259 for the image acquisition. Moreover, the same plant (beet and weed) was imaged multiple
260 times under different ranges of overlap and occlusion between beets and weeds.

261

262 **Figure 4. Examples of seedling dataset for weeds and beets**

263

264 *Multi-modality image fusion*

265 In general, object detection methods based on deep learning aim to understand the
266 distribution of basic data via a large amount of training data and subsequently induce the
267 optimizer to adjust the parameters of the network¹⁸. At present, RGB images on weeds
268 and beets are commonly used to train deep learning models. However, RGB images are
269 sensitive to variations in light, resulting in the loss of important information on the shape,
270 color and texture of target objects¹⁹. Therefore, the performance of such models is poor
271 under complex backgrounds at night.

272 In order to solve this problem, in the current study, the NIR and visible images of beets
273 and weeds were fused into multi-channel images. In particular, multi-modality image
274 fusion spatially registers the data of the same image from different sources, and
275 subsequently combines the information in each image to generate an integrated data set
276 of all the images²⁰.

277

278 *Deep fusion algorithm frame diagram*

279

280 **Figure 5. Deep fusion method for multi-modality images**

281

282 Essentially, the deep fusion algorithm extracts features from a denseblock composed of
283 convolutional layers, and then sends the visible and near-infrared images to a fusion layer
284 for pixel-level superposition. These fusion images are reconstructed through a decoding
285 network also composed of convolutional layers to obtain the three-channel multi-
286 modality images. As shown in figure 2, the encoding network consists of two sections:
287 C1 and denseblock. In the C1 section, a 3×3 convolutional kernel is used to extract the
288 rough features of the image. The denseblock section then uses three convolutional layers
289 (the output of each layer is cascaded as the input of the next layer) to extract the high-
290 level abstract semantic features of the image. The denseblock adopted in this research can
291 retain image features as much as possible to ensure that all significant features can be
292 used in the fusion strategy. In the fusion layer, $l_1 - norm$ fusion strategy is selected to
293 fuse the visible and near-infrared images. Finally, four convolutional layers (3×3
294 convolutional kernel) are used to reconstruct the final fusion image in the decoding
295 network. What's more, in this framework, the size of the input and output images is 1296
296 $\times 1296$ pixels, and the number of feature mapping channels per convolutional layer is 16.
297 More details of the deep fusion algorithm can be found in the study²¹.

298

299 *Multi-modality fusion image*

300 Figure 3 presents the three-channel multi-modality image obtained via the depth fusion
301 algorithm. Following the fusion of the data, the Labellmg software
302 (<https://tzutalin.github.io/Lab-eImg/>) was used by experts in the agricultural field to label

303 the beet seedlings and weeds based on the PASCAL Visual Object Classes Challenge²².
304 A python script was then used to randomly divide the images and the corresponding tag
305 files into training and testing sets at a ratio of 4:1.

306

307 **Figure 6. Examples of RGB and multi-modality images. The image triplet (a)**
308 **shows the original visible dataset, and the image triplet (b) shows the**
309 **corresponding fusion dataset.**

310

311 *Improved weeds and beets detection model*

312 *Deformable Convolution*

313 Recently, the use of CNNs has made significant breakthroughs in many vision
314 applications. However, due to the regular grid sampling and the fixed geometric structure
315 in traditional convolution methods, it is difficult for networks to deal with geometric
316 deformations. The adaptability of the existing models to process the geometric
317 deformation of objects almost comes from the diversity of the data itself, and there is no
318 mechanism to adapt to geometric deformation in the model. Thus, the ability of geometric
319 transformation modeling is limited, and it cannot be adjusted adaptively according to the
320 image content²³. In order to overcome this limitation, a new module, deformable
321 convolution, was adopted to improve the modeling ability of CNN for transformations in
322 the current study. More specifically, an offset variable is added to the position of each
323 sampling point in the convolutional kernels. The kernel with the offset variable can then
324 be sampled randomly near the current position, and is thus not restricted to the previous
325 regular grid points. Moreover, the offset variable can be obtained by learning within the

326 target task without the need of an additional monitoring signal, improving the traditional
327 convolution.

328

329 **Figure 7. Schematic diagram of deformable convolution. This figure shows the**
330 **sampling methods of traditional convolution and deformable convolution with**
331 **convolution kernel size 3×3 . Figure 4(a) demonstrates the regular sampling grid**
332 **(green points) of traditional convolution, while (b) present the deformed sampling**
333 **locations (black points) with augmented offsets (blue arrows) in deformable**
334 **convolution. Then (c) and (d) are special cases of (b), showing that deformable**
335 **convolution generalizes scale, aspect ratio and rotation transformations.**

336

337 **Figure 8. Illustration of 3×3 deformable convolution**

338

339 Figure 5 presents the internal structure of deformable convolution. First, the
340 displacement required for deformable convolution is obtained through the output of a
341 small convolutional layer, and the displacement is then applied to the convolutional kernel
342 in order to achieve the effect of deformable convolution. This operation is able to add the
343 offsets to the regular grid sampling locations in the standard convolution, thus enabling
344 the free form deformation of the sampling grid. The offsets are learned from the preceding
345 feature maps via additional convolutional layers.

346 Traditional convolution consists of two steps: 1) Sampling using a regular grid R over
347 the input feature map x ; and 2) the summation of sampled values weighted by w . For
348 each location p_0 on the output feature map y , traditional convolution is then performed
349 as follows:

$$350 \quad y(p_0) = \sum_{p_n \in R} w(p_n) \cdot x(p_0 + p_n) \quad (1)$$

351 where the grid R defines the receptive field size and dilation of a 3×3 kernel with

352 dilation 1, $R = \{(-1, -1), (-1, 0), \dots, (0, 1), (1, 1)\}$, and p_n enumerates the locations in
353 R .

354 In deformable convolution, the grid R has the offsets $\{\Delta p_n | n = 1, 2, \dots, N\}$, and $N =$
355 $|R|$, thus, formula (1) becomes:

$$356 \quad y(p_0) = \sum_{p_n \in R} w(p_n) \cdot x(p_0 + p_n + \Delta p_n) \quad (2)$$

357 As the offset Δp_n is typically fractional, formula (2) is implemented via bilinear
358 interpolation.

359

360 *The network structure of the improved model*

361

362 **Figure 9. Key idea of the improved R-FCN for weeds and beets detection.**

363

364 The network structure of the improved model is described as follows. Following the input
365 of the image, feature extraction was performed, resulting in a feature map of $k^2 \times (C+1)$
366 dimensions. The region proposal network (RPN)²⁴ was then used to extract the regions of
367 interest (ROIs) in the feature map, with C denoting the number of categories. The
368 extracted ROIs were divided into $k \times k$ regions, with k generally equal to 3, corresponding
369 to 9 regions: top-left, top-center, and bottom-right. Finally, the score of each region was
370 determined by a pooling operation, and the output feature vector of the ROIs were then
371 obtained by voting. This vector was subsequently used for classification and regression
372 of the weeds and beets.

373

374 *Experimental environment*

375 *Equipment and platform*

376 The Ubuntu 16.04 system is used as the operating platform, and MXNet is adopted as the
377 deep learning framework to train the network. The computer memory of the system is
378 32GB, with a 3.6 GHZ i7-9700k CPU processor. Additionally, the 11 GB GeForce
379 GTX1080Ti GPU with Pascal architecture was used.

380

381 *Model parameter setting*

382 In order to reduce variation of parameter updates and to stabilize the convergence of the
383 model, a mini-batch stochastic gradient descent (SGD) was used to train the network²⁵.

384 The parameters were set as follows: the number of each mini-batch of samples was 128,
385 the momentum factor was fixed to 0.9, and the weight attenuation factor was 0.0005 to
386 avoid over-fitting. Finally, the gradient descending learning rate was applied to all layers
387 of the network, and was gradually reduced in stages to 0.1 times of the current learning
388 rate. Additionally, the initial learning rate during the training process was set to 0.005,
389 and the model was iterated for 100 epochs.

390

391 **Acknowledgements**

392 We are very grateful to Professor Sun Jun, Wu Xiaohong and Shen Jifeng for their
393 suggestions to the manuscript.

394 **Authors' contributions**

395 Yang Yu and He Xiaofei conceived and designed research. Yang Yu conducted
396 experiments. All authors read and approved the final manuscript.

397 **Funding**

398 This work is partially supported by a project funded by the Priority Academic Program
399 Development of Jiangsu Higher Education Institutions (PAPD-2018-87), Synergistic
400 Innovation Center of Jiangsu Modern Agricultural Equipment and Technology
401 (4091600002). Project of Faculty of Agricultural Equipment of Jiangsu University
402 (4121680001)

403 **Availability of data and materials**

404 The datasets used and/or analyzed during the current study are available from the
405 corresponding authors on reasonable request.

406 **Declarations**

407 Ethics approval and consent to participate

408 Not applicable.

409 Consent for publication

410 Not applicable.

411 Competing interests

412 The authors declare that they have no competing interests.

413

414 **References**

415 1. Bauer M V, Marx, C, Bauer F V, Flury D M, Ripken T, Streit B. Thermal weed control

- 416 technologies for conservation agriculture—a review. *Weed Res.*2020;00:1–10.
417 [https://doi.org/ 10.1111/wre.12418](https://doi.org/10.1111/wre.12418)
- 418 2. Hu K, Coleman G, Zeng S, Wang Z, Walsh M. Graph weeds net: a graph-based deep
419 learning method for weed recognition. *Comput Electron Agric* 2020;174, 105520.
420 [https://doi.org/ 10.1016/j.compag.2020.105520](https://doi.org/10.1016/j.compag.2020.105520)
- 421 3. Saichao Gao, Guobin Wang, Yangyang et al. Water-soluble food dye of allura red as a
422 tracer to determine the spray deposition of pesticide on target crops. *Pest Manag Sci.*
423 2019;75:2592-2597. [https://doi.org/ 10.1111/wre.12418](https://doi.org/10.1111/wre.12418)
- 424 4. Gai J, Tang L, Steward B L. Automated crop plant detection based on the fusion of
425 color and depth images for robotic weed control. *J. Field Robot.* 2019; 37:35-52.
426 [https://doi.org/ 10.1002/rob.21897](https://doi.org/10.1002/rob.21897)
- 427 5. Raja R, Nguyen TT, Slaughter DC, Fennimore SA. Real-time robotic weed knife
428 control system for tomato and lettuce based on geometric appearance of plant labels.
429 *Biosyst Eng.*2020;194:152-164. [https://doi.org/ 10.1016/j.biosystemseng.2020.03.022](https://doi.org/10.1016/j.biosystemseng.2020.03.022)
- 430 6. Sun, J., He, X., Tan, W., Wu, X., Shen, J., Lu, H. Recognition of crop seedling and
431 weed recognition based on dilated convolution and global pooling in CNN. *Trans Chin.*
432 *Soc Agric Eng.* 2018;34(11):159-165. [https://doi.org/ 10.11975/j.issn.1002-](https://doi.org/10.11975/j.issn.1002-6819.2018.1-1.020)
433 [6819.2018.1-1.020](https://doi.org/10.11975/j.issn.1002-6819.2018.1-1.020)
- 434 7. Zhao P, Wei Z. Weed Recognition in Agricultural Field Using Multiple Feature Fusions.
435 *Trans Chin Soc Agric Mach.* 2014;45(3):275-281. [https://doi.org/10.60-41/j.issn.1000-](https://doi.org/10.60-41/j.issn.1000-1298.2014.03.045)
436 [1298.2014.03.045](https://doi.org/10.60-41/j.issn.1000-1298.2014.03.045)

- 437 8. Yan B. Identification of Weeds in Maize Seedling Stage by Machine Vision Technology.
438 J. Agric. Mechanization Res. 2018;40(3): 212-216. <https://doi.org/10.3969/j.issn.1003->
439 188X.2018.03.043
- 440 9. Bakhshipour A, Jafari A, Nassiri S M, Zare D. Weed segmentation using texture
441 features extracted from wavelet sub-images. Biosyst Eng. 2017; 157: 1-12.
442 <https://doi.org/10.1016/j.biosyste-mseng.2017.02.002>
- 443 10. Nogueira K , Penatti O A B, Santos J A D. Towards Better Exploiting Convolutional
444 Neural Networks for Remote Sensing Scene Classification. Pattern Recognition.
445 Pattern Recognit. 2017;61:539-556. <https://doi.org/10.1016/j.patcog.2016.07.001>
- 446 11. Palanisamy T, Sadayan G, Pathinetampadiyan N. Neural network-based leaf
447 classification using machine learning. Concurrency Computat Pract Exper.
448 2019;(15):e5366.
- 449 12. Rehman T U, Zaman Q U, Chang Y K, Schumann A W, Corscadden K W.
450 Development and field evaluation of a machine vision based in-season weed detection
451 system for wild blueberry. Comput Electron Agric. 2019; 162:1-13. [https://doi.org/](https://doi.org/10.1016/j.compag.2019.03.023)
452 10.1016/j.compag.2019.03.023
- 453 13. Sun, J., Tan, W., Wu, X., Shen, J., Lu, B., Dai, C. Real-time recognition of sugar beet
454 and weeds in complex backgrounds using multi-channel depth-wise separable
455 convolution model. Trans Chin Soc Agric Eng. 2019; 35(2).
456 <https://doi.org/10.11975/j.issn.1002-6819.2019.12.000>
- 457 14. Jiang H, Wang P, Zhang Z, Mao W, Zhao B, Qi P. Fast Identification of Field Weeds

- 458 Based on Deep Convolutional Network and Binary Hash Code. *Trans Chin Soc Agric*
459 *Eng.* 2018;49(11):37-45. [https://doi.org/ 10.6041/j.issn.1000-1298.2018.11.004](https://doi.org/10.6041/j.issn.1000-1298.2018.11.004)
- 460 15. Andrea, C. C., Daniel, B. B. M., Misael, J. B. J. Precise weed and maize classification
461 through convolutional neuronal networks. *IEEE Second Ecuador Technical Chapters*
462 *Meeting (ETCM).* 2017; pp. 1-6. <https://doi.org/10.1109/ETCM.2017.8247469>
- 463 16. Joseph E Hunter, Travis W Gannon, Robert J Richardson, Fred H Yelverton, Ramon
464 G Leon. Integration of remote-weed mapping and an autonomous spraying unmanned
465 aerial vehicle for site-specific weed management. *Pest Manag Sci.* 2020; 76: 1386–
466 1392. [https://doi.org/ 10.1002/ps.5651](https://doi.org/10.1002/ps.5651)
- 467 17. Milioto A, Lottes P, Stachniss C. Real-time semantic segmentation of crop and weed
468 for precision agriculture robots leveraging background knowledge in CNNs. *Proc Intl*
469 *Conf on Robotics and Automation.* 2018; 2229-2235. [https://doi.org/1-](https://doi.org/10.1109/ICRA.2018.8460962)
470 [0.1109/ICRA.2018.8460962](https://doi.org/10.1109/ICRA.2018.8460962)
- 471 18. Shin H C, Roth H R, Gao M, Lu L, Xu Z, Noguees I. Deep convolutional neural
472 networks for computer-aided detection: CNN architectures, dataset characteristics and
473 transfer learning. *IEEE transactions on medical imaging.* 2016;35(5):1285-1298.
474 <https://doi.org/10.1109/TM-I.2016.2528162>
- 475 19. Ying, Z., Li, G., Ren, Y., Wang, R., Wang, W. A new image contrast enhancement
476 algorithm using exposure fusion framework. *Int. Conf on Comput Analysis of Images*
477 *and Pattern.* 2017; 36-46.
- 478 20. Ren X, Meng F, Hu T, Liu Z, Wang C. Infrared-Visible Image Fusion Based on

- 479 Convolutional Neural Networks. Int Conf on Intell Sci and Big Data Eng. 2018; 301-
480 307. https://doi.org/10.1007/978-3-030-02698-1_26
- 481 21. Li H, Wu X. DenseFuse: A Fusion Approach to Infrared and Visible Images. IEEE
482 Trans on Image Processing. 2019; 28(5):2614-2623.
483 <https://doi.org/10.1109/TIP.2018.2887342>
- 484 22. Everingham M, Winn J. The PASCAL Visual Object Classes Challenge 2007
485 (VOC2007) Development Kit. Int J of Comput Visi. 2006;111(1):98-136.
486 <https://doi.org/10.1-007/s11263-014-0733-5>
- 487 23. Dai J, Qi H, Xiong Y, Li Y, Zhang G, Hu H. Deformable Convolutional Networks.
488 Proc IEEE int conf on comput visi. 2017; 764-773. [https://doi-](https://doi.org/10.1109/ICCV.2017.89)
489 [i.org/10.1109/ICCV.2017.89](https://doi.org/10.1109/ICCV.2017.89)
- 490 24. Ren, S., He K., Girshick, R., Sun, J. Faster R-CNN: Towards real-time object
491 detection with region proposal networks. In Advances in neural information processing
492 systems. 2015;91-99.
- 493 25. Zhou Y, Xu T, Deng H, Miao T. Real-time recognition of main organs in tomato based
494 on channel wise group convolutional network. Trans Chin Soc Agric Eng.
495 2018;34(10):153-162. <https://doi.org/10.11975/j.issn.1002-6819.2018.10.019>
- 496 26. Raghavendra, R., Dorizzi, B., Rao, A., Kumar, G.H.. Particle swarm optimization
497 based fusion of near infrared and visible images for improved face verification. Pattern
498 Recognit 2011;44(2). <https://doi.org/10.1016/j.patcog.2010.08.006>
- 499 27. Shrivastava A, Gupta A, Girshick R. Training region-based object detectors with

500 online hard example mining. Proc IEEE Conf Comput Vis and Pattern Recognit.

501 2016;761-769. [http:// dx.doi.org/10.1109/CVPR.2016.89](http://dx.doi.org/10.1109/CVPR.2016.89)

502

Figures

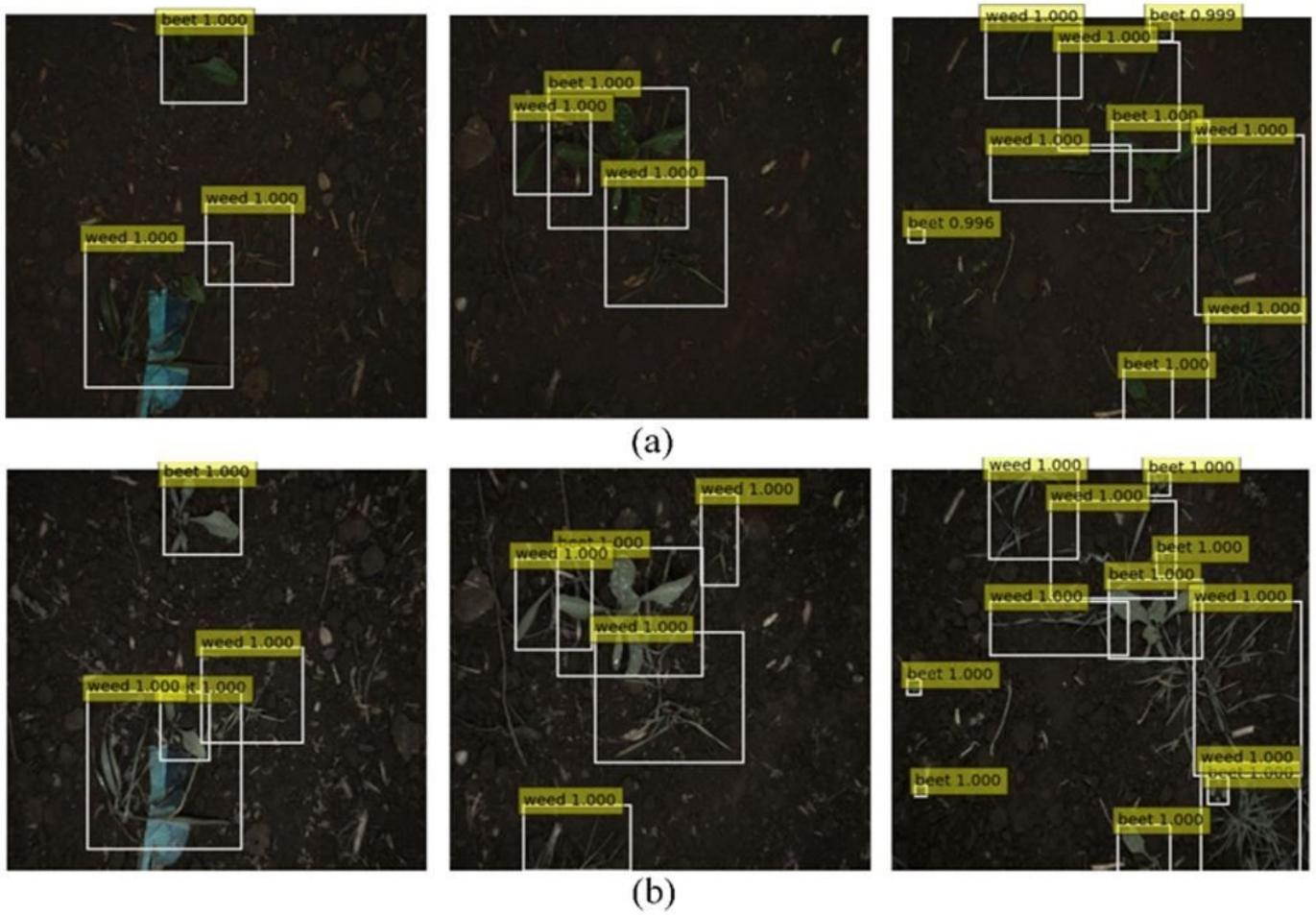


Figure 1

Testing results using different types of the images

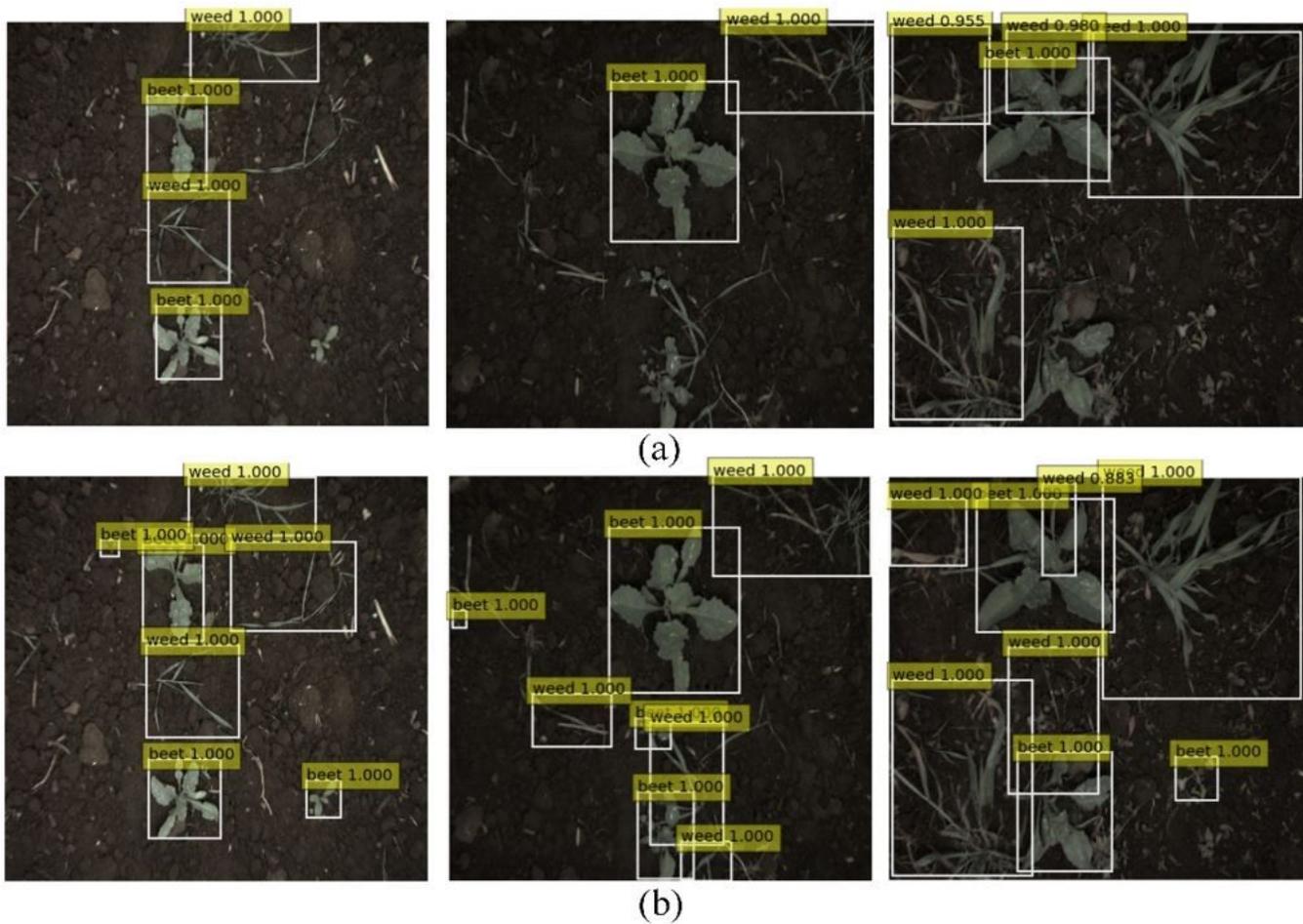


Figure 2

Testing results using deformable convolution and traditional convolution

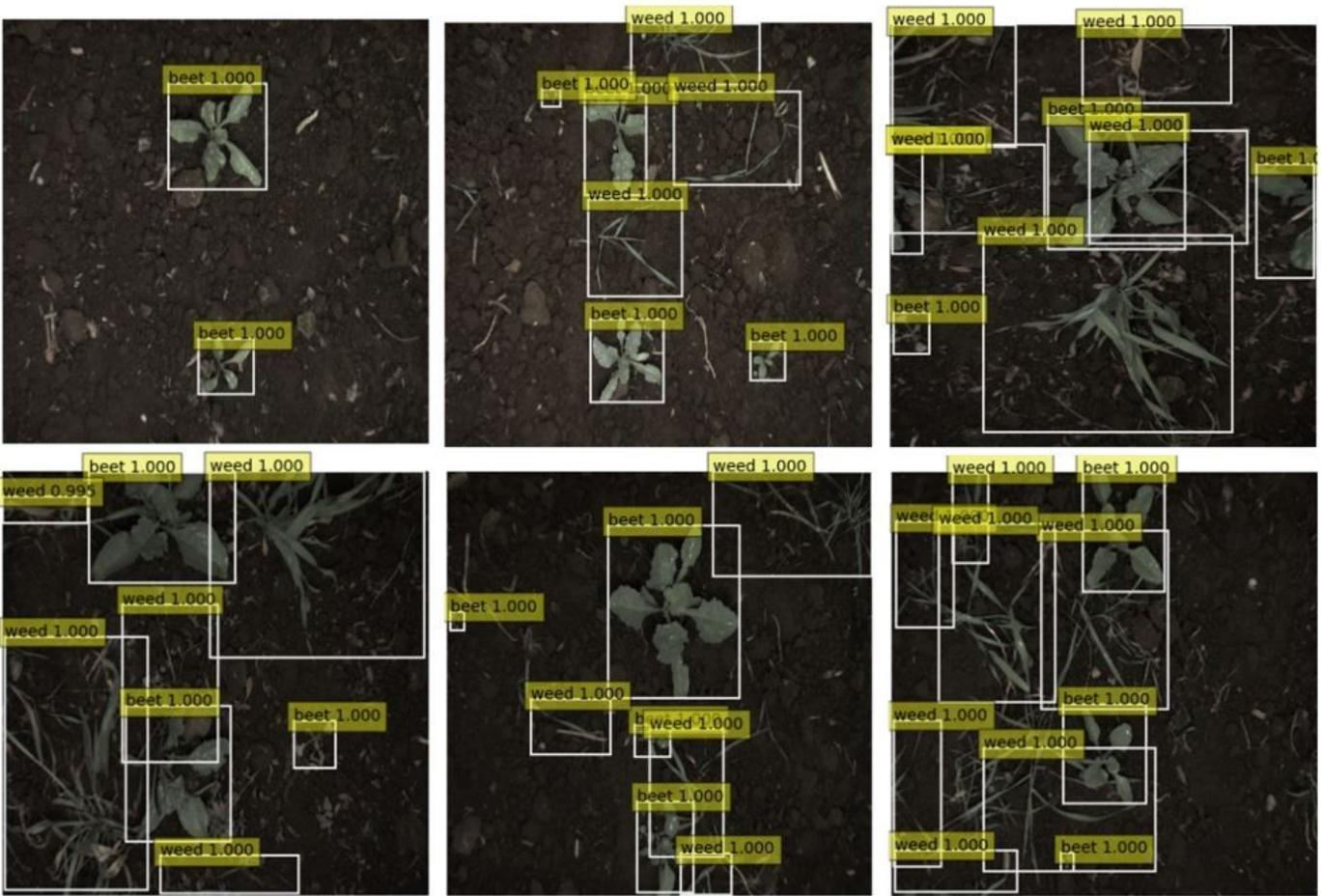


Figure 3

Detection results of the optimal model



Figure 4

Examples of seedling dataset for weeds and beets

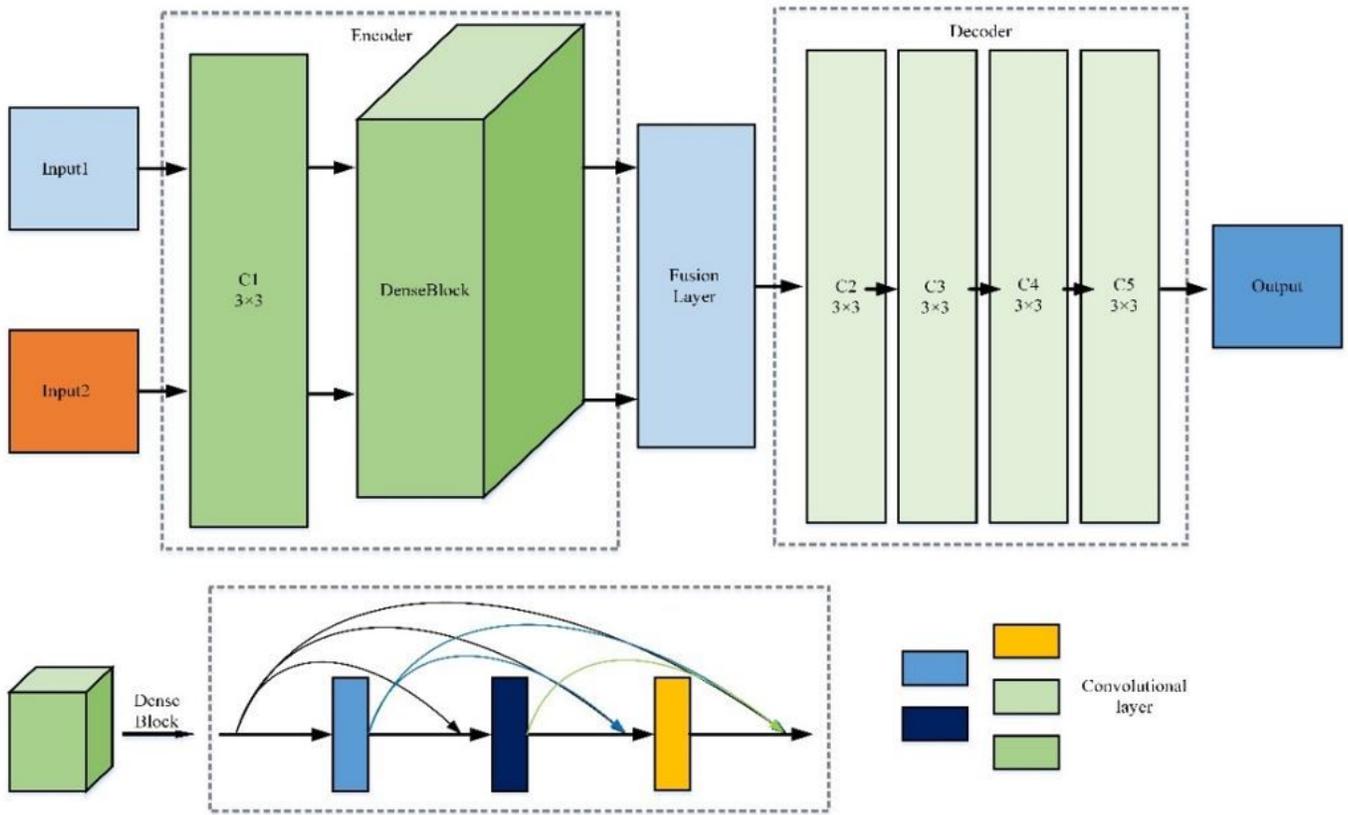


Figure 5

Deep fusion method for multi-modality images



(a)



(b)

Figure 6

Examples of RGB and multi-modality images. The image triplet (a) shows the original visible dataset, and the image triplet (b) shows the corresponding fusion dataset.

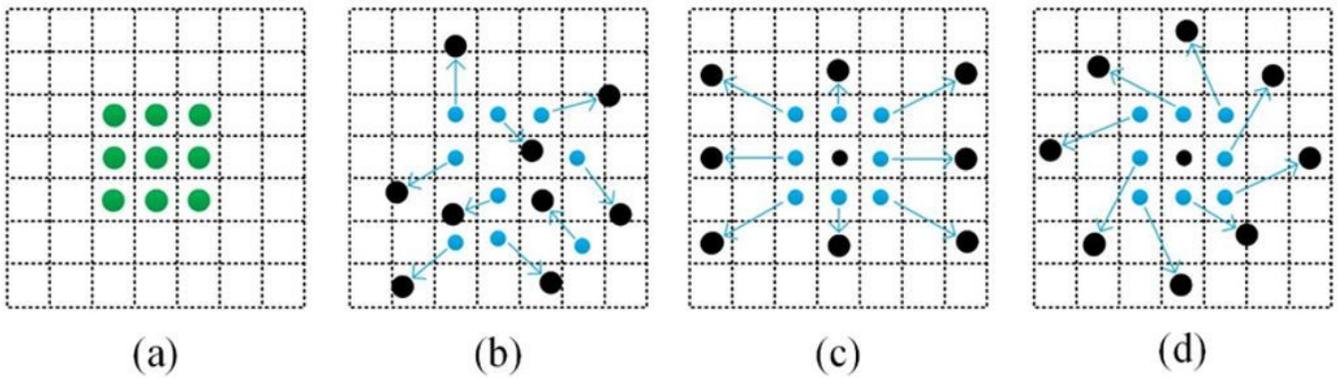


Figure 7

Schematic diagram of deformable convolution. This figure shows the sampling methods of traditional convolution and deformable convolution with convolution kernel size 3×3 . Figure 4(a) demonstrates the regular sampling grid (green points) of traditional convolution, while (b) present the deformed sampling locations (black points) with augmented offsets (blue arrows) in deformable convolution. Then (c) and (d) are special cases of (b), showing that deformable convolution generalizes scale, aspect ratio and rotation transformations.

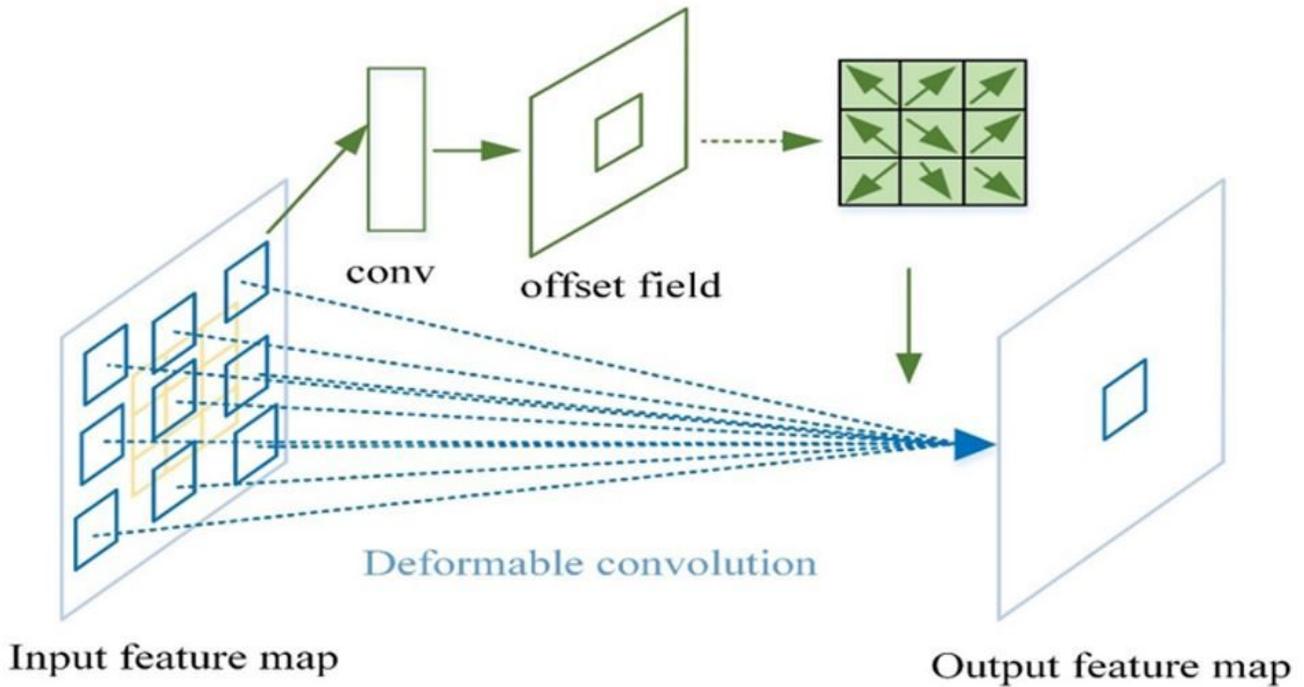


Figure 8

Illustration of 3×3 deformable convolution

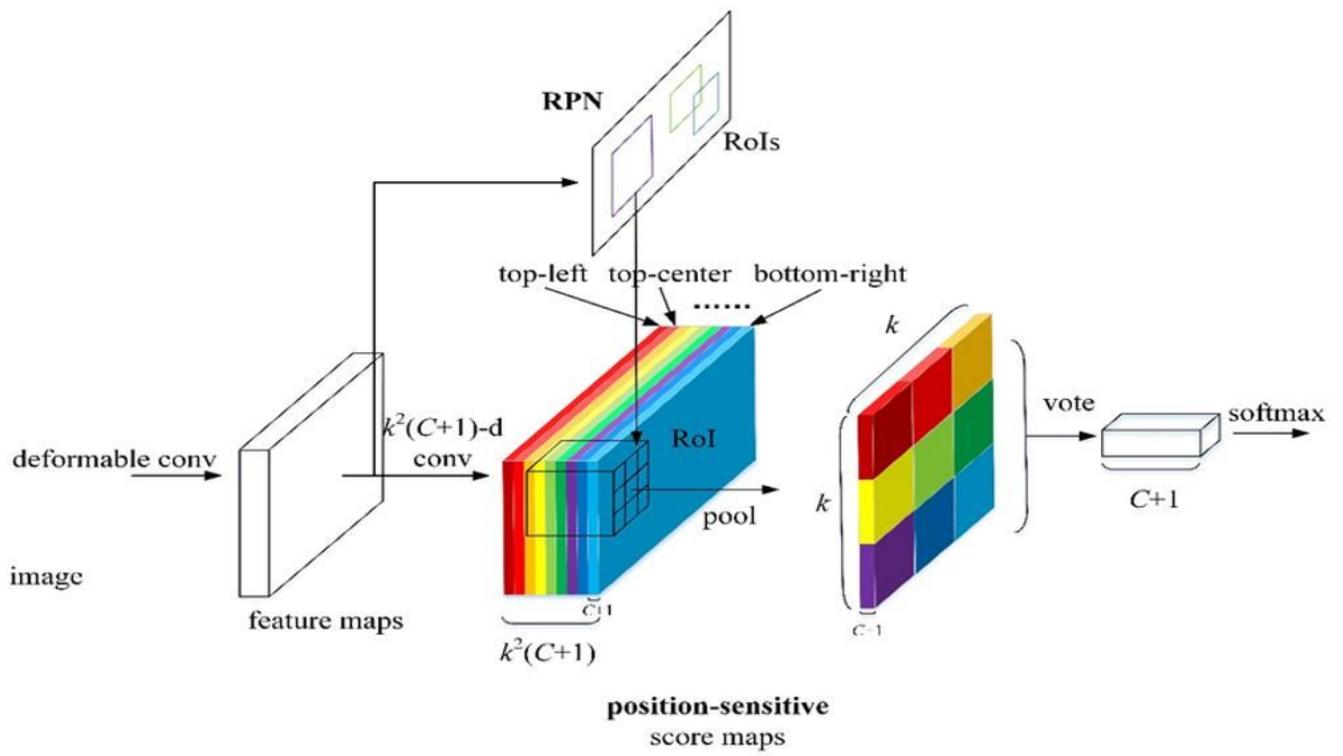


Figure 9

Key idea of the improved R-FCN for weeds and beets detection.