1    Implementation of Deep-learning Algorithm for Obstacle Detection and Collision Avoidance

2                          for Robotic Harvester

3        Li Yang[a], Iida Michihisa[a,*], Suyama Tomoya[a], Suguri Masahiko[a], Masuda Ryohei[a]

4    [a] Graduate School of Agriculture, Kyoto University, Kitashirakawa Oiwake-cho, Sakyo-ku, Kyoto, 606-

5                            8502, Japan

6
7    **Abstract:** Convolutional neural networks (CNNs) are the current state of the art systems in

8    image semantic segmentation (SS). However, because it requires a large computational cost, it

9    is not suitable for running on embedded devices, such as on rice combine harvesters. In order

10   to detect and identify the surrounding environment for a rice combine harvester in real time, a

11   neural network using Network Slimming to reduce the network model size, which takes wide

12   neural networks as the input model, yielding a compact model (hereafter referred to as "pruned

13   model") with comparable accuracy, was applied based on an image cascade network (ICNet).

14   Network Slimming performs channel-level sparsity of convolutional layers in the ICNet by

15   imposing L1 regularization on channel scaling factors with the corresponding batch

16   normalization layer, which removes less informative feature channels in the convolutional

17   layers to obtain a more compact model. Then each of the pruned models were evaluated by

18   mean intersection over union (IoU) on the test set. When the compaction ratio is 80 %, it gives

19   a 97.4 % reduction of model volume size, running 1.33 times faster with comparable accuracy

20   as the original model. The results showed that when the compaction ratio is less than 80 %, a

21   more efficient (less computational cost) model with a slightly reduced accuracy in comparison

22   to the original model was achieved. Field tests were conducted with the pruned model (80 %

23   compaction ratio) to verify the performance of obstacle detection. Results showed that the

24   average success rate of collision avoidance was 96.6% at an average processing speed of 32.2

* Corresponding author. Tel.: +81-75-753-6166. E-mail address: iida@elam.kais.kyoto-u.ac.jp.

25 FPS (31.1 ms per frame) with an image size of 640 ×480 pixels on a Jetson Xavier. It shows

26 that the pruned model can be used for obstacle detection and collision avoidance in robotic

27 harvesters.

28 **Keywords:** robotic combine harvester; deep learning; human detection; image cascade

29 network; network slimming.

30

## 1 Introduction

32     To ensure the safety and precision operation of autonomous combine harvesters it is

33 important to identify obstacles quickly and accurately in the surrounding paddy. When a

34 combine is working in a paddy, it should avoid colliding with paddy field ridges and humans,

35 and it should also go along the navigation line between harvested and unharvested areas. Our

36 laboratory has developed algorithms to determine the path between harvested and unharvested

37 areas (Cho et al., 2014a; Cho et al., 2014b), the identification of ridges (Takagaki et al., 2013),

38 and the detection of humans in the field (Hisae et al., 2017). However, the paddy field

39 environment is complex, and with many different objects (the harvested area, unharvested area,

40 ridges, and humans) need to be detected simultaneously. Traditional image recognition

41 methods are based on hand-crafted features, such as HOG, LBP and Haar features (Yao et al,

42 2015, Singh, et al., 2015, Cabrera et al., 2011). Since it is tedious to design features manually

43 and susceptible to the effects of light, vibration, and dust, it is difficult to mine deep-feature

44 information and obtain accurate results. One approach to addresses these challenges is semantic

45 segmentation (SS), which can realize pixel-by-pixel identification in an image.

46     Recently, SS has become a popular approach for a variety of computer vision tasks in

47 agriculture. For example, Yang et al. (2017) employed a SS method to recognize lactating sows.

48 Milioto et al. (2018) proposed a SS model for crop and weed. McCool et al. (2017) proposed

49 an approach for training SS that can be used to derive compact models for robotic platforms.

50   These research results indicate that SS can be used to achieve good results for processing

51   agricultural images (Kamilaris et al., 2018), thereby reducing manual preprocessing and

52   subsequent processing to obtain the final segmentation result directly from the original input

53   image (Tang et al., 2016). However, the large computational cost of SS models still makes it

54   difficult to apply to embedded devices in real-time. Our objective is to make a SS model

55   compact to implement for embedded devices and apply it for obstacle detection of robotic

56   combine harvester.

57   To achieve this objective, on the one hand, many scholars have proposed different real-

58   time SS models. For example, Yu (Yu et al., 2018) proposed a bilateral segmentation network,

59   which used affluent spatial details and large receptive field to improve the speed and accuracy

60   of SS. Wang (Wang et al., 2019) designed an asymmetric encoder-decoder architecture for SS.

61   Zhao (Zhao et al., 2018) proposed ICNet, which uses an image cascade to speed up the SS

62   method. On the other hand, many methods to compress large CNNs have been developed for

63   fast inference. These include low-rank approximation (Denton et al., 2014), network

64   quantization (Chen et al., 2015; He et al., 2015) and binarization (Rastegari et al., 2016;

65   Courbariaux et al., 2016), weight pruning (Han et al., 2015), dynamic inference (Huang et al.,

66   2017), etc. Network Slimming is a simple yet effective compaction approach (Liu et al., 2017),

67   and more importantly, it is convenient to obtain the pruned model just by modifying the number

68   of corresponding channels in the configuration files.

69   Considering the speed and accuracy in the CamVid (Atlas, 2018), the network used for

70   rice field images were based on ICNet. This method incorporates effective strategies to

71   accelerate network inference speed without sacrificing much performance (Zhao et al., 2018).

72   In this study, a ICNet that maintains a high accuracy was trained first with paddy field images.

73   Paddy field image are, however, not common in public data sets, such as the CamVid Dataset

74   (Brostow et al., 2009). When the ICNet that performs well with public datasets is applied to

75  paddy field images, high segmentation accuracy was obtained. Since the network was designed

76  manually, the importance of each component in the network cannot be determined before

77  training. During training, it could learn the importance of each component through adjusting

78  the weights in trainable layers automatically. After training, some connections and

79  computations in the model would become redundant or non-critical (Ye et al., 2018).

80  Consequently, the redundant or non-critical connections and computations in the network can

81  be removed without significant degradation in performance (Ye et al., 2018). Based on this

82  assumption, we removed these redundant parameters in the model while ensuring similar

83  accuracy, thereby increasing the speed of the model.

84  Since Network Slimming method is a simple yet effective compaction approach (Liu et

85  al., 2017), the pruned SS models were obtained based on this method in the convolutional

86  layers of ICNet. To this end, we enforced channel-level sparsity of convolutional layers by

87  imposing L1 regularization on channel scaling factors $\gamma$ in batch normalization (BN) layer (the

88  latter in formula (3)), then removed the less informative channels in the convolutional layers,

89  which correspond to the small $\gamma$ to obtain the pruned models. The models and methods were

90  introduced firstly in Section 2; then the pruned models were evaluated on test dataset and in

91  the field. Then the results and discussion were presented in Section 4. Finally, we made a

92  conclusion in Section 5.

93  **2 Materials and Methods**

94  **2.1 Semantic segmentation model**

95  In order to achieve SS in real-time, ICNet was used for paddy field images in this study,

96  and its structure is shown in Fig. 1. In this figure, numbers in parentheses are feature map size

97  ratios to the full-resolution input (640 $\times$480 pixels). Operations are highlighted in brackets.

98  The final $\times$4 upsampling in the bottom branch is only used during testing. The ICNet takes

99  cascade image inputs (i.e., medium- and high-resolution images), and it adopts a pyramid

100 pooling module (PPM) and cascade feature fusion (CFF) unit in Fig. 2. It was trained with

101 cascade label guidance. Different-scale (e.g., 1/16, 1/8, and 1/4) ground truth labels were

102 utilized to guide the learning stage of low, medium and high-resolution input.

103 As shown in Fig. 2a, the PPM fuses four different pyramid scale features, and 'POOL'

104 means pooling layer in the figure. First, it separates the feature map into different sub-regions

105 by using an operation called adaptive average pool. Then upsampling the low-dimension

106 feature maps to get the same size feature as the original feature map via bilinear interpolation.

107 Finally, different levels of the features are summed as a final pyramid pooling global feature.

108 To combine cascade features from different resolution inputs, 2 CFF units were used in the

109 ICNet. Details of the structure is shown in Fig. 2b, the sizes of feature maps F1 and F2 are C1

110 $\times$ H1 $\times$ W1 and C2 $\times$ H2 $\times$ W2, respectively, and the resolution of the label is 1 $\times$ H2 $\times$ W2,

111 where H2 = 2 $\times$ H1. It combines feature maps F1 and F2. In order to enhance the learning of

112 F1, auxiliary label on the upsampled feature of F1 is applied.

113 **2.2 Network Slimming algorithm**

114 The algorithm used in this paper to prune the network model was based on the principle

115 of Network Slimming method (Liu et al., 2017). The method could remove the less important

116 connections with small weights in each convolution layer. As we know that, batch

117 normalization (BN) layer performs the following transformation after each convolution layer

118 in the model:

119
$$\hat{z} = \frac{z_{in} - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}}$$
(1)

120
$$z_{out} = \gamma\hat{z} + \beta$$
(2)

121   where $z_{in}$ and $z_{out}$ are the input and output of a BN layer, $\mu_B$ and $\sigma_B$ are the mean and standard

122   deviation values of input activations over B, B denotes the current minibatch, $\gamma$ and $\beta$ are

123   trainable affine transformation parameters (scale and shift).

124       As $\gamma$ in BN layers corresponds to a specific convolutional channel, $\gamma$ was used for channel

125   scaling factors. The approach imposes L1 regularization (the latter part in formula (3)) on the

126   channel scaling factors $\gamma$ in BN layers for each channel. Pushing the values of channel scaling

127   factors towards zero with L1 regularization enables insignificant channels to be identified. The

128   network weights and these channel scaling factors were trained with sparsity regularization

129   (the latter part in formula (3)). The training objective of our approach is given by

130   $$L = \sum_{(x,y)} l(f(x,W),y) + \lambda \sum_{\gamma \in T} g(\gamma) \tag{3}$$

131   where $(x, y)$ denotes the training input and target, $W$ denotes the trainable weights, the first

132   sum-term corresponds to the normal training loss of a CNN, $T$ denotes the gradient of each

133   convolution layer, $g(\cdot)$ is a sparsity-induced penalty on the channel scaling factors, and $\lambda$

134   balances the two terms. In our experiment, we chose $g(s)=|s|$, which is known as L1-norm and

135   widely used to achieve sparsity. Subgradient descent was adopted as the optimization method

136   for the non-smooth L1 penalty term. The channel scaling factors act as the agents for channel

137   selection. As they were jointly optimized with the network weights, the network can

138   automatically identify insignificant channels, which can be safely removed without greatly

139   affecting the generalization performance. Channels with small factors $\gamma$ removed (all their

140   incoming and outgoing connections), then we could get the pruned network model.

141   **2.3 Dataset for semantic segmentation models**

142       All the images in the training set, validation set and test set were derived from

143   experimental videos from December 2016 to August 2019. A detailed description of the above

144   data set is shown in Table 1, which includes sample number, rice variety, field type, weather,

145 camera angle, camera depression angle, etc. Since some of the scenes in the video are not

146 related to the field scene, and sometimes some areas in the image are not clear enough, so some

147 clear images of appropriate size were cut out from the original images. Then the cut images

148 were rotated (±15°), and flipped horizontally. Finally, a total of 5000 images (jpeg format)

149 were obtained. The size of all images was 640×480 pixels, and the mean value of the RGB

150 channels of the images were 0.485, 0.456, and 0.406, and the standard deviation were 0.229,

151 0.224 and 0.225, respectively, when these images were transformed to the range of [0, 1].

152 According to a previous field trial video, a training set and test set were prepared, including

153 4,000 and 1,000 images, respectively, which were selected up from the data set of 5,000 images

154 mentioned before; Then the data was normalized to reduce the negative effects of uneven

155 brightness.

156 **2.4 The procedure of getting pruned (segmentation) model**

157 During the training process, a stochastic gradient descent method was used for backward

158 propagation of the learning phase to obtain the best network parameters. The initial learning

159 rate was 0.02, and the decay coefficient of the learning rate was 0.5. The decay frequency was

160 10 epochs, with the batch size of 4. The regularization parameter λ was 0.0001, with a penalty

161 factor 0.0001 to perform channel-level sparsity regularization. When the current loss function

162 converged and stabilized, training was halted.

163 After sparsity training, we removed channels with a global threshold γ1 across all layers

164 except for CFF, which was defined as a certain percentile of all the needed scaling factor values.

165 Such as 5-th percentile, corresponding to a 5 % compaction ratio. Then the compaction ratio is

166 defined as

167
$$\text{compaction ratio} = \frac{c_1}{c_2} \times 100\% \tag{4}$$

168  where $C_1$ is the numbers of removed channels, $C_2$ is the numbers of channels in the original

169  network. Two convolution layers before the 'sum' operation in the CFF unit were required to

170  have the same channel number. To match the feature channels of the 2 layers, we iterate through

171  the layers and perform the same percentile compaction operation to generate a pruning mask

172  for these connected layers, respectively. The percentile is same as the percentile used for the

173  global threshold γ1.

174      The channel pruning procedure is shown in Fig. 3. ICNet was initially trained with

175  channel-level sparsity regularization; sequentially, pruned ICNet was obtained by pruning

176  feature channels to a certain ratio according to their scaling factors in the ICNet; After channel

177  pruning, a fine-tuning operation was performed on pruned models to compensate potentially

178  temporary degradation in segmentation accuracy. The Network Slimming (training with

179  sparsity regularization, pruning, and fine-tuning) was repeated several times. The model was

180  pruned 10 % each time. In our experiments, we directly retrain using the same training hyper-

181  parameters as the initially training of ICNet.

182  **2.5 Experimental condition in field test**

183      In order to evaluate the performance of the pruned model, the pruned model which with

184  the compaction ratio of 80 % was used in the field test for human detection. Tests were

185  conducted in actual paddies, the places were in Kisosaki, Kuwana District, Mie, Japan

186  (35°05'19.9"N, on Aug. 24-25, 2019 and Nantan City, Kyoto, Japan (35°02'35.4"N,

187  136°46'16.8"E), on Sep. 22, 2019. The weather was sunny on Aug. 24-25, 2019 and cloudy on

188  Sep. 22, 2019. Fig. 4 shows the main devices used in this study. The base machine was a four-

189  row head-feeding combine harvester ER470 (Kubota, Osaka, Japan). Our laboratory has

190  developed an autonomous harvesting system based on ER470 (Iida et al., 2017), which could

191  follow a target path based on its absolute position and orientation, planning a counterclockwise

192  spiral path in a rectangular paddy field. An Intel RealSense D435 (Intel, Santa Clara, USA)

193  camera was mounted on the front of the harvester to capture color images and depth data in

194  real-time. It was mounted at a height of 1.65 m above the ground, and with its lens facing down

195  at an angle of 28 ° to the horizontal. A Jetson Xavier (NVIDIA, Santa Clara, USA) was used

196  for running segmentation models. Light levels were measured under different conditions with

197  a digital light meter KG-75 (Kaise, Nagano, Japan).

198  Since the pruned models could segment 5 classes (harvested area, unharvested area, ridges

199  area, human and background) and the combine could harvest rice automatically, the test was

200  conducted in this way. When the combine automatically harvested along the target path at a

201  speed of 1.0 m/s, a human would appear at different times on the target path in front of the

202  combine at different distances. In these conditions based on segmentation results and the

203  distance obtained by the depth camera, the combine took three actions timely, either stopping,

204  slowing down or continuing to work. The principal flow of the test algorithm for automatic rice

205  harvesting is shown in Fig. 5. When the combine begins harvesting, it captures color images

206  and depth data from the D435 camera, then inputs the RGB images into the segmentation model.

207  Based on the segmentation results, if there is a human in the image, it calculates the center of

208  the human area and gets the corresponding distance to the center from the depth data. Then

209  according to the distance between the combine and the human, it sends a control signal to stop,

210  slow down or continue to work, to the combine's electronic control unit through an RS-232

211  serial port. Two tests were conducted, in Test 1, a human appeared on the target path in front

212  of the combine at different distances. In Test 2, no human appeared on the target path in front

213  of the combine.

214  **3 Results and Discussion**
215  **3.1 Comparison of segmentation performance**
216  To evaluate the robustness of the models, each of the pruned models were validated on

217  the test set. Fig. 6 presents the mean intersection over union (IoU) at different compaction ratios.

218    Based on the data in Fig. 6, the following results were found. As the compaction ratio continues

219    to increase, there is a small loss in the accuracy of the model. In our experiments, the fine-

220    tuned pruned model could even achieve higher accuracy than the original unpruned model in

221    some cases (compaction ratio: 61.4 %). However, when the compaction ratio is greater than

222    80 %, the accuracy of the model seriously degrades. When the compaction ratio is less than

223    80 %, compelling results are achieved in comparison to the original counterpart.

224         When a combine harvester is working in the field, the harvested area, unharvested area,

225    ridge area, and human area occupy different ratios at different stages. Fig. 7 shows the

226    segmentation result for each class at different compaction ratios. It shows that models are more

227    inclined to predict pixels in the image as the harvested area and the unharvested area. This may

228    be due to data imbalance, because the harvested area and unharvested area occupy a bigger part

229    than the ridge area and human area for most images. It shows that when the compaction ratio

230    is less than 80 %, the mean IoU for each class at different compaction ratios is close to the

231    original counterpart. When the compaction ratio is greater than 80 %, the mean IoU for each

232    class decreases quickly.

233    **3.2 Inference run-time performance**

234         All these models were tested with an image size of 640 ×480 pixels, Table 2 shows the

235    frames per second (FPS) on the Jetson Xavier and model volume of different pruned models.

236    Because the accuracy of the model drops sharply when the compaction ratio is greater than

237    80%, only models with a compaction ratio of less than 80% were measured. The run-times

238    were achieved using CUDA 10.0.117 and cuDNN 7.3.07. As can be seen from Table 2 as the

239    compaction ratio increases, the size of the model volume decreases and the speed of the model

240    increases. When 80 % of the channel been pruned, the model has a 97.4% reduction of model

241    volume size, and ran 1.33 times faster with comparable detection accuracy to the original model.

242    It can be known from segmentation performance and the inference run-time performance that

243　when the compaction ratio is less than 80 %, the Network Slimming method could be used for

244　decreasing the computational cost of the ICNet for field image segmentation.

**3.4 Results and discussion in field test**

246　　　Because the combine harvester traveled counter-clockwise during the harvesting, as

247　shown in Fig. 8, the noise levels in the acquired images differed as the light conditions changed

248　depending on the direction of movement. All of the test scenes were categorized into four

249　scenes (A, B, C, D) according to the direction of harvester movement. In all scenarios, based

250　on the segmentation results of the model, the harvester would take the appropriate action (stop,

251　slow down or continue to work). Table 3 shows the results of Test 1 by using the pruned model.

252　Because a human always appeared on the target path in front of the combine at different

253　distances during Test 1, if the harvester slowed down, stopped and then continued to work, it

254　was regarded as a successful result.

255　　　The results show that the average success rate of collision avoidance was 96.6% at an

256　average processing speed of 32.2 FPS (31.1 ms per frame). The evaluation results show that

257　the proposed method is effective for human segmentation and collision avoidance regardless

258　of the movement direction of the combine harvest or the light conditions experienced, as shown

259　in Fig. 9. However, as shown in the last column of Fig. 9, the human is not successful

260　segmentation when the camera is backlighted (dataset B). Because the camera in scene B was

261　in backlight mode, the sunlight affected the image quality obtained by the camera, which

262　reduced the accuracy of model segmentation. Finally, it made the success rate in scene B lower

263　than that in other scenes.

264　　　Table 4 shows the result of Test 2. Because no humans appeared on the target path in front

265　of the combine in Test 2, the harvester should continue to work normally, so we focused on the

266　number of false results in this test. If the harvester slowed down or stopped, it was regarded as

267　a false result. The result in Table 4 indicate that the number of false detection was small under

268 various light conditions. However, the segmentation is not successful when the camera is

269 backlighted (first column in Fig. 10), or the shadow of rice is similar to that of a human (second

270 column in Fig. 10).

271     It can be known from two field tests that when the camera is in a backlight mode or some

272 objects are visually similar to a human in the image, the SS model that only relied on a color

273 image as input still has the probability of false detection. Since thermal images and Lidar data

274 are less affected by light than color images, which could provide additional information for

275 making detection. So, our future work is to fuse the thermal image or Lidar data for further

276 improving the accuracy of detection.

277 **4 Conclusion**

278     1) Network Slimming based on ICNet was proposed and evaluated as a means to compact

279 the semantic segmentation model. It directly imposes sparsity-induced regularization on the

280 scaling factors in batch normalization layers, and unimportant channels in convolutional layers

281 can thus be automatically identified during training.

282     2) The pruned models, which were achieved through channel pruning of the convolutional

283 layers, substantially decreased the computational cost of ICNet, with a slightly reduction in

284 accuracy. When the compaction ratio is 80 %, it gives a 97.4 % reduction of model volume

285 size, running 1.33 times faster with comparable detection accuracy as the original model.

286     3) A pruned model (with 80 % compaction ratio) was then tested in the field to validate

287 the feasibility of the method. Results showed that the average success rate of collision

288 avoidance was 96.6% at an average processing speed of 32.2 FPS (31.1 ms per frame) with an

289 image size of $640 \times 480$ pixels on a Jetson Xavier. Results demonstrate that with channel

290 reduction of the convolutional layer in the ICNet, a pruned (segmentation) model can be

291 successfully used in a rice combine harvester for obstacle detection and collision avoidance in

292 real time.

293 **References**

294 Atlas ML. 2018. Real-Time Semantic Segmentation on CamVid. Accessed March 15, 2020.

295 https://paperswithcode.com/sota/real-time-semantic-segmentation-on-camvid.

296 Brostow, G. J., Fauqueur, J., Cipolla, R. 2009. Semantic object classes in video: A high-

297 definition ground truth database. Pattern Recognition Letters, 30(2), 88-97.

298 Cabrera, R.R., Tuytelaars, T., Gool, L.V., 2011. Efficient multi-camera detection, tracking, and

299 identification using a shared set of haar-features. In Proceedings of the IEEE Conference on

300 Computer Vision and Pattern Recognition, pp. 65-71.

301 Chen, W., Wilson, J., Tyree, S., et al., 2015. Compressing neural networks with the hashing

302 trick. In International Conference on Machine Learning, pp. 2285-2294.

303 Cho, W., Iida, M., Suguri, M., et al., 2014a. Vision-based uncut crop edge detection for

304 automated guidance of head-feeding combine. Engineering in Agriculture, Environment and

305 Food, 7(2), 97-102.

306 Cho, W., Iida, M., Suguri, M., et al., 2014b. Using multiple sensors to detect uncut crop edges

307 for autonomous guidance systems of head-feeding combine harvesters. Engineering in

308 Agriculture, Environment and Food, 7(3), 115-121.

309 Courbariaux, M., Hubara, I., Soudry, D., et al., 2016. Binarized neural networks: Training deep

310 neural networks with weights and activations constrained to + 1 or -1. arXiv preprint

311 arXiv:1602.02830.

312 Denton, E. L., Zaremba, W., Bruna, J., et al., 2014. Exploiting linear structure within

313 convolutional networks for efficient evaluation. In Advances in neural information processing

314 systems, pp. 1269-1277.

315 Huang, G., Chen, D., Li, T., et al., 2017. Multi-scale dense convolutional networks for efficient

316    prediction. arXiv preprint arXiv:1703.09844, 2.

317    Han, S., Pool, J., Tran, J., et al., 2015. Learning both weights and connections for efficient

318    neural network. In Advances in neural information processing systems. pp. 1135-1143.

319    He, K., Zhang, X., Ren, S., et al., 2015. Delving deep into rectifiers: Surpassing human-level

320    performance on imagenet classification. In Proceedings of the IEEE international conference

321    on computer vision, pp. 1026-1034.

322    Hisae, T., Masuda, R., Suguri, M., et al., 2017. Human detection in paddy field by image

323    processing (in Japanese). Journal of the Japanese Society of Agricultural Machinery and Food

324    Engineers, 79(1), 49-58.

325    Iida, M., Harada, S., Sasaki, R., et al., 2017. Multi-combine robot system for rice harvesting

326    operation. ASABE paper No.1700321. DOI:10.13031/aim.201700321

327    Kamilaris, A., Prenafeta-Boldu, F.X., 2018. Deep learning in agriculture: a survey. Computers

328    and Electronics in Agriculture, 147, 70-90.

329    Liu, Z., Li, J., Shen, Z., et al., 2017. Learning efficient convolutional networks through

330    Network Slimming. In Proceedings of the IEEE International Conference on Computer Vision,

331    pp. 2736-2744.

332    McCool, C., Perez, T., Upcroft, B., 2017. Mixtures of lightweight deep convolutional neural

333    networks: applied to agricultural robotics. IEEE Robotics and Automation Letters, 2(3), 1344-

334    1351.

335    Milioto, A., Lottes, P., Stachniss, C., 2018. Real-time semantic segmentation of crop and weed

336    for precision agriculture robots leveraging background knowledge in CNNs. IEEE

337    International Conference on Robotics and Automation, Brisbane, Australia, pp. 2229-2235.

338    Rastegari, M., Ordonez, V., Redmon, J., et al., 2016. Xnor-net: Imagenet classification using

339    binary convolutional neural networks. In European Conference on Computer Vision, pp. 525-

340    542.

341     Singh, S., Kaur, A., Taqdir, A., 2015. A face recognition technique using local binary pattern

342     method. International Journal of Advanced Research in Computer and Communication

343     Engineering, 4(3), 165-168.

344     Takagaki, A., Masuda, R., Iida, M., et al., 2013. Image processing for ridge/furrow

345     discrimination for autonomous agricultural vehicles navigation. IFAC Proceedings Volumes,

346     46(18), 47-51.

347     Tang, H., Wang, W., Gimpel, K., et al., 2016. End-to-end training approaches for

348     discriminative segmental models. Spoken Language Technology Workshop, pp. 496-502.

349     Wang, Y., Zhou, Q., Liu, J., et al., 2019. LEDNet: A Lightweight Encoder-Decoder Network

350     for Real-Time Semantic Segmentation. arXiv preprint arXiv:1905.02423.

351     Yang, A., Xue, Y., Huang, H., et al., 2017. Lactating sow image segmentation based on fully

352     convolutional networks. Transactions of the Chinese Society of Agricultural Engineering,

353     33(23), 219-225. (in Chinese with English abstract)

354     Yao, S., Pan, S., Wang, T., et al., 2015. A new pedestrian detection method based on combined

355     HOG and LSS features. Neurocomputing, 151, 1006-1014.

356     Ye J, Lu X, Lin Z, et al. 2018. Rethinking the smaller-norm-less-informative assumption in

357     channel pruning of convolution layers[J]. arXiv preprint arXiv:1802.00124.

358     Yu, C., Wang, J., Peng, C., et al., 2018. Bisenet: Bilateral segmentation network for real-time

359     semantic segmentation. In Proceedings of the European Conference on Computer Vision

360     (ECCV), pp. 325-341.

361     Zhao, H., Qi, X., Shen, X., et al., 2018. Icnet for real-time semantic segmentation on high-

362     resolution images. In Proceedings of the European Conference on Computer Vision (ECCV),

363     pp. 405-420.

364

365 **Figures**
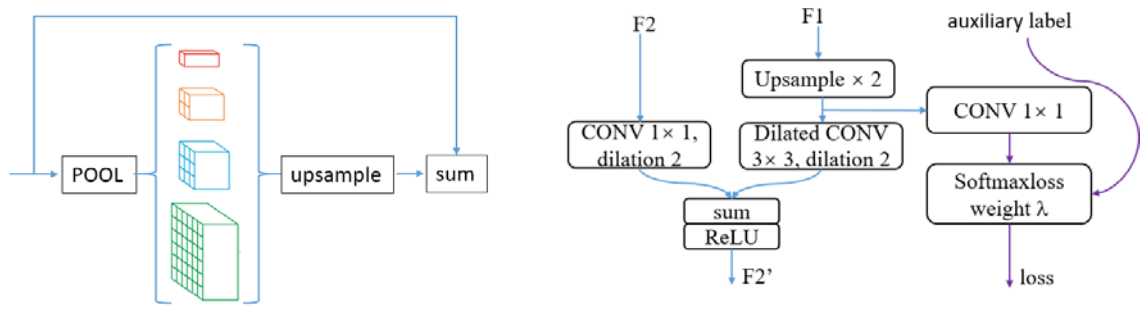
366



Fig. 1 ICNet structure for image segmentation.
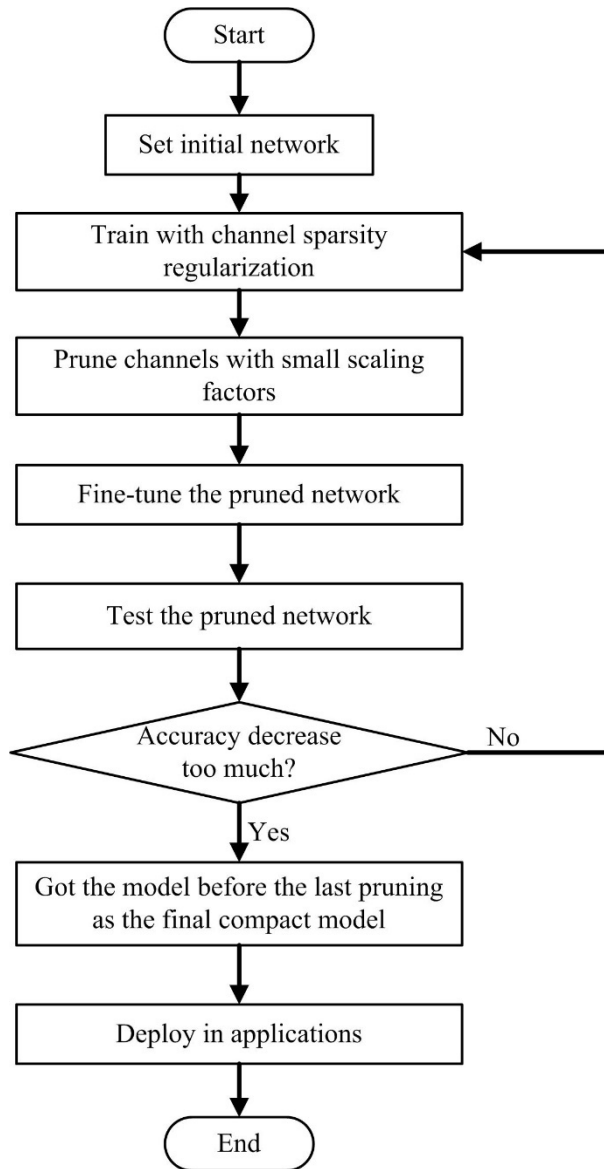
369

370



(a) Pyramid pooling module(PPM)          (b) Cascade feat fusion unit(CFF)

Fig. 2 Pyramid pooling module (PPM) and cascade feat fusion (CFF) unit in ICNet.
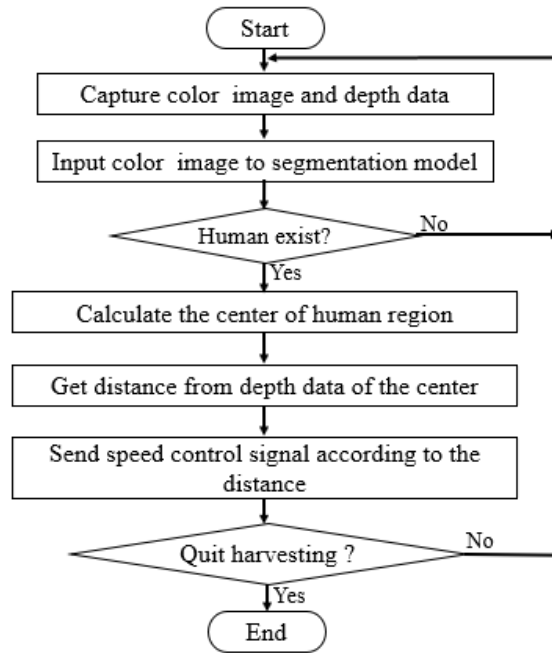
372

373

374

Fig. 3. An iterative procedure of getting efficient segmentation model through sparsity training and channel pruning.

377



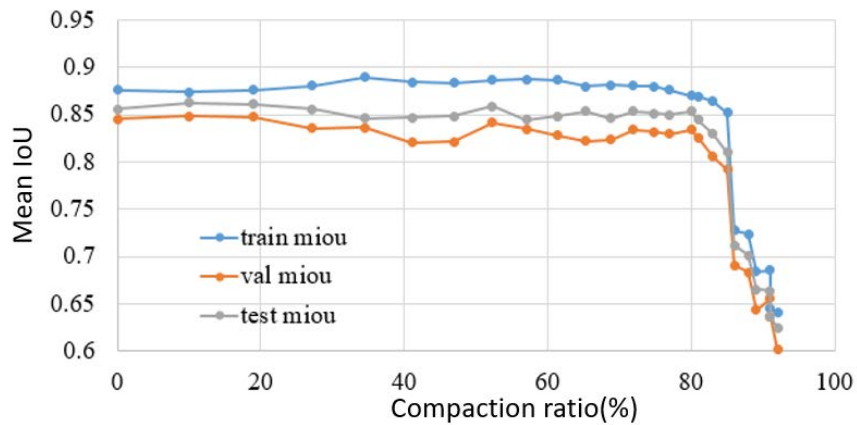Fig. 4. Robotic combine harvester and devices installed.

378

379



380

381        Fig. 5. The principal flow of the test algorithm for automatic rice harvesting.

382

383

384



385
386                Fig. 6. The mean IoU at different compaction ratio.

387

388
389
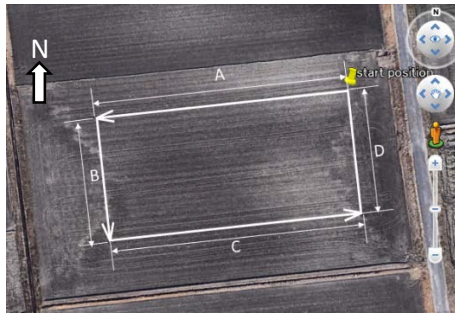390 Fig. 7. Segmentation accuracy of mean IoU for each class at different compaction ratio.
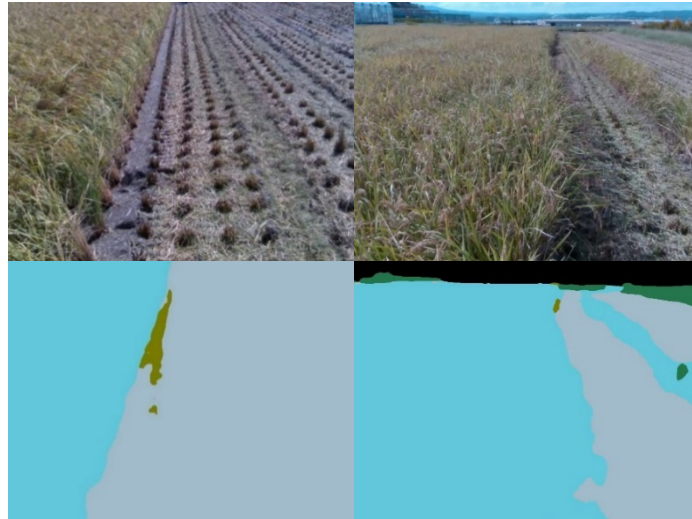391

392

393



Fig. 8. Movement direction of the robotic combine harvester in paddy field in Kisosaki.

394

395



396 Fig. 9. Examples of images(top) and segmentation results (bottom) of the model in Test 1.

397

398

399

400     Fig. 10. Examples of images (top) and outputs (bottom) from SS model in Test 2.

401

402 Tables
403
404                        Table 1 Description of the dataset.

| Items | Description of dataset source 1 | Description of dataset source 2 |
|---|---|---|
| Camera | GoPro HERO5 | Intel RealSense D435 |
| Time | Morning, Noon | Afternoon |
| Weather | Cloudy, Sunny | Cloudy, Sunny |
| Place | Nantan, Japan; Narita, Japan; | Narita, Japan; Kizu, Japan; |
| Source image Size(width × height) | 1920 × 1080 | 640 × 480 |
| Rice variety | KoshiHikari, Husakogane | Husakogane, HinoHikari |
| Source sample Number | 700 | 550 |
| Field type | paddy field | paddy field |
| Camera height | 1.75 m | 1.75 m |
| Camera depression angle | The lens is facing down and at an angle of 15 degrees to the horizontal | The lens is facing down and at an angle of 15 degrees to the horizontal |

405
406

407   Table 2 The frames per second (FPS) on the Jetson Xavier and model volume of different
408                                pruned models.

| compaction ratio (%) | FPS on Xavier | inference time (ms) | Volume size of model parameter file (MB) |
|---|---|---|---|
| 0 | 24.2 | 41.3 | 30.8 |
| 10.0 | 26.2 | 38.2 | 26.7 |
| 19.0 | 28.4 | 35.2 | 21.8 |
| 27.1 | 29.0 | 34.4 | 17.5 |
| 34.4 | 29.4 | 34.0 | 13.8 |
| 41.0 | 29.5 | 33.9 | 10.8 |
| 46.9 | 30.5 | 32.8 | 8.5 |
| 52.3 | 30.6 | 32.7 | 6.6 |
| 57.1 | 30.5 | 32.8 | 5.1 |
| 61.4 | 30.5 | 32.8 | 4.1 |
| 65.3 | 30.7 | 32.6 | 3.1 |
| 68.8 | 30.8 | 32.5 | 2.4 |
| 71.9 | 32.0 | 31.3 | 1.8 |
| 74.7 | 31.3 | 31.9 | 1.4 |
| 77.0 | 32.0 | 31.3 | 1.0 |
| 80.0 | 32.2 | 31.1 | 0.8 |

409

410

Table 3 Results of Test 1 by the pruned models.

| Illumination [lx] | Movement direction | Number of times human appeared | Number of successes |
|---|---|---|---|
| 32500 ~ 54850 | A | 10 | 10 |
| | B | 5 | 4 |
| | C | 10 | 10 |
| | D | 5 | 5 |
| 63210 ~ 79610 | A | 10 | 10 |
| | B | 5 | 4 |
| | C | 10 | 10 |
| | D | 5 | 5 |

411

412

413
414

415

Table 4 Results of Test 2 by the pruned models.

| Illumination [lx] | Movement direction | Travel distance[m] | Number of failures |
|---|---|---|---|
| 32500 ~ 54850 | A | 200 | 0 |
| | B | 100 | 0 |
| | C | 200 | 0 |
| | D | 100 | 0 |
| 63210 ~ 79610 | A | 200 | 0 |
| | B | 100 | 1 |
| | C | 200 | 1 |
| | D | 100 | 0 |

416