1  **Length estimation of fish detected as non-occluded using a smartphone application**

2  **and deep learning techniques**

3  Yasutoki Shibata[1], Yuka Iwahara[1], Masahiro Manano[1], Ayumi Kanaya[1], Ryota Sone[2],

4  Satoko Tamura[3], Naoya Kakuta[3], Tomoya Nishino[4], Akira Ishihara[4] and Shungo

5  Kugai[4]

6

7  1) Fisheries Resources Institute, National Research and Development Agency, Japan

8      Fisheries Research and Education Agency, Yokohama, Kanagawa, Japan

9  2) Marine Resources Research Center, Aichi Fisheries Research Institute, Toyohama,

10      Aichi, Japan

11  3) Fisheries Technology Center Sagami Bay Experiment Station Of Kanagawa

12      Prefectural Government, Odawara, Kanagawa, Japan

13  4) Computermind Corp., Nishi-Shinjyuku, Tokyo, Japan

14

15  **Correspondence**

16  Yasutoki Shibata, Fisheries Resources Institute, National Research and Development

17  Agency, Japan Fisheries Research and Education Agency, 2-12-4 Fukuura, Kanazawa,

18  Yokohama, Kanagawa 236-8648, Japan.

19    shibata_yasutoki20@fra.go.jp

20

**Abstract**

22    Uncertainty in stock assessment can be reduced if accurate and precise length

23    composition of catch is available. Length data are usually manually collected, although

24    this method is costly and time-consuming. Recently, some studies have estimated fish

25    species and length from images using deep learning by installing camera systems in

26    fishing vessels or a fish auction center. Once the deep learning model is properly trained,

27    it does not require expensive and time-consuming manual labor. However, several

28    previous studies have focused on monitoring fishing practices using an electronic

29    monitoring system (EMS); therefore, it is necessary to solve many challenges, such as

30    counting the total number of fish in the catch. In this study, we proposed a new deep

31    learning-based method to estimate fish length using images. Species identification was

32    not performed by the model, and images were taken manually by the measurers;

33    however, length composition was obtained only for non-occluded fish detected by the

34    model. A smartphone application was developed to calculate scale information

35    (cm/pixel) from a known size fish box in fish images, and the Mask R-CNN (Region-

36    based convolutional neural networks) model was trained using 76,161 fish to predict

37    non-occluded fish. Two experiments were conducted to confirm whether the proposed

38    method resulted in errors in the length composition. First, we manually measured the

39    total length (TL) for each of the five fish categories and estimated the TL using deep

40    learning and calculated the bias. Second, multiple fish in a fish box were photographed

41    simultaneously, and the difference between the mean TL estimated from the non-

42    occluded fish and the true TL from all fish was calculated. The results indicated that the

43    biases of all five species categories were within $\pm$ 3%. Moreover, the difference was

44    within $\pm$ 1.5% regardless of the number of fish in the fish box. In the proposed method,

45    deep learning was used not to replace the measurer but to increase their measurement

46    efficiency. The proposed method is expected to increase opportunities for the

47    application of deep learning-based fish length estimation in areas of research that are

48    different from the scope of conventional EMS.

49

50    **Keywords**

51    Mask R-CNN, mobile imaging, occlusion, total length composition in catch, stock

52    assessment

53

54    **1. Introduction**

55     It is important to accurately estimate current and future population of target fish

56     stocks to maintain maximum sustainable yield (MSY), and population dynamics models,

57     such as age-structured models, have been used to estimate abundance in stock

58     assessments (Ichinokawa et al., 2017; Privitera-Johnson and Punt, 2020). Length-at-age

59     relationships are especially important for these models because the observed fish length

60     compositions are often converted to age compositions using these relationships (Piner et

61     al., 2016). Age composition data are important for estimating recruits, selectivity,

62     fishing and natural mortality rates, and relative stock sizes (Lee et al., 2014; Piner et al.,

63     2016; Shibata et al., 2021). Thus, the length-at-age relationship that has been used to

64     transform length to age is very influential in stock assessments (Wang et al., 2014).

65     However, obtaining the length composition from landed catch is costly and time-

66     consuming because fish length is usually measured directly by hand (Palmer et al.,

67     2022). In other words, the sample size for fish length of the target species depends on

68     the number of measurers, i.e., the people that measure the fish. This cost often implies

69     reduced sample sizes, which may lead to information loss and negatively affect the

70     accuracy and precision of estimated abundance.

71     In recent years, some studies have estimated fish length using deep learning methods

72     (Álvarez-Ellacuría et al., 2020; Lekunberri et al., 2022; Ovalle et al., 2022; Palmer et al.,

4

73    2022). In these studies, cameras were set in fishing vessels or the fish auction centers,

74    and the obtained images of the catch were analyzed using deep learning techniques.

75    Once the deep learning model is properly trained, it does not require expensive and

76    time-consuming manual labor (Lekunberri et al., 2022). The results of these studies

77    could be applied to obtain fish length data and eliminate manual labor dependency.

78    However, three major issues must be considered when installing these camera systems.

79     First, power supplies and spaces for imaging are needed to setup cameras and secure

80    their field of view. Some large vessels or fish auction centers could carry out length

81    estimation, as in previous studies; however, that is not the case for majority. For

82    example, more than 75% of vessels in America and 80% in Asia, Europe, Africa, and

83    Oceania     are     less     than     12     m     long     (FAO,     2022

84    https://www.fao.org/3/cc0461en/online/sofia/2022/fishing-fleet.html).  The  availability

85    of stable electrical power in small vessels may be limited by the battery capacity when

86    an engine is not running (van Helmond et al., 2020). Although an electronic monitoring

87    system (EMS) for small vessels (i.e., 10–12 m) has been developed using portable

88    power  with  solar  panels,  only  a  single  camera  can  be  installed  because  of  space

89    limitations or vulnerability to hidden activity outside the field of view (Bartholomew et

90    al., 2018). Moreover, in some cases, some construction is needed to use a power supply

91    in the fish auction center because power supply is usually not available in these areas

92    (probably to prevent a short circuit due to weathering). In addition, imaging can often be

93    blocked by human movement when belt conveyors carrying fish in a fish auction center

94    or landing fishing port are not installed, which may be the case in most fisheries.

95    Second, there are difficulties in identifying fish species using deep learning

96    techniques. Identification of Fish species using deep learning techniques involves

97    families and genera (e.g., van Essen et al., 2021). Some examples have been reported as

98    possible situations to identify fish species, although the study is based on a few (ten

99    classes) fish species (Lu et al., 2020), or the surrounding environment of imaging is

100    fixed on the ship (Ovalle et al., 2022).

101    Finally, occlusion of fish bodies by other fish or objects (i.e., only a part of the fish

102    body is visible in the image) affects identification and length estimation. A previous

103    study estimated the total length through head length using Mask R-CNN (Region-based

104    convolutional neural networks) (He et al., 2017) and weight information (Álvarez-

105    Ellacuría et al., 2020). Another study estimated the total length using a deep learning

106    method corrected using Bayesian estimation (Palmer et al., 2022). A common aspect of

107    both studies was the estimation of the total length as a function of some other

108    information. This method is very useful because it can be used even when fish have

6

109    been occluded, but some effort may be required because the equation must change

110    depending on various factors, such as fish species and season.

111        These challenges need to be solved to collect fish length of all fish species using the

112    deep learning method as an alternative to measurers. However, none of them are

113    fundamentally solvable with current equipment and technology. In contrast, it is

114    possible to increase the sample size of fish length data with the burden of measurers

115    decreasing if the purpose of using the deep learning method was changed. Here, image

116    analysis using the deep learning method was used to obtain only the total length

117    composition. The identification of fish species and investigation of discarding or

118    bycatch on commercial vessels using EMS and deep learning methods were not the

119    objective of this study.

120        We propose a method that combines an application and a deep learning method to

121    assist the measurement of catch length composition data. Measurers identify the fish

122    species, hold a camera in their hands, and directly take photographs of the fish. When

123    photographing, an app called ToroCam (an app to record coordinates of the fish box

124    through the smartphone camera) developed in this study will be used. ToroCam adds the

125    coordinates of the four corners of a box containing fish (fish box) to the JPEG image as

126    pixel values. Subsequently, the value of the cm/pixel can be calculated if the size of the

127    fish box in centimeters is known. After photographing with ToroCam, only non-

128    occluded fish can be detected using deep learning techniques. The length composition in

129    the catch can be obtained by multiplying the value of cm/pixel with the lengths of non-

130    occluded fish in pixels. Although some limitations remain because measurers identify

131    species and take images manually, but the fact that fish length can be collected

132    automatically makes it more efficient because body length information can be typed into

133    a field notebook or computer, and the difficulty of cleaning measuring instruments can

134    be reduced. In addition, continuous security and power supply of the camera are not

135    required. The hardware required is only a smartphone, and a computer for analysis is

136    cheaper than camera systems. Because this is a sampling rather than a full count

137    measurement, fish length extracted from only non-occluded fish can be used for stock

138    assessment when the occlusion is random for actual fish size.

139    The purpose of this study is to develop a method where measurers carry out direct

140    imaging (mobile imaging) with a portable tool, such as a smartphone, even under

141    unsteady imaging conditions, and obtain the total length of non-occluded fish using

142    deep learning techniques. This method will enable the determination of fish length

143    composition by image analysis, even in situations where the methods of fish image

144    analysis developed based on conventional EMS have not been targeted.

145

## 2. Materials and Methods

147    2.1 Image acquisition for learning

148    We obtained 8,087 fish images from a fish market, two fishing ports, a bottom trawl

149    vessel, and a research vessel (R/V): Matsuura fish market, Odawara and Toyohama

150    fishing ports, bottom trawl vessel (Horyo-maru No. 18), and R/V Kaiyo-maru No. 5

151    (Table 1). Horyo-maru No. 18 and R/V Kaiyo-maru No. 5 operated commercial bottom

152    trawling and a scientific bottom trawling survey, respectively, in the southern Hokkaido

153    area (Fig. 1). Fish images were either captured with a fixed camera or manually

154    captured by researchers. The images were collected using a camera mounted directly

155    above a conveyor belt on the Matsuura fish market, Odawara fishing port, and the

156    bottom trawl vessel or the sorting platform of R/V Kaiyo-maru No. 5 (Fig. 2a-d). In

157    addition, at the Matsuura fish market, Odawara, and Toyohama fishing ports,

158    researchers used a camera and directly photographed landed fish (Fig. 2e-g). The

159    images were captured at 3:00–17:00 at the Matsuura fish market, 4:00–7:00 at the

160    Odawara fishing port, 16:00–17:00 at the Toyohama fishing port, 7:00–19:00 at the

161    Horyo-maru No. 5, and 7:30–17:00 at the R/V Kaiyo-maru No. 5.

162    The image sizes were as follows: $4800 \times 3200$, $2048 \times 1536$, $2704 \times 1520$, $1920 \times$

163  1080, 5184 × 3888, 960 × 1080, 3072 × 1728, 4608 × 3456, 424 × 480, and 640 × 720.

164  The images were taken from August 2020 to December 2021, whereas those from

165  Horyo-maru No. 18 were obtained from December 2015 to February 2016. Because the

166  images obtained from Horyo-maru No. 18 and Kaiyo-maru No. 5 were taken in fish-eye

167  mode, 25% of the left and right side of the images were removed, and the remaining

168  50% in the middle was extracted as the region of interest and used for the following

169  analysis.

170

171  2.2 Annotation and classes

172  All fish in the images were annotated using instance segmentation (n = 76,161).

173  Because fish identification was not the target of this study, detailed names of fish

174  species and their number of annotations are shown in Supplementary Table 1. All fish

175  bodies, including fins and spines, were annotated. The mantle was annotated for squids.

176  Exposure of the fish body was annotated as "*exposure*". The label was classified on

177  two scales as "F-100" and "F-other." For example, if the fish body is not occluded by

178  another fish body or object, the label *exposure* would be "F-100." If the body's visibility

179  is reduced from 99 to 1% of its whole body because of occlusion, then *exposure* was "F-

180  other." The numbers for each category are shown in Table 2. In addition, if fish could

10

181   not be annotated one by one and fish species could not be determined or the fish was

182   outside the conveyor, they were labeled as "Non-target" and those were annotated in a

183   same polygon (Fig. 3b). A difference between "Non-target" and "F-other" is that the

184   former does not provide instance segmentation for each individual and does not contain

185   information that can be used for species identification, while the latter can be used even

186   if it is occluded. This difference may not be very useful in this study because we did not

187   perform species identification. However, we anticipate that it will be useful to keep the

188   two separate considering future studies. If a fish is detected and identified as "F-100" by

189   a deep learning model in inference, the fish can be used to obtain the total length

190   composition.

191

192   2.3 Training and validation

193      Eighty percent of the annotated images were used for training, and the remaining

194   10% were used for validation. The remaining 10% was used as the test data to calculate

195   the confusion matrix. During training, the image size was resized to $800 \times 600$ pixels,

196   flipped upside down with 50% probability, and data expansion was performed. A

197   PyTorch library (Paszke et al., 2019) included in Python was used, and the estimation

198   was performed using Mask R-CNN (He et al., 2017). The development environment

11

199 was PyTorch 1.7.1, Python3.7, and a GPU NVIDIA 3800. The number of epochs was

200 fixed at 65, batch size was 4, and IoU (Intersection over Union) was 0.5. ResNet-50-

201 FPN was used as the backbone; the model classified only three classes: "F-100," "F-

202 other" and "Non-target." Here, only those with a probability of 0.8 or higher were used

203 to exclude ambiguous objects. The value of the loss function was calculated from the

204 data used for validation, and the weight parameters were adopted when the value was

205 the smallest.

206

207 2.4 Experiments and inference

208　　Two experiments were conducted from September 2022 to January 2023 to test the

209 performance of estimating fish length using "F-100" with ToroCam. A smartphone

210 (Sony, Xperia 5III, Android12, $0.82 \times 6.8 \times 15.7$ cm) was used to carry out the two

211 experiments and the size of image was fixed at $3024 \times 2268$. These image sets were not

212 included in the training or validation datasets.

213

214 *2.4.1 Experiment I. Difference between true and estimated length*

215　　The purpose of this experiment was to evaluate the difference in the total length

216 between the measured and estimated individuals. Each fish was photographed and

217    measured individually. At the Matsuura fish Market and Odawara fishing port, the

218    landed catch was randomly sampled and photographed after measuring the total length

219    of each fish. The measured total length was considered the true value. The shooting time

220    was within the same range as that when the training data were obtained. The species of

221    fish photographed were mackerels (*Scomber japonicus* and *Scomber australasicus*),

222    Japanese jack mackerels (*Trachurus japonicus*), Japanese sardines (*Sardinops*

223    *melanostictus*), red barracudas (*Sphyraena pinguis*), and bullet tuna (*Auxis rochei*).

224    Bullet tuna were only sampled at the Odawara fishing port, and the other five species

225    were sampled from the Matsuura fish market. *Scomber japonicus* and *Scomber*

226    *australasicus* were difficult to identify in the field survey; therefore, these two species

227    were treated under the same name category. Hereafter, when referring to both

228    simultaneously, they are denoted as mackerels. The sample size (*N*) of mackerels was

229    150, that of bullet tunas was 100, and that of the remaining species was 50.

230    The fish were photographed in a blue colored fish box (73 × 40.5 cm), although this

231    background color beneath the fish was not included in the training data. The fish were

232    photographed using ToroCam, a smartphone application (only for Android OS and the

233    app is now ready to be registered on Google Play) that displays a rectangle on the

234    smartphone screen and assigns the coordinates of the rectangle to the photographed

13

235 JPEG image (Fig. 4). In other words, the photographer must visually align the rectangle

236 on the screen with the corners of the fish box. Not only the coordinates, but also the

237 name of the fish, the true length of the fish box, and the location of the photograph can

238 be recorded automatically once declared at the time of shooting in the comment section

239 of the JPEG image.

240  ToroCam displays a rectangle on the camera screen of the smartphone, and after

241 aligning the rectangle with the four corners of the fish box containing the fish, the

242 photographer takes a picture. Because the distances between the four corners of the fish

243 box are easily known, the scale information (cm/pixel) was obtained by calculating the

244 number of pixels in the image between the coordinates of the four corners. The fish

245 determined as "F-100" in the image at the time of inference was enclosed by a rectangle,

246 and the number of pixels at long side was multiplied by the calculated cm/pixel value

247 and was converted to the total length. The relative differences $\hat{d}_n$ between the $n$th

248 observed total length $l_n$ and the calculated total length $\hat{l}_n$ and relative bias $(\hat{B})$ were

249 calculated for each fish species as follows:

250

251 $\hat{d}_n = \left(\hat{l}_n - l_n\right)/l_n \times 100,$     (1)

252 $\hat{B} = \sum_n^N \hat{d}_n/N.$     (2)

14

253

254

*2.4.2 Experiment II: Difference with increasing individuals in a fish box*

The purpose of this experiment was to evaluate the degree of change in the difference between the estimated and true length composition when fish in a fish box were increased. The experiment was conducted at the Matsuura fishing market using two fish categories: Japanese jack mackerels and mackerels. Fifty fish were prepared for each category. Fish were added individually to the blue fish box and photographed until $N = 10$ ($N=1, 2, ..., 10$). Thereafter, fish were added twice and photographed ($N=12, 14,...,50$). Here, we photographed three shots ($i=1,2,3$) of each $N$ with the fish randomly stirred by hand each time. In other words, the same fish were photographed three times, although their positions and degrees of occlusion differed each time. The experiment was conducted over two days, on January 17 and 18. Thus, 180 (($10 + 20$) $\times$ 3 $\times$ 2) photographs were taken for each fish category.

The $n$th fish ($n=1,…,N$) in the box were identified as three classes or missed to be detected by the deep learning model: "F-100" ($c=1$), "F-other" ($c=2$), "Non-target" or "not-detected," although "Non-target" fish were not included in this experiment and "not-detected" fish could not be counted. If the occluded area is sufficiently large, the

271 underlying fish are not visible, and the fish may not be detected. The number of

272 detected fish ($\widehat{N}_{N,c|i}$) as either "F-100" ($\widehat{N}_{N,c=1|i}$) or "F-other" ($\widehat{N}_{N,c=2|i}$) was counted,

273 where the value of $\widehat{N}_{N,c|i}$ was changed for each shot $i$ even if $N$ was the same because

274 the degree of occlusion was changed by hand. The detected rates of only non-occluded

275 individuals ($\hat{R}_{1,N,i}$) and both non-occluded and occluded individuals ($\hat{R}_{2,N,i}$) were

276 calculated for every shot $i$ under the true number $N$ using the following equation:

277

278 $\hat{R}_{1,N,i} = \widehat{N}_{N,c=1|i}/N \times 100,$ (3)

279 $\hat{R}_{2,N,i} = \sum_{c=1}^{2} \widehat{N}_{N,c|i}/N \times 100.$ (4)

280

281 The total length of $n$th individual identified as either $c=1$ or $c = 2$ for every shot $i$ under

282 the true number in the fish box $N$ ($\hat{l}_{n,N,c|i}$) was estimated using ToroCam and deep

283 learning methods. The method used to estimate the total length was the same as that in

284 Experiment 1, although the estimates were obtained for $c=1$ and $c=2$. The mean total

285 length of both non-occluded individuals ($\hat{L}_{1,N,i}$) and all detected individuals ($\hat{L}_{2,N,i}$) was

286 calculated for every shot $i$ and $N$ as follows:

287

288 $\hat{L}_{1,N,i} = \sum_{n=1}^{N} \left( \hat{l}_{n,N,c=1|i} \right)/\widehat{N}_{N,c=1|i},$ (5)

16

289   $\hat{L}_{2,N,i} = \sum_{n=1}^{N} \left(\hat{l}_{n,N,c|i}\right) / \sum_{c=1}^{2} \widehat{N}_{N,c|i}.$     (6)

290

291   Note that the $n$th estimated total length is not added if the individual is identified as $c=2$

292   in Eq. (5) but Eq. (6). Finally, a relative difference at each shot $i$ of only non-occluded

293   individuals $(\widehat{D}_{1,N,i})$ and all detected individuals $(\widehat{D}_{2,N,i})$ under the true fish number $N$

294   was calculated using the following equation:

295

296   $L_N = \sum_{n=1}^{N} l_{n,N}/N,$       (7)

297   $\widehat{D}_{1,N,i} = \left(\hat{L}_{1,N,i} - L_N\right)/L_N \times 100,$     (8)

298   $\widehat{D}_{2,N,i} = \left(\hat{L}_{2,N,i} - L_N\right)/L_N \times 100,$     (9)

299

300   where $L_N$ is the true mean total length in the fish box when the true sample size was $N$

301   and $L_N$ did not change for each shot $i$. A simple regression analysis was conducted

302   where the response and independent variables were $\widehat{D}_{1,N,i}$ and detected rate to confirm if

303   the difference was affected by the number of fish in the box.

304

305   **3. Results**

306   3.1 Confusion matrix

17

307    Because the loss function from the validation data was minimized at the $32^{nd}$ epoch,

308    weight parameters were used at that time. The confusion matrix obtained from the

309    dataset used for the validation is presented (Table 3). The precision and recall for "F-

310    100" and "F-other" were 0.91 and 0.61, and 0.87 and 0.63, respectively. Because we

311    only counted classes in which the probability was larger than 0.8, misses were high,

312    especially for "F-other."

313

314    3.2 Experiment I

315        The difference $\hat{d}_n$ between the estimated and measured values obtained for each fish

316    species is shown on the y-axis and the measured value on the x-axis of Fig. 5a–f. 99%,

317    98%, 92%, 96%, and 97% of $\hat{d}_n$ were included within 5% intervals for mackerels,

318    Japanese jack mackerels, Japanese sardines, red barracudas, and bullet tuna, respectively.

319    In addition, all the bias $\hat{B}$ was within $\pm$ 3%, where a maximum and minimum bias was

320    1.15 and -2.77, respectively. The smallest bias was 0.14 for mackerels.

321

322    3.3 Experiment II

323        Detected rates of fish that were identified as both "F-100" and "F-other" decreased as

324    the true number of fish in the box $N$ increased (Fig. 6a, b). The rates rapidly decreased

18

325    for the "F-100" individual, although the rates of "F-other" gradually decreased, and

326    some of them were over 100% because some parts of the fish body were detected twice.

327    The relationships between $\widehat{D}_{1,N,i}, \widehat{D}_{2,N,i}$, and the detected rates are shown in Fig. 7.

328    The value of $\widehat{D}_{2,N,i}$ took on larger negative values as the detected rate decreased,

329    although it did not decrease when only "F-100" was extracted. The variances of $\widehat{D}_{1,N,i}$

330    increased when the detected rates were low. The parameters from the simple regression

331    analysis are shown in Fig. 7. The estimates showed that $\widehat{D}_{1,N,i}$ was not affected by the

332    detected rates, and mean differences (i.e., regression line) were included within, at most,

333    $\pm$ 1.5%, where the detected rate (i.e., the independent variable $x$) was changed from 0 to

334    100. Examples of predicted labels and mask areas for every shot of Japanese jack

335    mackerel, where $N = 14$, are shown in Fig. 8.

336

337    **4. Discussion**

338    Few studies have used deep learning to estimate the body length of fish using 2D

339    images. Previous studies have estimated the total length from the size of a fish's head

340    (Álvarez-Ellacuría et al., 2020), estimated average body size from the weight of fish in a

341    fish box (Palmer et al., 2022), estimated size directly by deep learning (Ovalle et al.,

342    2022). In all these studies, the camera was fixed and was capable of capturing images of

19

343    the fish directly below. The distance between the camera and the target fish also does

344    not change from one image to the next; therefore, the value of cm/pixel remains

345    constant from image to image. In our study, we showed that total length estimation is

346    possible even in situations in which the camera cannot be fixed (e.g., no place or high

347    cost to setup, no power supply, and poor security). This will increase opportunities to

348    apply fish length estimation through deep learning, regardless of whether the camera

349    can be fixed.

350    The detected rate of fish detected as "F-100" decreased rapidly and eventually reached

351    zero (Fig. 6a, b). This is because as the number of fish in the fish box increases, the

352    probability of occlusion increases. This was also supported by the results shown in Fig.

353    7. Here, the difference between the average length of individuals detected as "F-100"

354    and their actual mean length (i.e., the regression lines) was within ±1.5%, regardless of

355    the detected rate. If the individuals detected as "F-100" included individuals that were

356    underlain by other fish (i.e., "F-other"), the difference would have increased to a

357    negative value as the detected rate decreased. In fact, the mean difference was negative

358    only in the results that included both individuals detected as "F-100" and "F-other."

359    This was because the total length was calculated from individuals smaller than the

360    actual size as "F-other." These results suggest that fish detected as "F-100" do not

361  include occluded fish, regardless of the degree of occlusion. In other words, when the

362  total length was estimated only from individual fish detected as non-occluded,

363  regardless of how many fish were placed in the fish box and photographed, the

364  difference from the actual mean length would be within ±1.5% on average, or all would

365  be processed as occluded individuals, which would result in an incorrect estimation of

366  the total length composition. Because it is difficult to completely control the number of

367  fish in the fish box, this fact simplifies the procedure by which a photographer can make

368  length estimates using the proposed method.

369  There is a trade-off between recall and precision. To obtain an accurate total length

370  composition of the target fish, the recall can be low, but the precision must be high. In

371  this study, when the probability was set to 0.8, 8% (1–0.92) of the fish were

372  misclassified as "F-100". To obtain the total length composition, it is desirable to

373  maintain a high probability (e.g., score = 0.8). For example, for species that are

374  abundant and frequently photographed, a more stringent probability value would

375  provide accurate total length composition from many images. For species with many

376  stock management constraints, an accurate total length composition is important;

377  therefore, a higher probability is desirable. However, when 99% of the fish body is

378  visible, it is necessary to label the *exposure* as "F-other," but this is often difficult for

379  annotators to determine. It is very important to determine the judgment criterion of the

380  label in advance to detect "F-100" accurately, even if the probability is high.

381     If total length estimation can be performed using ToroCam and deep learning

382  techniques, there are two advantages to conducting total length measurements. The first

383  is a reduction in the work time. The total length of catch at fishing ports or the fish

384  market is collected manually by measurers. In most cases, the data were recorded by

385  punching holes in the measurement paper using an eyeleteer. After returning to the

386  laboratory, the positions of the holes in the measurement paper are read and converted

387  to numerical values and entered in a spreadsheet such as Excel. However, the proposed

388  method can significantly reduce the labor hours of the measurers because the total

389  length estimate is output simply by extracting data from a smartphone, and a deep

390  learning model makes the inference. The authors experimentally calculated that the time

391  required to automatically output numerical values from a fish image to a spreadsheet

392  using this method was only 9% of the time required to read the values from the holes in

393  the measurement sheet and type them into a spreadsheet.

394     The second advantage is that the time required for measurement is reduced, allowing

395  the collection of the total length information for a greater number of fish. In this study,

396  we showed that we could estimate six fish species with an accuracy of $\pm$ 3%, and we

397     expect that fish species not considered in this study can also be detected as "F-100" in

398     the same manner if the sample size of the training data is increased. This will update the

399     conventional stock assessment methods by incorporating length-based models (Hordyk

400     et al., 2014), which are expected to improve the accuracy of stock assessments for

401     several fish species.

402     When a measurer takes a fish image manually with ToroCam, slight movements of

403     the hand would change the position of the rectangle, making it time-consuming for

404     sensitive users to align the four corners of the fish box precisely. However, in practice,

405     the size of the fish box was sufficiently large relative to the image size, and slight

406     shaking did not have a significant effect on the bias. The act of aligning the rectangle

407     with the four corners of the fish box also played a role in enhancing the effect of

408     shooting the fish directly above. It would be good to confirm that the difference

409     between the estimates and true length does not vary greatly from person to person

410     before actually taking the images. As a future issue, linking the transmission function of

411     the smartphone with ToroCam would further facilitate the process of obtaining length

412     compositions.

413     Future studies should include a combination of fish species identifications. There

414     have been reports of the impact of individual occlusions on fish species classification

23

415 (Ovalle et al., 2022). In a previous study, the degree of occlusion was manually

416 separated; however, it was suggested that learning the occlusion itself would reduce this

417 effort. Although our study did not classify fish species, it is possible to combine fish

418 species classification models, which will be the next step. In such cases, the burden on

419 measurers should be further reduced.

420

421 **5. Conclusion**

422 　In this study, a smartphone application, ToroCam, and deep learning method were

423 used to detect fish that were not occluded. The performance was within $\pm$ 1.5%, and

424 reliable total length estimates were obtained. Increasing the amount of stock used within

425 the maximum sustainable yield (MSY) is an international goal described in the

426 Sustainable Development Goals (SDGs). The results of this study will contribute to the

427 calculation of MSY for several fish species that conventional EMS cannot target,

428 because the total length composition is available simply by photographing.

429

433  would also like to thank Mr. Kazuharu Iwasaki for providing the original code for

434  carrying out the Mask R-CNN. Dr. Yutaka Osada worked with us to acquire images of

435  the fish and provided helpful comments early in the analysis. Moreover, the authors

436  would like to thank Editage (www.editage.com) for English language editing and TTPM

437  Inc. for providing high-quality annotated data. This study was funded by the Fisheries

438  Agency of the Ministry of Agriculture, Forestry and Fisheries of Japan.

439

440  **Figures and Table captions**

441  Fig. 1

442  Positions of the Matsuura fishing market, Odawara and Toyohama fishing ports, and

443  southern Hokkaido area where Kaiyo-maru No.5 and Horyo-maru No. 18 were operated.

444

445  Fig. 2

446  Examples of images that are used for this analysis; (a) images from a camera set up near

447  a conveyor belt at Matsuura fishing market, (b) Odawara fishing port, (c) Horyo-maru

448  No. 18, (d) images from a camera set up near a sorting table on Kaiyo-maru No. 5, (e) a

449  manually photographed image at Matsuura fishing market, (f) Odawara fishing port and

450  (g) Toyohama fishing port.

451

452  Fig. 3

453  (a) Examples of *exposure* where mask area is not drawn to keep visibility of original

454  fish body. An individual fish annotated as "F-100" is shown in the yellow colored

455  bounded boxes and text, whereas those of "F-other" are shown in white. (b) Examples

456  of the label "Non-target" where the polygon area includes parts of fish body but cannot

457  identify their species and number.

458

459  Fig. 4

460  Example of an actual screen through ToroCam. The red rectangle is shown in the

461  camera screen to manually align the corners of rectangle with the corners of the fish box.

462

463  Fig. 5

464  The relative difference $\hat{d}_n$ (%) between the estimates and observed total length (TL) for

465  each fish species: (a) Mackerels, (b) Japanese jack mackerels, (c) Japanese sardines, (d)

466  Red barracudas and (e) Bullet tunas. Black circles indicate the values of the difference.

467  The calculated relative bias $\hat{B}$ and sample size are also shown in legends. Black break

468  line shows ±5%.

469

470    Fig. 6

471    Relationship between the detected rates and the true number of fish in the box ($N$); (a)

472    Mackerels and (b) Japanese jack mackerels. Black circles indicate values of the rate

473    calculated from only non-occluded individuals ($\hat{R}_{1,N,i}$), and the white squares indicate

474    both non-occluded and occluded ($\hat{R}_{2,N,i}$) each shot $i$. Black bold and break lines indicate

475    those mean values at each $i$ and 50% of the detected rate, respectively.

476

477    Fig. 7

478    Relationship between the detected rates and mean difference of estimated and observed

479    total length composition; (a) Mackerels and (b) Japanese jack mackerels. Black circles

480    indicate values of the difference calculated from only non-occluded individuals ($\hat{D}_{1,N,i}$)

481    and the white squares indicate both non-occluded and occluded ($\hat{D}_{2,N,i}$) each shot $i$.

482    Black break line shows ±5%.

483

484    Fig. 8

485    Example of predictions for Japanese jack mackerels ($N$=14). The positions of the fish

486    body is different because it was photographed three times (a: first time, b: second time,

27

487    c: third time) after being stirred by hand. In the third shot, the image is slightly blurred

488    due to the camera shaking.

489

490    Table 1

491    The number of images by image size and obtained places that were used in this analysis.

492

493    Table 2

494    The number of images and individuals or objects that are classified as "F-100", "F-

495    other", and "Non-target" by obtained places. These are used for training data.

496

497    Table 3

498    Confusion matrix and calculated precision and recall each class. *Missing* indicates that

499    the fish were there but missed as they were considered as background. *Misdetection on*

500    *the background* indicates that fish were detected as present, even though there were no

501    fish in the background.

502

503    Supplementary Table 1

504    Although identification of species names using deep learning techniques was not carried

505    out in this study, the species name had been annotated; otherwise, their genus, family, or

506    order names were annotated if the species name of the fish could not be identified (e.g.,

507    their body position was not appropriate to identify). Species for which species

508    identification was not possible for some reason and that were similar in appearance

509    were treated as the same group of fish species, even if their genus, family, and order

510    were different.

511

512

513    **References**

514    Álvarez-Ellacuría, A., Palmer, M., Catalán, I.A., Lisani, J.L., 2020. Image-based,

515        unsupervised estimation of fish size from commercial landings using deep learning.

516        ICES Journal of Marine Science 77, 1330–1339.

517        https://doi.org/10.1093/icesjms/fsz216

518    Bartholomew, D.C., Mangel, J.C., Alfaro-Shigueto, J., Pingo, S., Jimenez, A., Godley,

519        B.J., 2018. Remote electronic monitoring as a potential alternative to on-board

520        observers    in    small-scale    fisheries.    Biol    Conserv    219,    35–45.

521        https://doi.org/10.1016/j.biocon.2018.01.003

522    He, K., Gkioxari, G., Dollár, P., Girshick, R., 2017. Mask R-CNN. Proceedings of the

523        IEEE international conference on computer vision 2961–2969.

524    Hordyk, A., Ono, K., Valencia, S., Loneragan, N., Prince, J., 2014. A novel length-based

525        empirical estimation method of spawning potential ratio (SPR), and tests of its

526        performance, for small-scale, data-poor fisheries, in: ICES Journal of Marine

527        Science.      Oxford      University      Press,      pp.      217–231.

528        https://doi.org/10.1093/icesjms/fsu004

529    Ichinokawa, M., Okamura, H., Kurota, H., 2017. The status of Japanese fisheries

530        relative to fisheries around the world. ICES Journal of Marine Science.

531        https://doi.org/10.1093/icesjms/fsx002

532    Lee, H.H., Piner, K.R., Methot, R.D., Maunder, M.N., 2014. Use of likelihood profiling

533        over a global scaling parameter to structure the population dynamics model: AN

534        example using blue marlin in the Pacific Ocean. Fish Res 158, 138–146.

535        https://doi.org/10.1016/j.fishres.2013.12.017

536    Lekunberri, X., Ruiz, J., Quincoces, I., Dornaika, F., Arganda-Carreras, I., Fernandes,

537        J.A., 2022. Identification and measurement of tropical tuna species in purse seiner

538        catches   using   computer   vision   and   deep   learning.   Ecol   Inform   67.

539        https://doi.org/10.1016/j.ecoinf.2021.101495

540    Lu, Y.C., Tung, C., Kuo, Y.F., 2020. Identifying the species of harvested tuna and

541 billfish using deep convolutional neural networks. ICES Journal of Marine Science

542 77, 1318–1329. https://doi.org/10.1093/icesjms/fsz089

543 Ovalle, J.C., Vilas, C., Antelo, L.T., 2022. On the use of deep learning for fish species

544 recognition and quantification on board fishing vessels. Mar Policy 139.

545 https://doi.org/10.1016/j.marpol.2022.105015

546 Palmer, M., Álvarez-Ellacuría, A., Moltó, V., Catalán, I.A., 2022. Automatic,

547 operational, high-resolution monitoring of fish length and catch numbers from

548 landings using deep learning. Fish Res 246.

549 https://doi.org/10.1016/j.fishres.2021.106166

550 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury Google, J., Chanan, G., Killeen, T.,

551 Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Xamla, A.K., Yang, E., Devito,

552 Z., Raison Nabla, M., Tejani, A., Chilamkurthy, S., Ai, Q., Steiner, B., Facebook,

553 L.F., Facebook, J.B., Chintala, S., 2019. PyTorch: An Imperative Style, High-

554 Performance Deep Learning Library.

555 Piner, K.R., Lee, H.H., Maunder, M.N., 2016. Evaluation of using random-at-length

556 observations and an equilibrium approximation of the population age structure in

557 fitting the von Bertalanffy growth function. Fish Res 180, 128–137.

558 https://doi.org/10.1016/j.fishres.2015.05.024

559    Privitera-Johnson, K.M., Punt, A.E., 2020. A review of approaches to quantifying

560        uncertainty in fisheries stock assessments. Fish Res 226.

561        https://doi.org/10.1016/j.fishres.2020.105503

562    Shibata, Y., Nagao, J., Narimatsu, Y., Morikawa, E., Suzuki, Y., Tokioka, S., Yamada,

563        M., Kakehi, S., Okamura, H., 2021. Estimating the maximum sustainable yield of

564        snow crab (Chionoecetes opilio) off Tohoku, Japan via a state-space stock

565        assessment model with time-varying natural mortality. Popul Ecol 63, 41–60.

566        https://doi.org/10.1002/1438-390X.12068

567    van Essen, R., Mencarelli, A., van Helmond, A., Nguyen, L., Batsleer, J., Poos, J.J.,

568        Kootstra, G., 2021. Automatic discard registration in cluttered environments using

569        deep learning and object tracking: class imbalance, occlusion, and a comparison to

570        human review. ICES Journal of Marine Science.

571        https://doi.org/10.1093/icesjms/fsab233

572    van Helmond, A.T.M., Mortensen, L.O., Plet-Hansen, K.S., Ulrich, C., Needle, C.L.,

573        Oesterwind, D., Kindt-Larsen, L., Catchpole, T., Mangi, S., Zimmermann, C.,

574        Olesen, H.J., Bailey, N., Bergsson, H., Dalskov, J., Elson, J., Hosken, M., Peterson,

575        L., McElderry, H., Ruiz, J., Pierre, J.P., Dykstra, C., Poos, J.J., 2020. Electronic

576        monitoring in fisheries: Lessons from global experiences and future opportunities.

577    Fish and Fisheries 21, 162–189. https://doi.org/10.1111/faf.12425

578    Wang, S.P., Maunder, M.N., Piner, K.R., Aires-da-Silva, A., Lee, H.H., 2014. Evaluation

579    of virgin recruitment profiling as a diagnostic for selectivity curve structure in

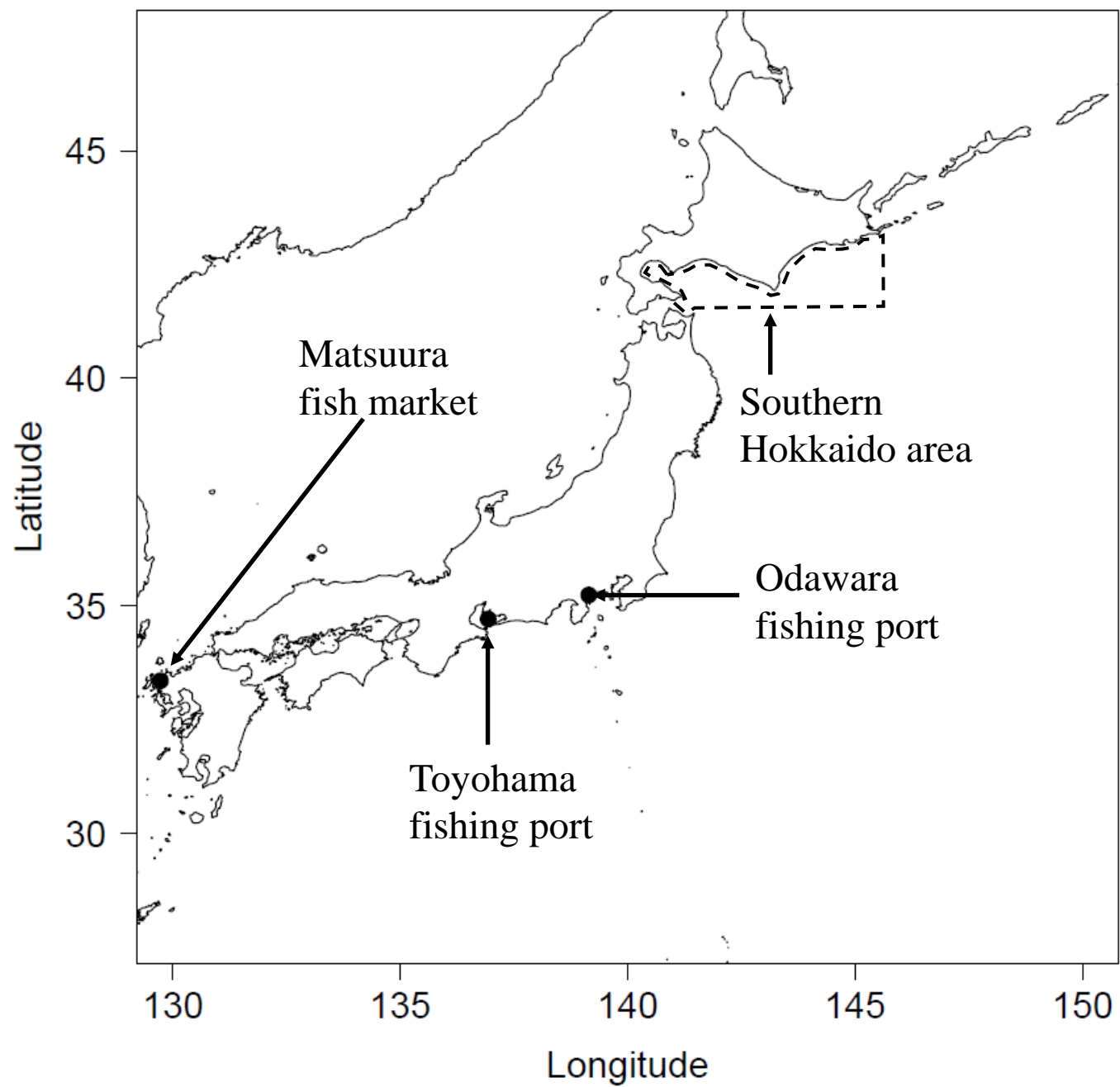580    integrated    stock    assessment    models.    Fish    Res    158,    158–164.

581    https://doi.org/10.1016/j.fishres.2013.12.009

582

Fig. 1

# Fig. 2

(a) On conveyor at Matsuura fishing market

# Fig. 2

(b) On conveyor at Odawara fishing port

# Fig. 2

(c) On conveyor at Horyo-maru No. 18

# Fig. 2

(d) On table at Kaiyo-maru No. 5

# Fig. 2

(e) Manually photographed at Matsuura fishing market

# Fig. 2

(f) Manually photographed at Odawara fishing port

# Fig. 2

(g) Manually photographed at Toyohama fishing port

# Fig. 3

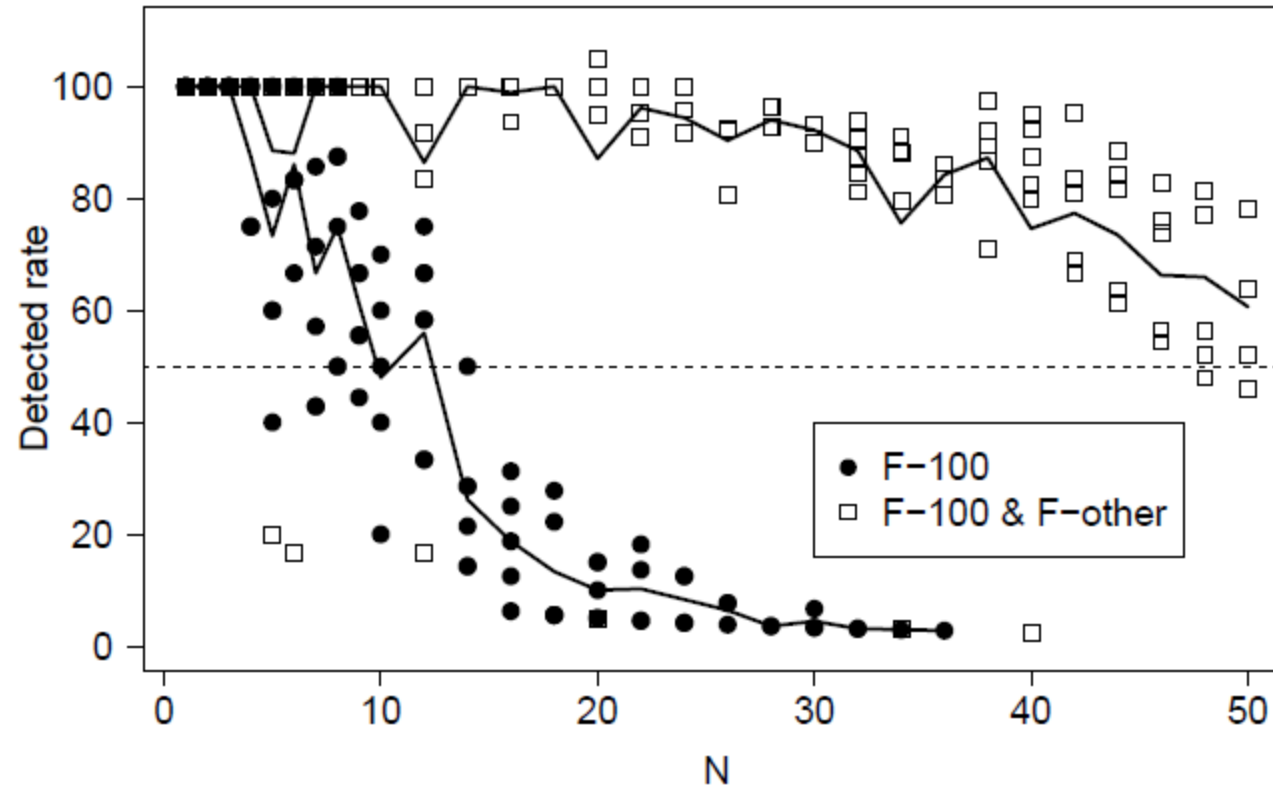# Fig. 3

(b) Example of "non-target".



non-target

Fig. 4

# Fig. 5

## (a) Mackerels

# Fig. 5

(b) Japanese jack mackerels

# Fig. 5

(c) Japanese sardines

# Fig. 5

(d) Red barracudas

# Fig. 5

(e) Bullet tunas

# Fig. 6

## (a) Mackerels

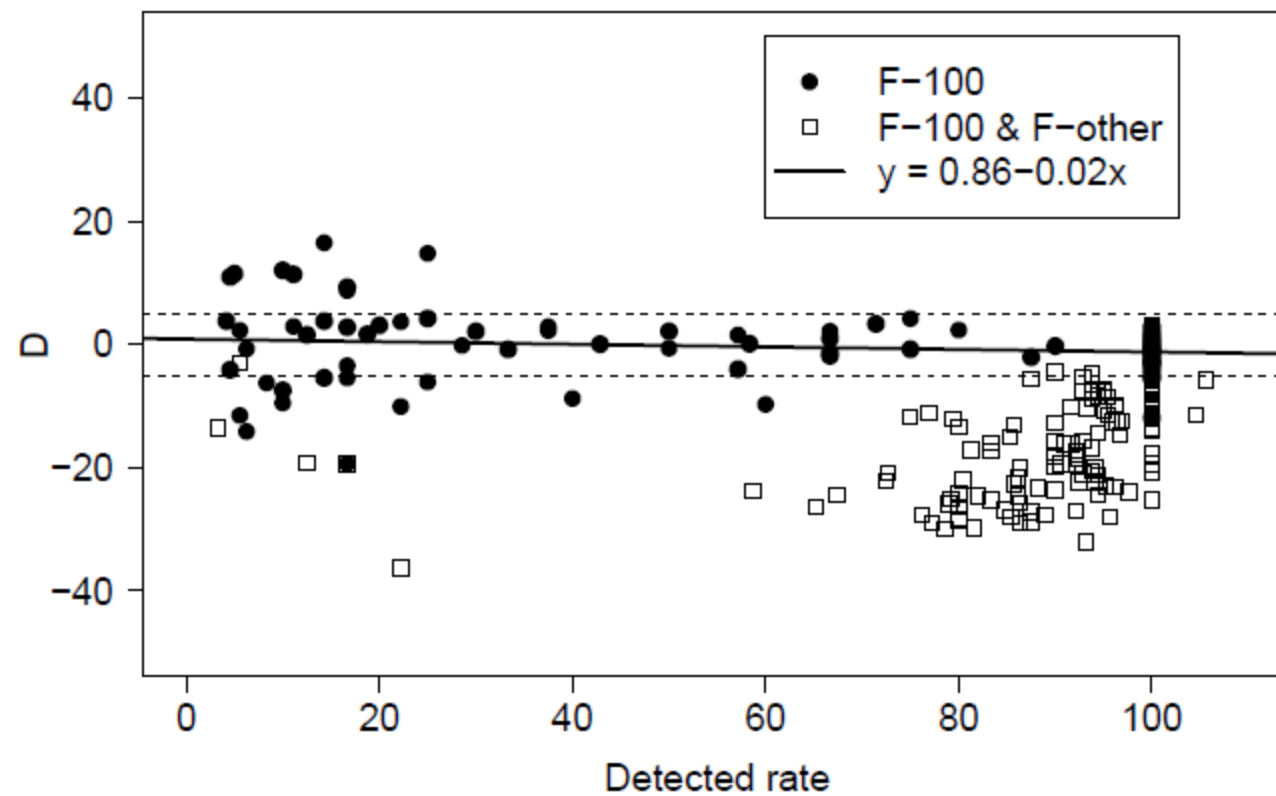# Fig. 6

(b) Japanese jack mackerels

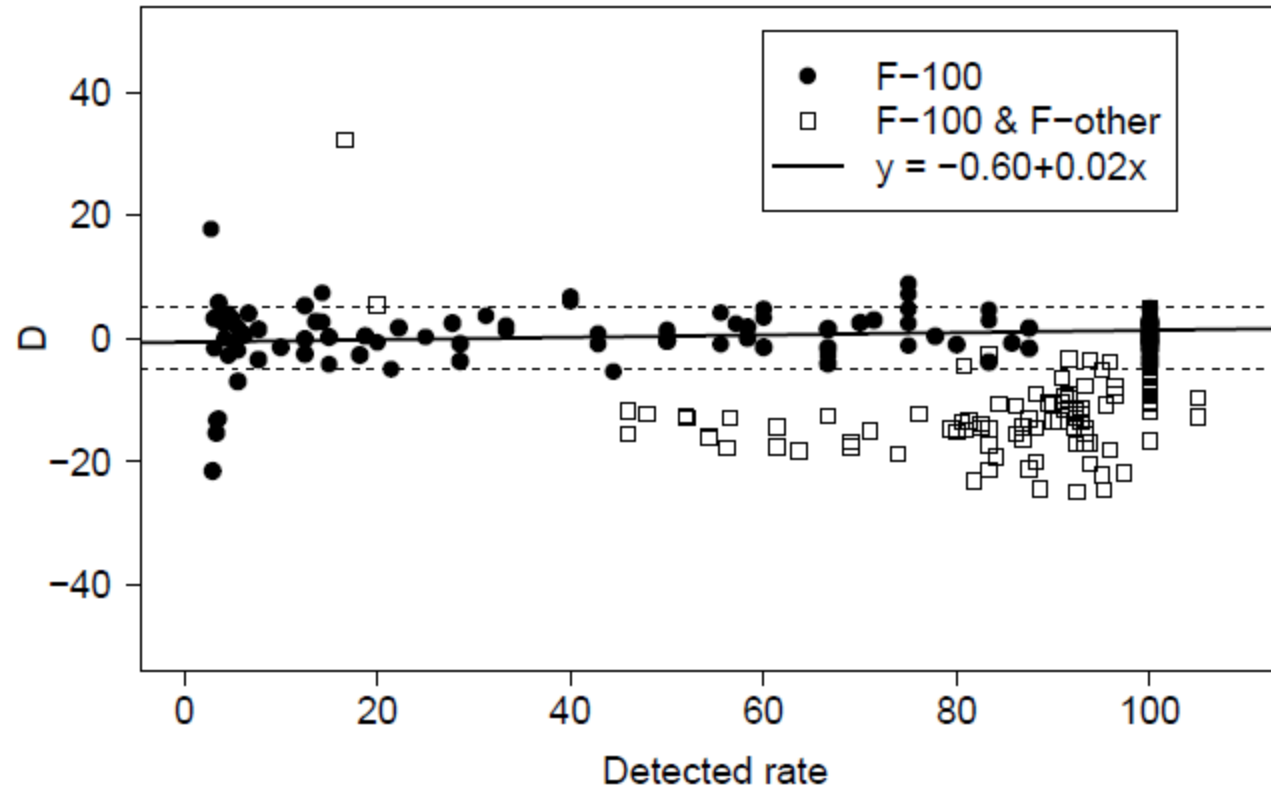# Fig. 7

## (a) Mackerels

# Fig. 7
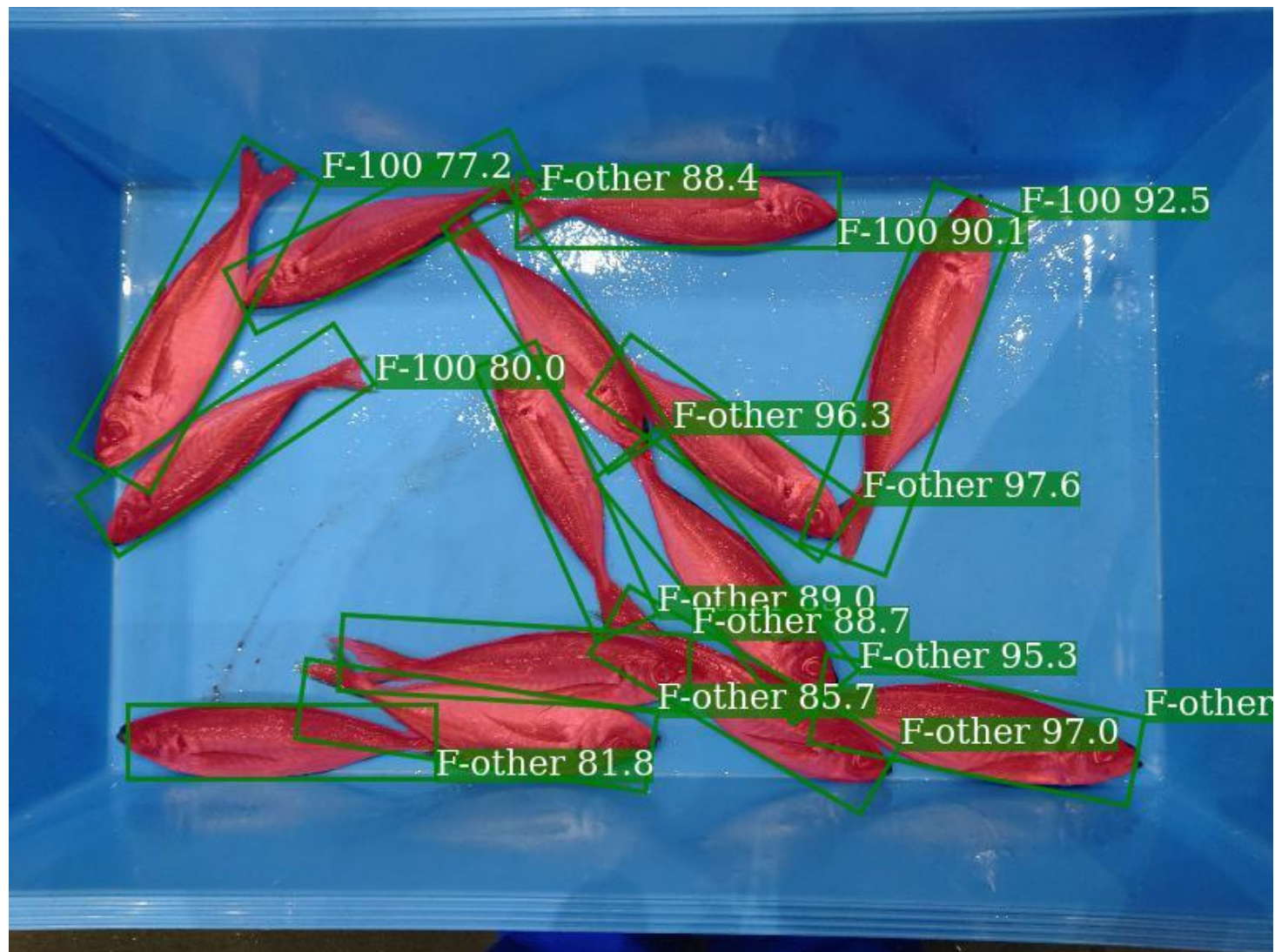
(b) Japanese jack mackerels
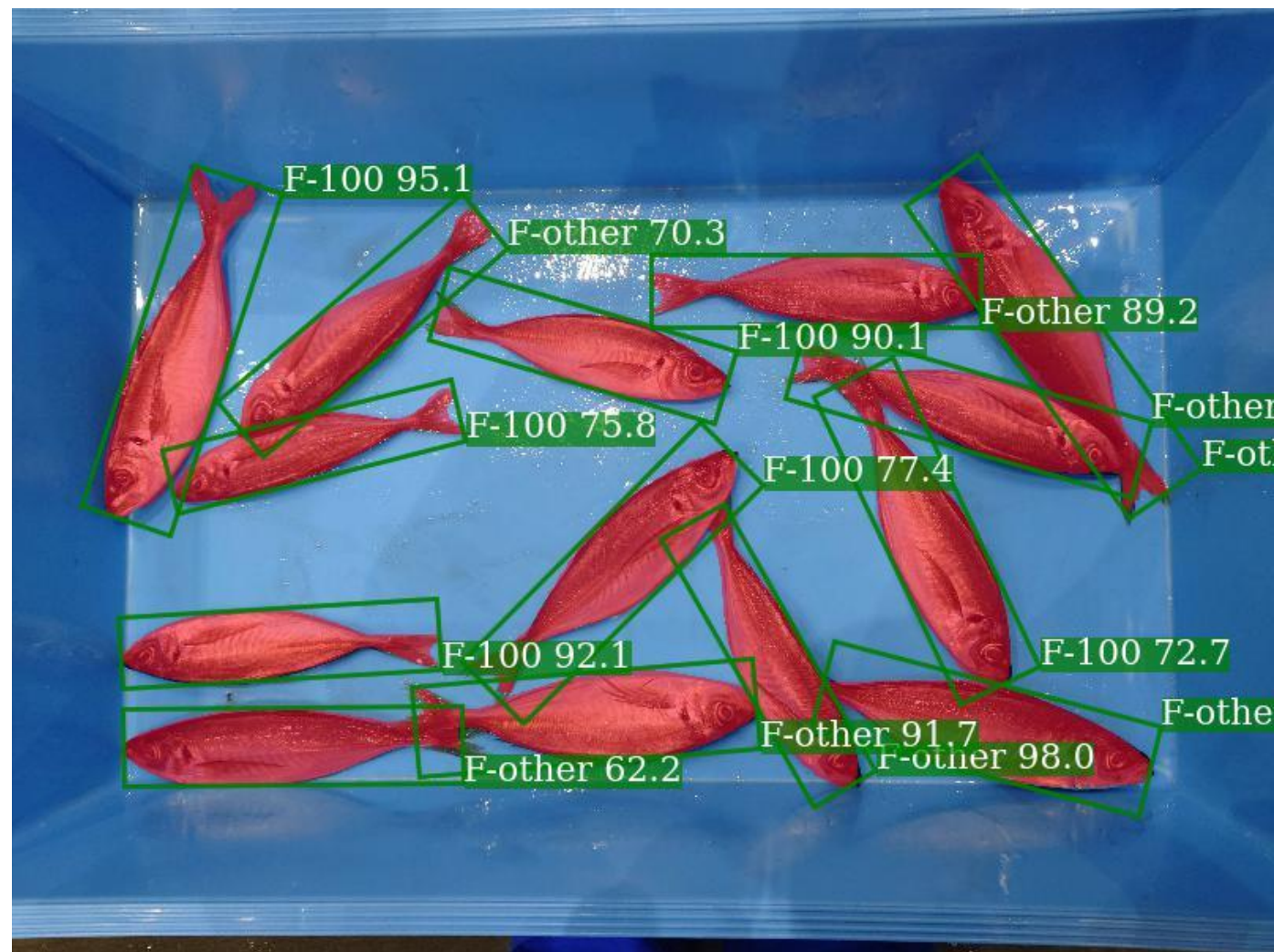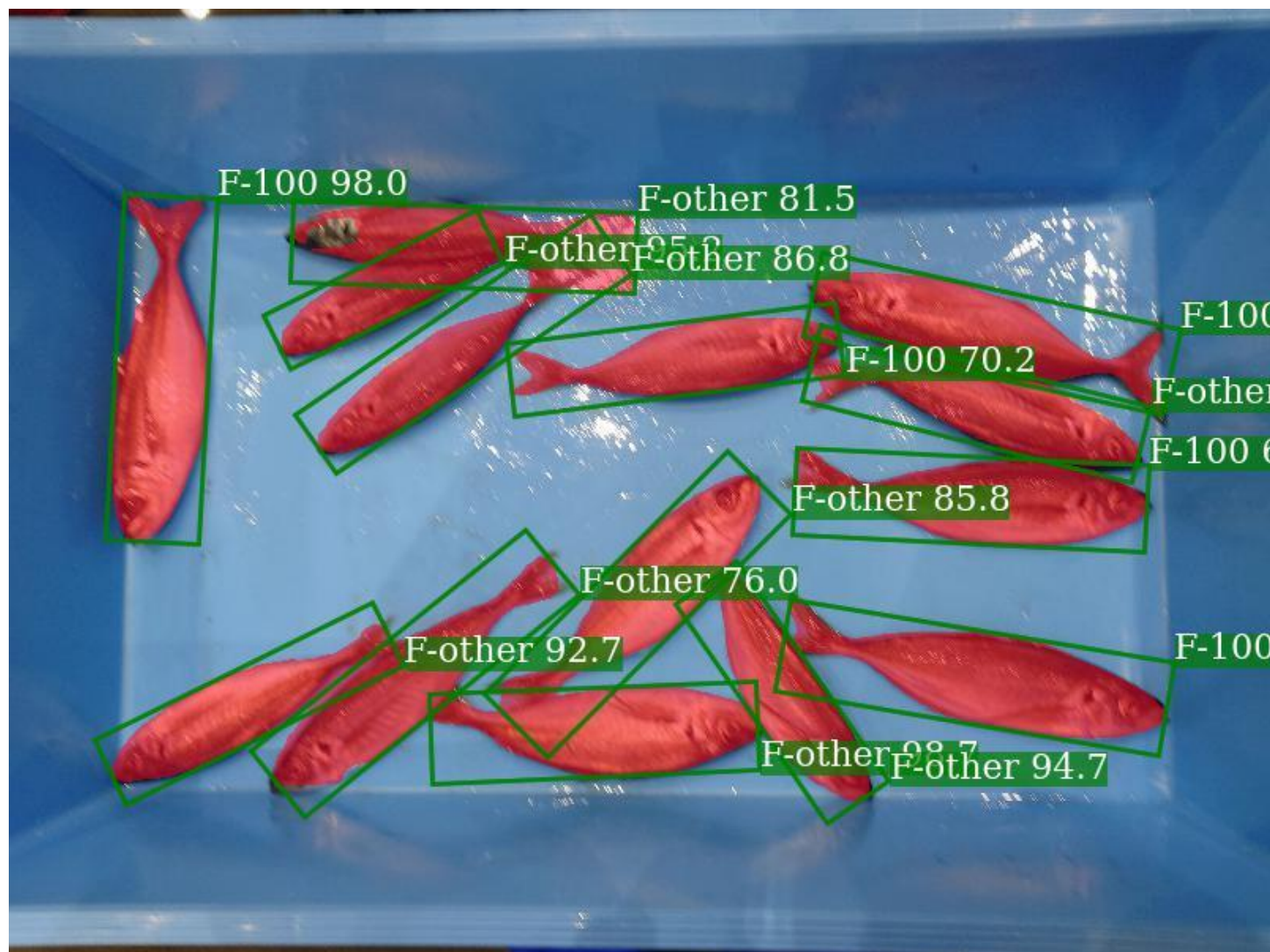
# Fig. 8

(a)

Fig. 8

(b)

# Fig. 8

(c)

Table. 1

| | Matuura fish market | Odawara fishing port | Toyohama fishing port | Kaiyo-maru No. 5 | Horyo-maru No.18 | Total |
|---|---|---|---|---|---|---|
| 4800x3200 | 119 | 3,891 | | | | 4,010 |
| 2048x1536 | | 1,769 | | | | 1,769 |
| 2704x1520 | | 1,229 | | | | 1,229 |
| 1920x1080 | | | 360 | 55 | | 415 |
| 5184x3888 | 193 | 115 | 80 | | | 388 |
| 960x1080 | | | | 41 | 61 | 102 |
| 3072x1728 | | | 101 | | | 101 |
| 4608x3456 | | | 40 | | | 40 |
| 424x480 | | | | | 27 | 27 |
| 640x720 | | | | | 6 | 6 |
| Total | 312 | 7,004 | 581 | 96 | 94 | 8,087 |

Table 2

| | Matuura fish market | Odawara fishing port | Toyohama fishing port | Kaiyo-maru No. 5 | Horyo-maru No.18 | Total |
|---|---|---|---|---|---|---|
| Images | 312 | 7,004 | 581 | 96 | 94 | 8,087 |
| F-100 | 2,407 | 17,468 | 571 | 240 | 159 | 20,845 |
| F-other | 4,675 | 39,753 | 6,504 | 1,209 | 3,175 | 55,316 |
| Non-target | 0 | 1,535 | 149 | 28 | 76 | 1,788 |

Table 3

| | | True class | | | | Recall |
|---|---|---|---|---|---|---|
| | | F-100 | F-other | Non-target | missing | |
| Predicted class | F-100 | 1,254 | 221 | 0 | 570 | 0.61 |
| | F-other | 109 | 3,405 | 0 | 1,924 | 0.63 |
| | Non-target | 0 | 0 | 125 | 41 | 0.75 |
| | misdetection on the background | 12 | 274 | 4 | | |
| | Precision | 0.91 | 0.87 | 0.97 | | |