

Using Machine Vision to Estimate Fish Length from Images using Regional Convolutional Neural Networks

Graham G. Monkman ^{a,*}, Kieran Hyder^{d,e}, Michel J. Kaiser ^c, Franck P. Vidal ^b

^a School of Ocean Sciences, Bangor University, Menai Bridge, Anglesey LL59 5AB, United Kingdom

^b School of Computer Science, Bangor University, Dean Street, Bangor LL57 1UT, United Kingdom

^c The Lyell Centre, Institute of Life and Earth Sciences (ILES), School of Energy, Geoscience, Infrastructure and Society, Heriot-Watt University, Riccarton, Edinburgh EH14 4AS, United Kingdom

^d Centre for Environment, Fisheries & Aquaculture Science, Pakefield Road, Lowestoft, Suffolk NR33 0HT, United Kingdom

^e School of Environmental Sciences, University of East Anglia, Norwich Research Park, Norwich, Norfolk NR4 7TJ, United Kingdom. Tel. +44 (0)1502 524501

* Corresponding author at: School of Ocean Sciences, Bangor University, Menai Bridge, Anglesey LL59 5AB, United Kingdom. Tel.: + 44 (0)1248 382842.

Email addresses: gmonkman@mistymountains.biz (G.G. Monkman); m.kaiser@hw.ac.uk (M.J. Kaiser), kieran.hyder@cefas.co.uk (K. Hyder); f.vidal@bangor.ac.uk (F.P. Vidal)

ORCID

G.G. Monkman <http://orcid.org/0000-0002-5645-1834>, K. Hyder <http://orcid.org/0000-0003-1428-5679>, M. J. Kaiser <http://orcid.org/0000-0001-8782-3621>, F.P. Vidal <https://orcid.org/0000-0002-2768-4524>

Keywords fiducial marker, photogrammetry, European sea bass, regional convolutional neural network, CNN, videogrammetry

30 **Summary**

- 31 1 An image can encode date-time, location and camera information as metadata and
32 implicitly encodes species information and data on human activity, e.g. the size
33 distribution of fish removals. Accurate length estimates can be made from images
34 using a fiducial marker however, their manual extraction is time consuming and
35 estimates are inaccurate without control over the imaging system. This article
36 presents a methodology which uses machine vision to estimate the total length (TL)
37 of a fusiform fish (European sea bass).
- 38 2 Three regional convolutional neural networks (R-CNN) were trained from public
39 images. Images of European sea bass were captured with a fiducial marker with 3
40 non-specialist cameras. Images were undistorted using the intrinsic lens properties
41 calculated for the camera in OpenCV, then TL was estimated using machine vision
42 (MV) to detect both marker and subject. MV performance was evaluated for the three
43 R-CNNs under downsampling and rotation of the captured images.
- 44 3 Each R-CNN accurately predicted the location of fish in test images (mean
45 intersection over union, 93%) and estimates of TL were accurate, with percent mean
46 bias error (%MBE [95% CIs]) = 2.2% [2.0, 2.4]). Detections were robust to
47 horizontal flipping and downsampling. TL estimates at absolute image rotations > 20°
48 became increasingly inaccurate but %MBE [95% CIs] was reduced to -0.1% [-0.2,
49 0.1] using machine learning to remove outliers and model bias.
- 50 4 Machine vision can classify and derive measurements of species from images
51 without specialist equipment. It is anticipated that ecological researchers and
52 managers will make increasing use of MV where image data is collected (e.g. in
53 remote electronic monitoring, virtual observations, wildlife surveys and

54 morphometrics) and MV will be of particular utility where large volumes of image
55 data will be gathered.

56 **1 Introduction**

57 Currently only a small proportion of the world's marine stocks are sufficiently data rich to
58 enable formal stock assessments to be performed, hence most marine fisheries are data poor
59 (Costello et al., 2012; Ricard et al., 2012). This is in spite of legislation (e.g. European
60 Commission Decision 2008/56/EC) which requires marine stocks to be exploited sustainably
61 and managed with consideration of their associated ecosystems. The potential for commercial
62 fisheries to negatively impact stocks and ecosystems is accepted, but recreational fishing can
63 also negatively impact fisheries and associated ecosystems (reviews Lewin et al., 2006;
64 Radford et al., 2018). Marine recreational fisheries in particular can lack current and historical
65 data even in developed countries and monitoring of the sector is poor (ICES, 2017; Hyder et
66 al., 2018).

67 Fisheries assessments have survey phases in which a metrological measurement of the target
68 species occurs (National Research Council, 2006; ICES, 2012). In both commercial and
69 recreational fisheries, measurement has traditionally involved observations by researchers,
70 fisheries managers or the fishers themselves. Observer costs are high in commercial monitoring
71 (e.g. Needle et al., 2015) and in the assessment of recreational fisheries (pers. observ. KH).
72 Hence, there has been an increasing interest in remote electronic monitoring (REM) (e.g. White
73 et al., 2006, Chang et al., 2010, Hold et al., 2015, Bartholomew et al., 2018). Videogrammetry
74 and photogrammetry (hereafter, photogrammetry) are becoming increasingly commonplace in
75 non-destructive observational marine research (e.g. Dunbrack, 2006, Deakos, 2010).

76 The use of REM and related approaches is likely to increase as camera technology improves
77 and equipment costs fall (reviews Struthers et al., 2015, Bicknell et al., 2016). Photogrammetry
78 can provide considerable savings when compared to observers (Chang et al., 2010; National

79 Oceanic and Atmospheric Administration, 2015). Capturing images produces vast volumes of
80 data which is time consuming to process (e.g. Needle et al., 2015, van Helmond et al., 2017).
81 This problem can be alleviated by using motion detection algorithm(s) to extract salient frames
82 from videos (e.g. Weinstein, 2015), but the extracted frames still require manual processing.
83 Object detection with machine vision (MV) could be used to automate the extraction of data
84 from images. Historically, MV has been used to analyse images which have been captured
85 under controlled conditions (e.g. fixed cameras, backgrounds and lighting). This control makes
86 the isolation of the subject from the background (segmentation) much easier, allowing
87 computationally inexpensive techniques to be applied, e.g. using optical flow (Zion et al., 2007;
88 Spampinato et al., 2010; Hsiao et al., 2014) and segmentation by pixel properties (e.g. White
89 et al., 2006, Jeong et al., 2013).

90 To date, photogrammetry has typically used multi-laser (e.g. Deakos, 2010, Bartholomew et
91 al., 2018) or multi-camera systems (e.g. Dunbrack, 2006, Rosen et al., 2013, Neuswanger et
92 al., 2016), but the equipment is comparatively bulky and expensive. Single camera systems and
93 a fiducial marker (i.e. an object of known scale placed in the camera's field of view) have been
94 used (Hold et al., 2015; van Helmond et al., 2017) but control of the camera model or the
95 framing of the fiducial marker and subject is usually required (e.g. Rogers, Cambiè, & Kaiser,
96 2017). Without this control, length estimates are subject to an unknown error because lenses
97 have different optical properties. The additional challenges in extracting quantitative data from
98 images taken by volunteers—or other scenarios where expensive or less portable equipment is
99 unsuitable—may explain the almost complete lack of a suitable solution. Convolution neural
100 networks (CNN) currently outperform all other methods at object detection and CNN
101 application programming interfaces (API) are now mature and stable enough to be viable
102 solutions for (merely) competent programmers to use regional CNNs (R-CNN) for object
103 detection.

104 This article explores the feasibility of using MV to automate the identification and size
105 estimation of an important species from images. The objectives are to (i) introduce the software
106 and methods to achieve length estimation with a cheap and portable fiducial marker; (ii) to
107 show that length estimates can be made with no control over the image background, lighting
108 or specialist cameras using a foreground fiducial marker; (iii) provide region of interest (RoI)
109 labelled images of the European sea bass, *Dicentrarchus labrax* (see Appendix S2 Supporting
110 Information); (iv) to compare the speed and performance of three state-of-the-art R-CNN
111 networks.

112 **2 Methods¹**

113 **2.1 Ethics**

114 European sea bass captures were made by recreational fishers and a commercial vessel as
115 part of their day-to-day activity. All reasonable measures were taken to minimise air exposure
116 time to the fish while photographs were taken. Ethical approval was granted by the Animal
117 Welfare and Ethical Review Board of Bangor University, Wales, UK.

118 **2.2 Training and validation image acquisition**

119 Training ($n = 734$) and validation ($n = 184$) images were obtained from online public sources.
120 The RoI for each image was drawn tight to the fish body, to the limits of the caudal fin tips and
121 the snout vertex (Fig. 1a). Training and inference were carried out in Tensorflow (Google,
122 2018) using transfer learning with the following pretrained R-CNNs; (i) ResNet-101 (He et al.,
123 2016), (ii) Single shot MobileNet detector (Howard et al., 2017) and (iii) NASNet (Zoph & Le,
124 2017), abbrevs. ResNet, MobileNet and NASNet respectively.

¹ Appendix S1 Supporting Information contains additional methodological detail.

125 **2.3 Fiducial marker selection and image acquisition**

126 Three ArUco fiducial markers (Garrido-Jurado et al., 2014) of side lengths 25 mm, 30 mm
127 and 50 mm were mounted on polypropylene sheets (Fig. 1b). Photographs of European sea
128 bass were taken on the shore and afloat, with the informed consent of fishers and with 3
129 different non-specialist cameras (henceforth *marker images*). Fish were posed to minimise
130 body distortion and occlusion. Fish total length (TL) was measured and recorded. The marker
131 was placed on the fish (Fig. 1c) and then photographed.

132 **2.4 Undistorting marker images**

133 Images from each camera were corrected for radial and tangential distortion with the
134 OpenCV API (OpenCV team, 2018). Lens calibration profiles were created in OpenCV for
135 each camera at each supported field of view and focal length (henceforth *undistorted images*).

136 **2.5 Length estimation**

137 An R-CNN predicts the rectangle which most accurately bounds the subject within the image
138 and then defines the detection as a rectangle defined by four vertices. Intersection over Union
139 (IoU) measures the accuracy of object localisation by comparing the area of a manually defined
140 ground truth rectangle which bounds the subject with the bounding rectangle predicted by the
141 R-CNN. Each model outputs an objectness score (*score*) which is interpreted as the probability
142 that the proposed region contains the predicted class (Ren et al., 2017).

143 When estimating TL, the pixel length of the long side of the detection rectangle approximates
144 to the TL (pixels) of the fish. The real-world length per pixel, \bar{l} was estimated from the four
145 sides of the detected ArUco marker according to, $\bar{l} = \frac{1}{n} \cdot \sum_1^n l/p_i$ where p_i is the i^{th} side length
146 in pixels, and l is the real-world side length (e.g. 50 mm). The accuracy of \bar{l} was validated
147 manually (Linear Regression, $b = 1.003$, $R^2 = 0.999$) using ImageJ (Schneider et al., 2012).
148 Mean absolute error (MAE) and mean bias error (MBE) are reported and are calculated as

149 follows, $MAE = \frac{1}{n} \cdot \sum_{i=1}^n |l_i - \hat{l}_i|$ and $MBE = \frac{1}{n} \cdot \sum_{i=1}^n l_i - \hat{l}_i$ where l_i is the i^{th} estimate of TL
150 and \hat{y}_i is the expected (i.e. actual) TL of the i^{th} element. Hence a negative bias represents an
151 underestimate of TL.

152 **2.6 Detection and length estimation with rotation, flipping and downsampling**

153 The accuracy of TL estimates under three translations were checked, these were; (i) image
154 rotation between -30° and 30° in increments of 1° ; (ii) horizontal flipping of the image by the
155 x-axis, i.e. the line $x = 0.5 \cdot width$; and (iii) image downsampling by a factor of 1.5, to a
156 minimum image height or width of 50 pixels. TL estimates for rotated images were corrected
157 based on the geometry of the detection box under increasing rotation in relation to the snout
158 and caudal vertices of the subject.

159 **2.7 Removing outliers and modelling bias**

160 NASNet R-CNN detections were split into training and test data. Training data were used to
161 identify biased outliers using an isolation forest (Liu et al., 2008; Pedregosa et al., 2011) with
162 the variables; (i) ratio of height to width of the detection, (ii) objectness score and (iii) % MBE.
163 Outliers were then removed from the training set and a gradient boost regressor (Friedman,
164 2002; Pedregosa et al., 2011) trained on the predictors (i) and (ii) above. Outliers were removed
165 from the test dataset and the gradient boost regressor model used to correct bias. Further
166 methodological details are given in Appendix S3 Supporting Information.

167 Several estimates of length measurements are reported and are listed in Table 1. Means
168 followed by square brackets or the \pm notation indicate 95% confidence intervals or standard
169 deviation respectively.

170 **3 Results**

171 For every non-transformed European sea bass image, each CNN generated region proposals
172 with objectness scores > 0.5 (with the exception of a single MobileNet score of 0.01). All

173 regional proposals were at least partially coincident with ground truth, with a minimum IoU of
174 45% (45% IoU detection shown in Fig. 1b). Negative images had no false detections under any
175 network (score mean of 0.005 ± 0.008 , $n = 30$, $\max = 0.04$).

176 Detection performance between networks was practically indistinguishable on
177 untransformed and horizontally flipped images (Table 2), hence detections were invariant to
178 horizontal flipping (IoU mean; horizontal flip, 93.2% [93.0, 93.4]; untransformed, 92.8% [92.5,
179 93.0]). This equivalence is despite the large differences in mean detection times (Table 2).
180 Nonetheless, when visualised it is apparent that the NASNet network delivered more consistent
181 object detections with no IoU outliers (Fig. 2). All single MobileNet detections had IoUs >
182 75% however, ResNet had 7 detections < 75% IoU (1.1% of all detections).

183 **3.1 Length estimates**

184 ArUco markers were consistently recognised using the OpenCV API under natural
185 conditions, with the marker successfully localised in 99.3% of untransformed images. Two
186 detection failures occurred because of over-exposure (Fig. 1e). *Corrected MV-TL* estimates had
187 a MBE of $5.9 \text{ mm} \pm 20$, compared with MBE derived from *corrected manual-TL* estimation of
188 $-0.5 \text{ mm} \pm 14.8$. *Corrected MV-TL* estimates showed consistent variance in bias across *physical*
189 *TL* (Fig. 3). On excluding TL estimates made under the noisier ResNet and MobileNet
190 networks, MBE for *corrected MV-TL* estimates was increased by 2 mm to 7.9 mm nevertheless,
191 S.D. decreased to 14.7 mm, matching the precision of manual estimates of TL (*corrected*
192 *manual-TL*).

193 *Corrected manual-TL* and *MV-TL* estimation errors tended to be less accurate and precise
194 (mean squared error, MSE) when made on the shore rather than afloat (Fig. 4, MSE; Afloat,
195 7.9; Shore, 25.9), and there was no apparent systematic bias in length estimation introduced by
196 the camera model when comparing *corrected manual-TL* estimates (which have lower variance
197 than *MV-TL* length estimates) with platform as a covariate (ANCOVA, $F_{(2, 1787)}$, $p = 0.15$).

198 Mean %MBE for *corrected manual-TL* estimates were $0.7\% \pm 4.6$, $1.1\% \pm 4.0$ and $0.7\% \pm 4.1$
199 for the GoPro Hero 5 action camera, Samsung s5690 smartphone and Fujifilm XP30 camera
200 respectively.

201 The increased %IoU outliers observed during detection with ResNet and—to a lesser
202 degree—the MobileNet single shot detector manifest as the %MBE outliers in Fig. 4. The
203 ResNet detector produced 9 of the top 10 MV associated underestimates (fully corrected
204 percent errors of -16.4% to -38.0%). These errors arose because detections followed the
205 approximate pattern observed in (Fig. 1d), with the ResNet detector occasionally truncating the
206 detection. This behaviour was not observed in the other detectors on untransformed images
207 (i.e. an image which has not been flipped, downsampled or rotated).

208 **3.2 Scale**

209 ArUco marker detection was robust to downsampling to approximately 30% of the original
210 image size (original image size, mean = 1355 by 1029 pixels, or 1.5M pixels²). ArUco markers
211 were approximately 18 pixels² at 30% of original image size and images were approximately
212 400 by 300 pixels (120k pixels²). At 30% image size the marker detection rate was 93%
213 however, this dropped to 53% at the next scaling factor of 20% (Table 3). The networks on
214 average, maintained objectiveness scores of $\sim 98\%$ at the 20% scaling factor, where the mean
215 image size was 41.4k pixels² (i.e. ~ 203 pixels²). At this image size, the average ground truth
216 RoI was 158 by 23 pixels. NASNet produced marginally more accurate TL estimates under
217 downsampling. For each network %MAE increased in increments of between 1% and 2% until
218 the downsampling factor exceeded $\sim 30\%$ (mean ground truth width = 238 pixels), after which
219 %MAE began to increase in larger increments. Each network responded similarly to
220 downsampling (Fig. 5), at 20% image size, %MAE = $9.9\% \pm 7.8$ which increased markedly to
221 $15.9\% \pm 8.4$ at 13% of the original image size at ~ 153 pixels².

222 3.3 Rotation

223 The NASNet and ResNet networks behaved similarly under image rotation (Fig. 6) and
224 detection was robust to small rotations, with over 90% of objectiveness scores greater than
225 50% at absolute rotation $\leq 20^\circ$ for the NASNet and ResNet networks. At 20° absolute rotation
226 the MobileNet network had 67% of objectiveness scores below 50%. As the absolute rotation
227 angle increased beyond $\sim 15^\circ$, NASNet and ResNet predictions of *corrected MV-TL* exceeded
228 5% %MBE however, %MBE was 2.5% for the MobileNet network (Fig. 6, absolute rotation =
229 15° , %MBE; NASNet, -5.0% [-5.3, -4.6]; ResNet, -5.3% [-5.9, -4.7]; MobileNet, 2.7% [2.2,
230 3.3]). This apparently good performance of the MobileNet CNN masks the greatly decreased
231 confidence in regional proposals under this network (score series, Fig. 6) and a corresponding
232 loss of valid detections.

233 The geometric rotation correction (variable *rotation corrected MV-TL*) did not consistently
234 decrease bias for all rotations (see Appendix S1 Supporting Information) and bias reduction
235 was only marginally improved for the NASNet and ResNet networks (1.2% and 0.5%
236 respectively) however, bias was increased for the MobileNet network (1.0%). The NASNet
237 and ResNet networks displayed a consistent hyperbolic pattern in TL estimation bias through
238 the rotation range and prediction error was consistent across rotations (Fig. 6).

239 Combining outlier removal and adjusting *rotation corrected MV-TL* per sample with the
240 trained gradient descent regressor model produced a marked reduction in %MBE across
241 rotations. This correction centred bias at $\sim 0\%$ for absolute rotations $\leq 20^\circ$ (Fig. 7; Table 4). The
242 overall improvement on applying all corrections to MV estimates following lens correction
243 only are unambiguous, with unadjusted *MV-TL* estimates of %MBE = -11.4% [-11.6, -11.2].

244 4 Discussion

245 This study introduced a methodology to estimate fish TL—a crucial measurement in multiple
246 stock assessment methods (Pauly & Morgan, 1987)—using recent advances in open-source

247 R-CNNs and associated software applications (e.g. Abadi et al., 2015, OpenCV, 2018). By
248 using these resources, it was shown that the position of an organism in an image could be
249 accurately predicted without strict control over lighting conditions or subject background. The
250 high degree of accuracy of the predicted RoI ($> 90\%$ IoU) enabled the accurate estimation of
251 TL. TL estimation was achieved without reliance on specialist cameras, multi-camera systems
252 (e.g. Dunbrack, 2006; Rosen et al., 2013) or paired lasers (e.g. Deakos, 2010, Rogers et al.,
253 2017).

254 Photographing a well-posed subject with a foreground fiducial marker is faster and more
255 convenient than measuring and recording the subject length manually (pers. observ.).
256 Possessing photographs of subjects provides a persistent record which can be used to derive
257 additional measurements, to cross check data and for validation by third parties. In volunteer
258 based research (e.g. diary surveys) additional data are typically required by research
259 programmes (e.g. GPS position, date/time, species) and these data can be automatically
260 captured at image acquisition. The potential for automatic recording of much of the required
261 data—including the onerous task of physically recording a dimension—reduces the recording
262 burden on volunteers which can improve participant retention, the volume of data submissions
263 and data quality as observed in surveys (Galesic, 2006; Hoerger, 2010).

264 **4.1 Networks**

265 Of the three networks, NASNet outperformed the ResNet-101 and MobileNet networks.
266 NASNet was particularly effective at limiting outlier detections. However, the NASNet
267 network had the slowest detection speeds of the three and was the most resource intensive.
268 During learning, NASNet had to be limited to a batch size of 1 to fit within the 6 Gb of memory
269 of the NVIDIA 1060 GTX card (configuration files are available in the Supporting
270 Information). This is unsurprising as the NASNet has many more parameters than ResNet
271 (Zoph & Le, 2017).

272 Neither ResNet nor NASNet are currently capable of performing real-time detections
273 however, MobileNet can be deployed on mobile devices. The performance of MobileNet in
274 this task was arguably better than ResNet and real time detection would be of particular benefit
275 in volunteer based data collection applications where users could be given immediate feedback
276 on the success or failure of a particular recognition task (e.g. Kumar et al., 2012, Fishbrain,
277 2018, International Game Fish Association, 2018, Stowell, 2018).

278 **4.2 Length estimation**

279 The standard fish length measurements (TL, fork length FL and standard length SL) are
280 particularly suited to estimation by R-CNN based networks because the longitudinal dimension
281 of an ideal detection corresponds with the distal extremes of the morphological features which
282 delineate these lengths. In this manuscript, TL was used to demonstrate the methodology, but
283 other measurements (including FL and SL) may be estimated by changing the RoIs defined in
284 the training and test images or using previously determined morphometric relationships (e.g.
285 Needle et al., 2015). To date, rectangular ROIs—the detection outputs of R-CNNs—have no
286 history of providing length data in fisheries assessments because R-CNNs are a recent
287 development in MV. However, our results demonstrate the accuracy which can be achieved
288 where body distortion can be controlled.

289 Where curvature of the subject body cannot be controlled, lengths have been estimated by
290 identifying depth midpoints and calculating the line bisecting these midpoints (Strachan, 1993;
291 White et al., 2006) or line fitting to subject contours (Miranda & Romero, 2017). Both
292 approaches require accurate segmentation of the subject from the background, which is
293 problematic in uncontrolled environments. Nevertheless, Tensorflow provides pretrained
294 R-CNNs capable of segmentation (He et al., 2017; Google, 2018) and real-time segmentation
295 is possible (e.g. Paszke et al., 2016). Segmentation using CNNs is a promising research area
296 for improving length estimates under body curvature as commonly observed in live shark

297 specimens, remote electronic monitoring and terrestrial species. Alternative approaches may
298 include training the network to identify sub regions of the subject which are subject to less
299 distortion (e.g. the head) and keypoint detection (e.g. Kazemi and Sullivan, 2014, Vandaele et
300 al., 2018) from which length estimates can be derived from known morphometric relationships.

301 There is evidence that using cameras instead of active observers could reduce biases with
302 the presence of observers affecting fisher behaviour (Benoît & Allard, 2009; Faunce &
303 Barbeaux, 2011). Automating length estimates can mitigate self-sampling biases, remove
304 recording errors and also reduce the perception of researchers that measurements reported by
305 non-scientists are biased or not recorded as rigorously as by trained observers or researchers
306 (Kraan et al., 2013). Any digit biases—a persistent problem in human measurement (Tarrant
307 & Manfredi, 1993)—will be negated. Increased sample throughput could reduce biases which
308 can arise through subsampling processes (e.g. Kraan et al., 2013). Images also have the benefit
309 of providing an archive of data. However, problems can occur when deploying virtual
310 observations and images or movies can suffer from poor quality, obscuring of the lens, or
311 equipment failure, which can result in missing observations (e.g. Needle et al., 2015, van
312 Helmond et al., 2017).

313 The fiducial marker deployed was particularly easy to identify in fully automated MV
314 processing pipelines and performed well as evidenced by the low bias and high detection rates.
315 Length was more accurately estimated on afloat platforms than on the shore, this is because the
316 afloat platforms provide a flat surface on which to measure and photograph the subject. Across
317 both platforms and all camera models there was a small but consistent overestimate of size
318 (mean bias error, 1.6%; 6 mm). Possible explanations include an underestimate of lens-subject
319 distance during the camera calibration process which did not account for the internal distance
320 between the lens and the glass cover of the cameras, or incorrect estimation of the parameters
321 (e.g. mean profile height) used in the length correction calculation.

322 Bias magnitude was consistent across the range of fish lengths measured (25 cm to 65 cm)
323 hence a correction could be estimated empirically during training. The model used for rotation
324 correction was successful in eliminating bias (%MBE = -0.1%), which brought the error
325 magnitude in line with methods which control the imaging conditions (Hold *et al.* 2015, 0.6%
326 *in lobster*; White *et al.* 2006, 0.3%, in halibut), use paired lasers (Deakos 2010, 0.4% in manta
327 rays) or multiple cameras (Rosen *et al.* 2013 1.0% across 3 fusiform fish species).

328 Despite bias being largely eliminated, outliers in TL estimates were observed (minimised
329 under NASNet). Without rotation, this error was largely attributable to errors arising from the
330 subject pose in the image. Parallax errors arising through depth differences across the fiducial
331 marker and the subject will be a major source of error which are typically dealt with by
332 excluding images following manual review (e.g. Deakos, 2010, Rogers *et al.*, 2017). Correction
333 for tangential deflection of MV designed fiducial markers is generally supported (e.g.
334 Bergamasco *et al.*, 2011, Garrido-Jurado *et al.*, 2014), but this is unlikely to be a consistent
335 correction (*pers. observ.*) for foreground fiducial markers where the tangential displacement of
336 the marker can differ from that of the subject.

337 **4.3 Transformations**

338 Detections and length estimations were robust to flipping and downsampling. Under
339 decreasing image size the fiducial marker was found to be the limiting factor for the automatic
340 extraction of TL. This is an intrinsic limitation of using a foreground fiducial marker where
341 increasing marker size could obscure salient features. The lowest IoU was observed on the
342 smallest European sea bass sampled, where the marker occluded a comparatively large
343 proportion of the subject (Fig. 1d). The effectiveness of the CNN under substantial
344 downsampling indicates that image sizes can be significantly reduced prior to inference to
345 improve speed and reduce memory requirements.

346 Length estimates were unbiased and acceptably precise at small degrees of rotation. The
347 bounding box under rotation generally predicted the x-coordinates of the snout and caudal
348 vertices reasonably well, particularly under the NASNet network (examples in Supporting
349 Information S4). However, the simplistic geometric model used (Appendix S1 Supporting
350 Information, 1.4.3) largely failed to adequately correct length estimates under rotation. This
351 failure is attributable to the divergence of the geometric model (detailed in Appendix S1
352 Additional Methods) from the bounding features of the subject which the CNNs “chose” under
353 rotation. In essence, the CNN detections cannot be represented by the simple geometry of a
354 rotating rectangle (Appendix S4 Supporting Information). Development of a more accurate
355 geometric correction model would be possible should the use case demand it.

356 Failure to demonstrate generalisability through all rotations poses a serious limitation in
357 some deployment scenarios. Under volunteer targeted image collection, a significant
358 proportion of subject rotations would exceed the experimental rotation limits. A common and
359 trivially implemented (but computationally expensive) approach to achieving rotation
360 invariance is the brute force repetition of detection through incremental rotations. The optimal
361 detection among all rotations is then determined by some combination of metrics, e.g.
362 objectness score, detection height to width maxima or maximal agreement with a secondary
363 detection technique such as template matching (Brunelli, 2009). In this article accurate
364 detections were achieved at absolute rotations to $\sim 15^\circ$ which suggests that steps of 15° could be
365 used to reduce the search space. However, it may prove to be more efficient to train the network
366 on incrementally rotated images. This training is relatively trivial to do and has native support
367 in most CNN APIs. Nonetheless, data on rotation invariance under rotated training images was
368 not published by Zoph and Le, (2017) and R-CNNs are not intrinsically rotation invariant,
369 hence further empirical investigation is required.

370 **4.4 Applications and Real-World Deployment**

371 **4.4.1 Other applications**

372 The choice of a foreground marker was driven by the use case, i.e. in large scale volunteer
373 based surveys and data gathering exercises. A foreground marker is cheap and portable, and
374 volunteers cannot deliberately inflate size estimates by moving the marker further away from
375 the subject as would occur with a marker on which the subject is placed. The principles of the
376 method are applicable to any type of marker (which can be detected) and multicamera systems,
377 and to any organism for which morphological estimates are made, provided sufficient volumes
378 of data are to be collected to justify technical development. Difficulties will arise in
379 unconstrained camera systems where the scale indicator is difficult to distinguish in the image,
380 such as lasers in strong sunlight (pers. observ.). None specialist markers can be segmented and
381 length estimated using machine vision, such as a standard ruler (Konovalov et al., 2017).
382 Opportunistic fiducial markers could also be segmented (e.g. human face) and used to produce
383 estimates of fish size from historical images (with reduced accuracy and precision) as has been
384 done manually to provide ecological data on some species (McClenachan, 2009, Canese and
385 Bava, 2015, Belhabib et al., 2016, Rizgalla et al., 2017).

386 The recent advances in the accuracy of MV mean that commercial length estimates may be
387 made without the need for complex and costly mechanical pre-sorting and identification may
388 be possible under occlusion by combining R-CNNs and deformable parts models and landmark
389 detection (Felzenszwalb et al., 2010, Zhang et al., 2016, Ouyang et al., 2018) and estimating
390 length from morphometric relationships.

391 **4.4.2 Real-world deployment**

392 Correction for lens distortion is critical for accurate photogrammetry as show in this article,
393 particularly with increased use of robust and waterproof action cameras (Struthers et al., 2015,
394 Rogers et al., 2017, Schmid et al., 2017, Claassens and Hodgson, 2018) which have significant
395 radial distortion. In small scale projects or where the camera model can be restricted then it

396 may be practical for images to be undistorted on an ad hoc basis, for example by using the
397 manufacturers own software. However, to deploy large scale volunteer based metrological data
398 gathering it will be necessary to build a repository of lens correction profiles for each camera
399 model where radial distortion is above an acceptable threshold. If a camera supports multiple
400 focal lengths and field of views then each unique combination would require a separate profile.
401 Fortunately cameras typically embed state data (e.g. focal length) and camera model in image
402 metadata which can be used to retrieve the correct profile to remove radial distortion. Profile
403 creation is a relatively straight forward process from the photographer's / volunteer's point of
404 view and involves taking multiple images of a regular pattern (e.g. a chessboard). For
405 smartphones, profile creation could be embedded in the data gathering application itself and
406 for non-smartphone cameras a web application could allow images to be submitted for profile
407 creation. OpenCV (OpenCV, 2018) provides the open-source code (used in this article) to
408 undistort images.

409 This article presents a closed problem with *a priori* knowledge that only a single class would
410 occur in the image, this may not be unusual where interest is in a single species. CNNs are
411 adept at discriminating between object classes (e.g. COCO, 2018, IMAGENET, 2018) and
412 improved predictive models are frequently released as can be seen on Google's model zoo page
413 (Google, 2018). The task of generalizing to additional species using R-CNN detectors and the
414 combination of approaches outlined is eminently achievable for many species and CNNs have
415 been used in fine grained species classification (e.g. Sun et al., 2016, Tamou et al., 2018).

416 Good results were obtained with fewer than 1000 training images and this may be sufficient
417 for fine grained species classification. CNNs have performed well in classifying images
418 according to bird species with fewer than 100 examples per class (Lin et al., 2015).
419 Nonetheless, data augmentation can be employed to improve the models (Ding et al., 2016,
420 Wong et al., 2016, Perez and Wang, 2017). Augmentation transforms training images as part

421 of the training pipeline to artificially boost the number of training images. Common
422 transformations include rotation, blurring and elastic transformations, and CNN APIs usually
423 have native support for augmentation. Alternatively augmentation can be managed prior to use
424 in a preferred image processing API (e.g. Jung, 2018). It will be extremely difficult to use MV
425 to discriminate between some species without large numbers of high resolution images. For
426 example, identifying the flatfishes *Pleuronectes platessa*, *Limanda limanda* and *Platichthys*
427 *flesus* is challenging even for postgraduate marine biologists (pers. observ.).

428 It will be impossible to obtain perfect object detections and length estimations, particularly
429 in diary like volunteer applications. Pragmatically, users could be prompted to provide “hints”
430 to any application to improve detection. For example, the IGFA fish catch log smartphone
431 application (International Game Fish Association, 2018) prompts users to identify the snout
432 and tail of the fish in an image to improve detection. This process could be used to determine
433 subject rotation. Users could also be prompted to identify species where there may be
434 uncertainty and these images can contribute to the training image set. Another smartphone
435 application has used user contributed images to train a species classifier from submitted images
436 (Fishbrain, 2018). Uncertain classifications and length estimations could be clarified by the
437 general public by crowd sourcing as in other successful citizen science projects (e.g. Joly et al.,
438 2014, Silvertown et al., 2015, Zooniverse, 2017) or by using paid-for crowdsourcing services
439 (e.g. Amazon, 2017).

440 Collecting species and environmental data is a core task in marine and terrestrial ecology,
441 and images are being used to monitor disease occurrence (e.g. Boesea et al., 2008, Barbedo,
442 2017), for species identification (Nilsback and Zisserman, 2008, Branson et al., 2014, Joly et
443 al., 2014) and a host of other applications (Kühl and Burghardt, 2013). It is clear that images
444 potentially encode much valuable data which is time consuming to process manually. CNNs

445 are transforming image classification and object detection and excellent detection results can
446 now be achieved from most images.

447 **4.5 Conclusion**

448 Automatically extracting metrological data from images provides opportunities to greatly
449 increase the volume and type of data that can be collected in many data gathering scenarios
450 such as national citizen science programmes, directed surveys, remote electronic monitoring
451 (e.g. camera traps), virtual observers with camera traps and other applications. Further research
452 is needed to reduce the potential bias and increase precision in extracted data in automated
453 machine vision systems to achieve mainstream adoption, but continued advances in the
454 technology will make machine vision approaches to data processing in ecology and fisheries
455 an increasingly viable option without needing a computer science expert to develop bespoke
456 MV solutions.

457 **5 Funding and Acknowledgements**

458 Graham Monkman was supported by the Fisheries Society of the British Isles under a PhD
459 Studentship. KH was supported by CEFAS Seedcorn (DP227AE).

460 **6 Data Accessibility**

461 Code, Tensorflow configuration files, data and images with the ground truth rectangles
462 defined in the VGG Image Annotator (<http://www.robots.ox.ac.uk/~vgg/software/via>) are
463 published at <https://github.com/seabass-detection/seabass-detection>. Training and object
464 detection made use of the Tensorflow object detection API, available at
465 https://github.com/tensorflow/models/tree/master/research/object_detection.

466 **7 Author Contribution Statement**

467 GM designed the methodology, collected and analysed all data and authored all software
468 routines for the analysis (excepting 3rd party APIs as noted). FV provided guidance on the

469 methodological approaches. All authors contributed to idea conception and contributed
470 critically to the drafts and gave final approval for publication.

471 **8 References**

- 472 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X. (2015).
473 TensorFlow: Large-scale machine learning on heterogeneous systems [Web Page].
474 Retrieved 10 March 2018, from <https://www.tensorflow.org/>
- 475 Amazon, Amazon Mechanical Turk: Artificial Artificial Intelligence [online] (2017).
476 Available from: <https://www.mturk.com/mturk/welcome> [Accessed 2 Mar 2017].
- 477 Barbedo, J.G.A., A new automatic method for disease symptom segmentation in digital
478 photographs of plant leaves. *European Journal of Plant Pathology*, 147 (2), 349–364
479 (2017).
- 480 Bartholomew, D. C., Mangel, C., Alfaro-shigueto, J., Pingo, S., Jimenez, A., Godley, B. J.,
481 ... Godley, B. J. (2018). Remote electronic monitoring as a potential alternative to on-
482 board observers in small-scale fisheries. *Biological Conservation*, 219(May 2017), 35–
483 45. doi:10.1016/j.biocon.2018.01.003
- 484 Belhabib, D., Campredon, P., Lazar, N., Sumaila, U.R., Baye, B.C., Kane, E.A., and Pauly,
485 D., Best for pleasure, not for business: evaluating recreational marine fisheries in West
486 Africa using unconventional sources of data. *Palgrave Communications*, 2 (15050), 1–
487 10 (2016).
- 488 Benoît, H. P., & Allard, J. (2009). Can the data from at-sea observer surveys be used to make
489 general inferences about catch composition and discards? *Canadian Journal of Fisheries*
490 *and Aquatic Sciences*, 66(12), 2025–2039. doi:10.1139/F09-116
- 491 Bergamasco, F., Albarelli, A., Rodol, E., Torsello, A., Ambientali, S., Statistica, I., &
492 Venezia, F. (2011). RUNE-Tag : A high accuracy fiducial marker with strong occlusion
493 resilience. In *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE.
494 doi:10.1109/CVPR.2011.5995544
- 495 Bicknell, A. W. J., Godley, B. J., Sheehan, E. V., Votier, S. C., & Witt, M. J. (2016). Camera
496 technology for monitoring marine biodiversity and human impact. *Frontiers in Ecology*
497 *and the Environment*, 14(8), 424–432. doi:10.1002/fee.1322
- 498 Boesea, B.L., Clinton, P.J., Dennis, D., Golden, R.C., and Kim, B., Digital image analysis of
499 *Zostera marina* leaf injury. *Aquatic Botany*, 88 (1), 87–90 (2008).
- 500 Branson, S., Van Horn, G., Belongie, S., and Perona, P., *Bird Species Categorization Using*
501 *Pose Normalized Deep Convolutional Nets* (2014).

502 Brunelli, R. (2009). *Template matching techniques in computer vision : Theory and practice*.
503 Book, Wiley.

504 Canese, S. and Bava, S., The decline of top predators in deep coral reefs. *In: Proceedings of*
505 *the 1st Mediterranean Symposium on the Conservation of Dark Habitats*. Portorož,
506 Slovenia: UNEP, 67–68 (2015).

507 Chang, S.-K., DiNardo, G., & Lin, T.-T. (2010). Photo-based approach as an alternative
508 method for collection of albacore (*Thunnus alalunga*) length frequency from longline
509 vessels. *Fisheries Research*, 105(3), 148–155. doi:10.1016/J.FISHRES.2010.03.021

510 Claassens, L. and Hodgson, A.N., Gaining insights into in situ behaviour of an endangered
511 seahorse using action cameras. *Journal of Zoology*, 304 (2), 98–108 (2018).

512 Costello, C., Ovando, D., Hilborn, R., Gaines, S. D., Deschenes, O., & Lester, S. E. (2012).
513 Status and solutions for the world’s unassessed fisheries. *Science*, 338, 517–520.
514 doi:10.1126/science.1223389

515 Deakos, M. H. (2010). Paired-laser photogrammetry as a simple and accurate system for
516 measuring the body size of free-ranging manta rays *Manta alfredi*. *Aquatic Biology*,
517 10(1), 1–10. doi:10.3354/ab00258

518 Ding, J., Chen, B., Liu, H., and Huang, M., Convolutional neural network with data
519 augmentation for SAR target recognition. *IEEE Geoscience and remote sensing letters*,
520 13 (3), 364–368 (2016).

521 Dunbrack, R. L. (2006). In situ measurement of fish body length using perspective-based
522 remote stereo-video. *Fisheries Research*, 82(1–3), 327–331.
523 doi:10.1016/J.FISHRES.2006.08.017

524 Faunce, C. H., & Barbeaux, S. J. (2011). The frequency and quantity of Alaskan groundfish
525 catcher-vessel landings made with and without an observer. *ICES Journal of Marine*
526 *Science*, 68(8), 1757–1763. doi:10.1093/icesjms/fsr090

527 Felzenszwalb, P.F., Girshick, R.B., McAllester, D., and Ramanan, D., Object detection with
528 discriminatively trained part-based models. *IEEE transactions on pattern analysis and*
529 *machine intelligence*, 32 (9), 1627–1645 (2010).

530 Fishbrain. (2018). Fishbrain [Web Page]. Retrieved 19 July 2018, from
531 <https://fishbrain.com/mission/>

532 Friedman, J. H. (2002). Stochastic gradient boosting. *Computational Statistics & Data*
533 *Analysis*, 38(4), 367–378. doi:10.1016/S0167-9473(01)00065-2

534 Galesic, M. (2006). Dropouts on the web: Effects of interest and burden experienced during
535 an online survey. *Journal of Official Statistics*, 22(2), 313–328.

536 Garrido-Jurado, S., Muñoz-Salinas, R., Madrid-Cuevas, F. J., & Marín-Jiménez, M. J. (2014).
537 Automatic generation and detection of highly reliable fiducial markers under occlusion.
538 *Pattern Recognition*, 47(6), 2280–2292. doi:10.1016/j.patcog.2014.01.005

539 Google. (2018). Tensorflow detection model zoo [Web Page]. Retrieved 1 May 2018, from
540 [https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/dete](https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/detection_model_zoo.md)
541 [ction_model_zoo.md](https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/detection_model_zoo.md)

542 He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2017). Mask R-CNN. In *Proceedings of the*
543 *IEEE International Conference on Computer Vision* (pp. 2980–2988). Venice, Italy.
544 doi:10.1109/ICCV.2017.322

545 He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition.
546 In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp.
547 770–778). Retrieved from <http://arxiv.org/abs/1512.03385>

548 Hoerger, M. (2010). Participant dropout as a function of survey length in Internet-mediated
549 university studies: Implications for study design and voluntary participation in
550 psychological research. *Cyberpsychology, Behavior, and Social Networking*, 13(6), 697–
551 700. doi:10.1089/cyber.2009.0445

552 Hold, N., Murray, L. G., Pantin, J. R., Haig, J. A., Hinz, H., & Kaiser, M. J. (2015). Video
553 capture of crustacean fisheries data as an alternative to on-board observers. *ICES*
554 *Journal of Marine Science*, 72(6), 1811–1821. doi:10.1093/icesjms/fsv030

555 Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., ... Adam, H.
556 (2017). MobileNets: Efficient convolutional neural networks for mobile vision
557 applications. *ArXiv Preprint, 1704.04861*. Retrieved from
558 <http://arxiv.org/abs/1704.04861>

559 Hsiao, Y. H., Chen, C. C., Lin, S. I., & Lin, F. P. (2014). Real-world underwater fish
560 recognition and identification, using sparse representation. *Ecological Informatics*, 23,
561 13–21. doi:10.1016/j.ecoinf.2013.10.002

562 Hyder, K., Weltersbach, M. S., Armstrong, M., Ferter, K., Townhill, B., Ahvonen, A., ...
563 Strehlow, H. V. (2018). Recreational sea fishing in Europe in a global context –
564 participation rates, fishing effort, expenditure, and implications for monitoring and
565 assessment. *Fish and Fisheries*, 19(2), 225–243. doi:10.1111/faf.12251

566 ICES. (2012). *Report on the Classification of Stock Assessment Methods developed by*
567 *SISAM. ICES CM 2012/ACOM/SCICOM:01* (Report). Retrieved from
568 [http://www.ices.dk/community/Documents/SISAM/Report on the Classification of](http://www.ices.dk/community/Documents/SISAM/Report%20on%20the%20Classification%20of%20Stock%20Assessment%20Methods%20developed%20by%20SISAM.pdf)
569 [Stock Assessment Methods developed by SISAM.pdf](http://www.ices.dk/community/Documents/SISAM/Report on the Classification of Stock Assessment Methods developed by SISAM.pdf)

570 ICES. (2017). *Report of the Working Group on Recreational Fisheries Surveys (WGRFS)*, 6–
571 *10 June 2016. ICES CM 2016/SSGIEOM:10* (Report). Nea Peramos, Greece. Retrieved
572 from [https://www.ices.dk/sites/pub/Publication Reports/Expert Group](https://www.ices.dk/sites/pub/Publication%20Reports/Expert%20Group%20Report/SSGIEOM/2016/WGRFS/WGRFS_2016.pdf)
573 [Report/SSGIEOM/2016/WGRFS/WGRFS_2016.pdf](https://www.ices.dk/sites/pub/Publication%20Reports/Expert%20Group%20Report/SSGIEOM/2016/WGRFS/WGRFS_2016.pdf)

574 Joly, A., Goëau, H., Bonnet, P., Bakić, V., Barbe, J., Selmi, S., Yahiaoui, I., Carré, J.,
575 Mouysset, E., Molino, J.F., Boujemaa, N., and Barthélémy, D., Interactive plant
576 identification based on social image data. *Ecological Informatics*, 23, 22–34 (2014).

577 IMAGENET, IMAGENET Large Scale Visual Recognition Challenge (ILSVRC) [online]
578 (2018). Available from: <http://www.image-net.org/challenges/LSVRC/> [Accessed 6 Jun
579 2018].

580 International Game Fish Association. (2018). IGFA Catch Log [Web Page]. Retrieved 19
581 July 2018, from <http://www.igfacatchlog.org/Default.aspx>

582 Jeong, S. J., Yang, Y. S., Lee, K., Kang, J. G., & Lee, D. G. (2013). Vision-based automatic
583 system for non-contact measurement of morphometric characteristics of flatfish. *Journal*
584 *of Electrical Engineering and Technology*, 8(5), 1194–1201.
585 doi:10.5370/JEET.2013.8.5.1194

586 Jung, A., imgaug: Image augmentation for machine learning experiments (2018). Available
587 from: <https://github.com/aleju/imgaug> [Accessed 6 Jun 2018].

588 Kazemi, V., & Sullivan, J. (2014). One millisecond face alignment with an assemble of
589 regression trees. In *27th IEEE Conference on Computer Vision and Pattern Recognition*
590 (pp. 1867–1874). Columbus: IEEE Computer Society. doi:10.13140/2.1.1212.2243

591 Konovalov, D.A., Domingos, J.A., Bajema, C., White, R.D., and Jerry, D.R., Ruler Detection
592 for Automatic Scaling of Fish Images. In: *Proceedings of the International Conference*
593 *on Advances in Image Processing*. New York, NY, USA: ACM, 90–95 (2017).

594 Kraan, M., Uhlmann, S., Steenbergen, J., Helmond, A. T. M. Van, Van Helmond, A. T. M.,
595 & Van Hoof, L. (2013). The optimal process of self-sampling in fisheries: Lessons
596 learned in the Netherlands. *Journal of Fish Biology*, 83(4), 963–973.
597 doi:10.1111/jfb.12192

598 Köhl, H.S. and Burghardt, T., Animal biometrics: quantifying and detecting phenotypic
599 appearance. *Trends in Ecology & Evolution*, 28 (7), 432–441 (2013).

600 Kumar, N., Belhumeur, P. N., Biswas, A., Jacobs, D. W., Kress, W. J., Lopez, I. C., &
601 Soares, V. B. (2012). A Computer Vision System for Automatic Plant Species
602 Identification What Plant Species is this ? In *European Conference on Computer Vision*
603 (pp. 1–14). Springer-Verlag. doi:10.1007/978-3-642-33709-3_36

604 Lewin, W.-C., Arlinghaus, R., & Mehner, T. (2006). Documented and potential biological
605 impacts of recreational fishing: Insights for management and conservation. *Reviews in*
606 *Fisheries Science*, 14(4), 305–367. doi:10.1080/10641260600886455

607 Lin, T., RoyChowdhury, A., and Maji, S., Bilinear CNN Models for Fine-grained Visual
608 Recognition. In: *IEEE International Conference on Computer Vision*. Santiago: IEEE,
609 1–14 (2015).

610 Liu, F. T., Ting, K. M., & Zhou, Z.-H. (2008). Isolation Forest. In *Eighth IEEE International*
611 *Conference on Data Mining* (pp. 413–422). IEEE Computer Society.
612 doi:<http://doi.ieeecomputersociety.org/10.1109/ICDM.2008.17>

613 McClenachan, L., Historical declines of goliath grouper populations in South Florida, USA.
614 *Endangered Species Research*, 7 (3), 175–181 (2009).

615 Miranda, J. M., & Romero, M. (2017). A prototype to measure rainbow trout’s length using
616 image processing. *Aquacultural Engineering*, 76, 41–49.
617 doi:10.1016/J.AQUAENG.2017.01.003

618 National Oceanic and Atmospheric Administration. (2015). *A Cost Comparison of At-Sea*
619 *Observers and Electronic Monitoring for a Hypothetical Midwater Trawl Herring /*
620 *Mackerel Fishery*. (Report). Retrieved from
621 [https://www.greateratlantic.fisheries.noaa.gov/fish/em_cost_assessment_for_gar_herring](https://www.greateratlantic.fisheries.noaa.gov/fish/em_cost_assessment_for_gar_herring_150904_v6.pdf)
622 [_150904_v6.pdf](https://www.greateratlantic.fisheries.noaa.gov/fish/em_cost_assessment_for_gar_herring_150904_v6.pdf)

623 National Research Council. (2006). *Committee on the Review of Recreational Fisheries*
624 *Survey Methods: Review of recreational fisheries survey methods*. (Report). Washington
625 D.C.: The National Academies Press. doi:/doi.org/10.17226/11616

626 Needle, C. L., Dinsdale, R., Buch, T. B., Catarino, R. M. D., Drewery, J., & Butler, N.
627 (2015). Scottish science applications of Remote Electronic Monitoring. *ICES Journal of*
628 *Marine Science*, 72(4), 1214–1229. doi:<https://doi.org/10.1093/icesjms/fsu225>

629 Neuswanger, J. R., Wipfli, M. S., & Rosenberger, A. E. (2016). Measuring fish and their
630 physical habitats : Versatile 2-D and 3-D video techniques with user-friendly software.
631 *Canadian Journal of Fisheries and Aquatic Sciences*, 13(June), 1–48. doi:10.1139/cjfas-
632 2016-0010

633 Nilsback, M.E. and Zisserman, A., Automated flower classification over a large number of
634 classes. In: *6th Indian Conference on Computer Vision, Graphics and Image Processing,*
635 *ICVGIP 2008*. 722–729 (2008).

636 OpenCV team. (2018). OpenCV: Camera Calibration and 3D Reconstruction [Web Page].
637 Retrieved 23 April 2018, from

638 https://docs.opencv.org/master/d9/d0c/group__calib3d.html

639 Ouyang, W., Zhou, H., Li, H., Li, Q., Yan, J., and Wang, X., Jointly Learning Deep Features,
640 Deformable Parts, Occlusion and Classification for Pedestrian Detection. *IEEE*
641 *Transactions on Pattern Analysis and Machine Intelligence*, 40 (8), 1874–1887 (2018).

642 Paszke, A., Chaurasia, A., Kim, S., & Culurciello, E. (2016). ENet: A deep neural network
643 architecture for real-time semantic segmentation. *ArXiv Preprint*. Retrieved from
644 <http://arxiv.org/abs/1606.02147>

645 Pauly, D., & Morgan, G. R. (1987). Length-Based Methods in Fisheries Research. In *The*
646 *Theory and Application of Length-Based Stock Assessment Methods* (pp. 1–459).
647 Mazara del Vallo, Sicily.

648 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ...
649 Duchesnay, E. (2011). Scikit-learn: Machine learning in python. *Journal of Machine*
650 *Learning Research*, 12, 2825–2830.

651 Perez, L. and Wang, J., The Effectiveness of Data Augmentation in Image Classification
652 using Deep Learning. *CoRR*, abs/1712.0, 8 (2017).

653 Radford, Z., Hyder, K., Mugerza, E., Ferter, K., Prellezo, R., Townhill, B., ... Weltersbach,
654 M. S. (2018). The impact of marine recreational fishing on key fish stocks in European
655 waters. *PloS One*, 13(9). doi:<https://doi.org/10.1371/journal.pone.0201666>

656 Ren, S., He, K., Girshick, R., & Sun, J. (2017). Faster R-CNN: Towards real-time object
657 detection with region proposal networks. *IEEE Transactions on Pattern Analysis and*
658 *Machine Intelligence*, 39(6), 1137–1149. doi:10.1109/TPAMI.2016.2577031

659 Ricard, D., Minto, C., Jensen, O. P., & Baum, J. K. (2012). Examining the knowledge base
660 and status of commercially exploited marine species with the RAM Legacy Stock
661 Assessment Database. *Fish and Fisheries*, 13(4), 380–398. doi:10.1111/j.1467-
662 2979.2011.00435.x

663 Rizgalla, J., Shinn, A.P., Ferguson, H.W., Paladini, G., Jayasuriya, N.S., and Bron, J.E., A
664 novel use of social media to evaluate the occurrence of skin lesions affecting wild dusky
665 grouper, *Epinephelus marginatus* (Lowe, 1834), in Libyan coastal waters. *Journal of*
666 *Fish Diseases*, 40 (5), 609–620 (2017).

667 Rogers, T. D., Cambiè, G., & Kaiser, M. J. (2017). Determination of size, sex and maturity
668 stage of free swimming catsharks using laser photogrammetry. *Marine Biology*, 164(11),
669 1–11. doi:10.1007/s00227-017-3241-7

670 Rosen, S., Jörgensen, T., Hammersland-White, D., & Holst, J. C. (2013). DeepVision: a
671 stereo camera system provides highly accurate counts and lengths of fish passing inside

672 a trawl. *Canadian Journal of Fisheries and Aquatic Sciences*, 70(10), 1456–1467.
673 doi:10.1139/cjfas-2013-0124

674 Schmid, K., Reis-Filho, J.A., Harvey, E., and Giarrizzo, T., Baited remote underwater video
675 as a promising nondestructive tool to assess fish assemblages in clearwater Amazonian
676 rivers: testing the effect of bait and habitat type. *Hydrobiologia*, 784 (1), 93–109 (2017).

677 Schneider, C. A., Rasband, W. S., & Eliceiri, K. W. (2012). NIH Image to ImageJ: 25 years
678 of image analysis. *Nature Methods*, 9(7), 671–5. Retrieved from
679 <http://www.ncbi.nlm.nih.gov/pubmed/22930834>

680 Silvertown, J., Harvey, M., Greenwood, R., Dodd, M., Rosewell, J., Rebelo, T., Ansine, J.,
681 and McConway, K., Crowdsourcing the identification of organisms: A case-study of
682 iSpot. *ZooKeys*, (480), 125–146 (2015).

683 Spampinato, C., Giordano, D., Salvo, R. Di, Fisher, R. B., & Nadarajan, G. (2010).
684 Automatic Fish Classification for Underwater Species Behavior Understanding
685 Categories and Subject Descriptors. In *Proceedings of the first ACM international*
686 *workshop on Analysis and retrieval of tracked events and motion in imagery streams*
687 (pp. 45–50). Firenze, Italy. doi:10.1145/1877868.1877881

688 Stowell, D. (2018). Warblr, The Birdsong Recognition App [Web Page]. Retrieved 1 August
689 2018, from <https://www.warblr.co.uk/warblrteam/>

690 Strachan, N. J. C. (1993). Length measurement of fish by computer vision. *Computers and*
691 *Electronics in Agriculture*, 8(2), 93–104. doi:10.1016/0168-1699(93)90009-P

692 Struthers, D. P., Danylchuk, A. J., Wilson, A. D. M., & Cooke, S. J. (2015). Action cameras:
693 Bringing aquatic and fisheries research into view. *Fisheries*, 40(10), 502–512.
694 doi:10.1080/03632415.2015.1082472

695 Sun, Y., Wang, X., and Tang, X., Deep convolutional network cascade for facial point
696 detection. *Proceedings of the IEEE Computer Society Conference on Computer Vision*
697 *and Pattern Recognition*, 3476–3483 (2013).

698 Tamou, A. Ben, Benzinou, A., Nasreddine, K., and Ballihi, L., Underwater Live Fish
699 Recognition by Deep Learning. Springer, Cham, 275–283 (2018).

700 Tarrant, M. A., & Manfredo, M. J. (1993). Digit preference, recall bias, and nonresponse bias
701 in self reports of angling participation. *Leisure Sciences*, 15(3), 231–238.
702 doi:10.1080/01490409309513202

703 van Helmond, A. T. M., Chen, C., & Poos, J. J. (2017). Using electronic monitoring to record
704 catches of sole (*Solea solea*) in a bottom trawl fishery. *ICES Journal of Marine Science*,
705 74(5), 1421–1427. doi:10.1093/icesjms/fsw241

706 Vandaele, R., Aceto, J., Muller, M., Péronnet, F., Debat, V., Wang, C. W., ... Marée, R.
707 (2018). Landmark detection in 2D bioimages for geometric morphometrics: A multi-
708 resolution tree-based approach. *Scientific Reports*, 8(1), 1–13. doi:10.1038/s41598-017-
709 18993-5

710 Weinstein, B. G. (2015). MotionMeerkat: Integrating motion video detection and ecological
711 monitoring. *Methods in Ecology and Evolution*, 6(3), 357–362. doi:10.1111/2041-
712 210X.12320

713 White, D. J., Svellingen, C., & Strachan, N. J. C. (2006). Automated measurement of species
714 and length of fish by computer vision. *Fisheries Research*, 80(2–3), 203–210.
715 doi:10.1016/j.fishres.2006.04.009

716 Wong, S.C., Gatt, A., Stamatescu, V., and McDonnell, M.D., Understanding Data
717 Augmentation for Classification: When to Warp? *2016 International Conference on*
718 *Digital Image Computing: Techniques and Applications, DICTA 2016* (2016).

719 Zion, B., Alchanatis, V., Ostrovsky, V., Barki, A., & Karplus, I. (2007). Real-time
720 underwater sorting of edible fish species. *Computers and Electronics in Agriculture*,
721 56(1), 34–45. doi:10.1016/j.compag.2006.12.007

722 Zoph, B., & Le, Q. V. (2017). Neural Architecture Search with Reinforcement Learning. In
723 *International Conference on Learning Representations*. Toulon, France. Retrieved from
724 <http://arxiv.org/abs/1611.01578>

725 Zhang, J., Kan, M., Shan, S., and Chen, X., Occlusion-Free Face Alignment: Deep
726 Regression Networks Coupled with De-Corrupt AutoEncoders. *In: 2016 IEEE*
727 *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 3428–3437
728 (2016).

729 Zooniverse, Zooniverse: The list of active projects [online] (2017). Available from:
730 <https://www.zooniverse.org/projects?status=live> [Accessed 10 Feb 2017].
731

732 **9 Tables**

Table 1. Description of variables used in this article.

Variable	Derived From	Comment
<i>Physical TL</i>	N/A	The direct measurement of the physical fish with a measure.
<i>Corrected manual-TL</i>	Undistorted image	Manual estimation of the marker and fish length from the undistorted image with ImageJ. Parallax corrections applied (Appendix S1 Supporting Information, 1.4.1 & 1.4.2).
<i>MV-TL</i>	Undistorted image	Machine vision estimates of TL from undistorted images with no other corrections.
<i>Corrected MV-TL</i>	MV-TL	MV TL, corrected for parallax errors (Appendix S1 Supporting Information, 1.4.1 & 1.4.2).
<i>Rotation corrected MV-TL</i>	MV-TL	Corrected MV TL plus a geometric correction based on the height and width of the detected region (Appendix S1 Supporting Information, 1.4.3) to adjust for detections under rotation.
<i>Model corrected MV-TL</i>	MV-TL	Rotation corrected MV TL plus correction with machine learnt models generated from training data to remove outliers and correct bias in test data (Appendix S1 Supporting Information, 1.6). Only test data reported.

733

Table 2. Mean percentage intersection over union (IoU) with standard deviation (S.D.) for NASNet (Zoph & Le, 2017), ResNet-101 (He et al., 2016) and single shot MobileNet detector (Howard et al., 2017). Relative detection time (Rel. Det. Time) compares the relative detection speeds where raw detection speeds were calculated per 1000 pixels².

	Untransformed		Flipped		Rel. Det. Time
	Mean IoU	S.D.	Mean IoU	S.D.	
NASNet	93.5	2.5	93.3	2.2	1.00
ResNet	92.5	6.2	93.4	5.1	0.36
MobileNet	92.2	3.5	92.8	3.0	0.10

734

Table 3. ArUco fiducial marker (Garrido-Jurado et al., 2014) detection rates under image scaling (factor = 1.5) with width and height minimum limit of 50 pixels. Marker size is the average side length of the marker in the image. G.T. width is the ground truth horizontal length. Columns are means \pm S.D. Obj. score is the mean objectness score across all networks. ND = no detections, px = pixels. % Det. is percentage of markers detected. Scale factor is the proportion by which an image was reduced in size.

Scale factor	N	Width (px)	Height (px)	Marker size (px)	G.T. width (px)	Obj. score	% Det.
1	921	1,355	1,029	63 \pm 15	874 \pm 132	1.00 \pm 0.04	100.0
0.67	921	903	685	42 \pm 10	536 \pm 79	1.00 \pm 0.02	99.3
0.44	921	601	456	28 \pm 6	357 \pm 53	1.00 \pm 0.04	98.7
0.30	921	400	303	18 \pm 4	238 \pm 35	0.99 \pm 0.04	92.8
0.20	921	266	201	13 \pm 3	158 \pm 23	0.98 \pm 0.10	52.8
0.13	921	177	133	10 \pm 3	105 \pm 15	0.91 \pm 0.21	13.0
0.09	921	118	88	7 \pm 1	70 \pm 10	0.77 \pm 0.34	1.3
0.06	918	78	58	ND	47 \pm 7	0.55 \pm 0.39	ND
0.04	3	62	50	ND	26 \pm 0	0.005 \pm 0.007	ND

735

Table 4. Mean bias error percentage with 95% confidence intervals (CIs) for fish total length estimates made under NASNet (Zoph & Le, 2017) after corrections for lens distortion only (lens only), parallax and geometric correction (corrected) and application of machine learning to remove outliers and model errors (model corrected). The || notation is the modulus function.

	All rotations		Rotation \leq 20°	
	Mean	95% CIs	Mean	95% CIs
Lens only	-11.4	-11.6, -11.2	-9.3	-9.4, -9.1
Corrected	-4.1	-4.3, -3.9	-0.2	-2.2, -1.9
Model Corrected	-0.5	-0.6, -0.3	-0.1	-0.2, 0.1

736