



PRIFYSGOL
BANGOR
UNIVERSITY

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

Monkman, Graham G.; Hyder, Kieran; Kaiser, Michel J.; Vidal, Franck P.

Methods in Ecology and Evolution

DOI:

[10.1111/2041-210X.13282](https://doi.org/10.1111/2041-210X.13282)

Published: 01/12/2019

Peer reviewed version

[Cyswllt i'r cyhoeddiad / Link to publication](#)

Dyfyniad o'r fersiwn a gyhoeddwyd / Citation for published version (APA):

Monkman, G. G., Hyder, K., Kaiser, M. J., & Vidal, F. P. (2019). Using Machine Vision to Estimate Fish Length from Images using Regional Convolutional Neural Networks. *Methods in Ecology and Evolution*, 10(12), 2045-2056. <https://doi.org/10.1111/2041-210X.13282>

Hawliau Cyffredinol / General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

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 and Electronic Engineering, 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

29 **Summary**

- 30 1 An image can encode date-time, location and camera information as metadata and
31 implicitly encodes species information and data on human activity, e.g. the size
32 distribution of fish removals. Accurate length estimates can be made from images
33 using a fiducial marker however, their manual extraction is time consuming and
34 estimates are inaccurate without control over the imaging system. This article
35 presents a methodology which uses machine vision to estimate the total length (TL)
36 of a fusiform fish (European sea bass).
- 37 2 Three regional convolutional neural networks (R-CNN) were trained from public
38 images. Images of European sea bass were captured with a fiducial marker with 3
39 non-specialist cameras. Images were undistorted using the intrinsic lens properties
40 calculated for the camera in OpenCV, then TL was estimated using machine vision
41 (MV) to detect both marker and subject. MV performance was evaluated for the three
42 R-CNNs under downsampling and rotation of the captured images.
- 43 3 Each R-CNN accurately predicted the location of fish in test images (mean
44 intersection over union, 93%) and estimates of TL were accurate, with percent mean
45 bias error (%MBE [95% CIs]) = 2.2% [2.0, 2.4]. Detections were robust to
46 horizontal flipping and downsampling. TL estimates at absolute image rotations $> 20^\circ$
47 became increasingly inaccurate but %MBE [95% CIs] was reduced to -0.1% [-0.2,
48 0.1] using machine learning to remove outliers and model bias.
- 49 4 Machine vision can classify and derive measurements of species from images
50 without specialist equipment. It is anticipated that ecological researchers and
51 managers will make increasing use of MV where image data is collected (e.g. in
52 remote electronic monitoring, virtual observations, wildlife surveys and

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

55 **1 Introduction**

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

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

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

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

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

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

110 **2 Methods¹**

111 **2.1 Ethics**

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

116 **2.2 Training and validation image acquisition**

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

¹ Appendix S1 Supporting Information contains additional methodological detail.

123 **2.3 Fiducial marker selection and image acquisition**

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

130 **2.4 Undistorting marker images**

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

134 **2.5 Length estimation**

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

141 When estimating TL, the pixel length of the long side of the detection rectangle approximates
142 to the TL (pixels) of the fish. The real-world length per pixel, \bar{l} was estimated from the four
143 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
144 in pixels, and l is the real-world side length (e.g. 50 mm). The accuracy of \bar{l} was validated
145 manually (Linear Regression, $b = 1.003$, $R^2 = 0.999$) using ImageJ (Schneider et al., 2012).
146 Mean absolute error (MAE) and mean bias error (MBE) are reported and are calculated as

147 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
148 and \hat{y}_i is the expected (i.e. actual) TL of the i^{th} element. Hence a negative bias represents an
149 underestimate of TL.

150 **2.6 Detection and length estimation with rotation, flipping and downsampling**

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

157 **2.7 Removing outliers and modelling bias**

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

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

168 **3 Results**

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

171 regional proposals were at least partially coincident with ground truth, with a minimum IoU of
172 45% (45% IoU detection shown in Fig. 1b). Negative images had no false detections under any
173 network (score mean of 0.005 ± 0.008 , $n = 30$, $\max = 0.04$).

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

181 **3.1 Length estimates**

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

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

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

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

206 **3.2 Scale**

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

220 3.3 Rotation

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

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

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

242 4 Discussion

243 This study introduced a methodology to estimate fish TL using state-of-the-art open-source
244 R-CNNs and associated software applications (e.g. Abadi et al., 2015, OpenCV, 2018). It was

245 shown that the position of an organism in an image could be accurately predicted without strict
246 control over lighting conditions or subject background. The high degree of accuracy of the
247 predicted RoI (> 90% IoU) enabled the accurate estimation of TL. Estimation was achieved
248 without reliance on specialist cameras, multi-camera systems (e.g. Dunbrack, 2006; Rosen et
249 al., 2013) or paired lasers (e.g. Deakos, 2010, Rogers et al., 2017).

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

260 **4.1 Networks**

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

268 Neither ResNet nor NASNet are currently capable of performing real-time detections
269 however, MobileNet can be deployed on mobile devices. The performance of MobileNet in

270 this task was arguably better than ResNet and real time detection would be of particular benefit
271 in volunteer based data collection applications where users could be given immediate feedback
272 on the success or failure of a particular recognition task (Fishbrain, 2018; International Game
273 Fish Association, 2018).

274 **4.2 Length estimation**

275 Fish length measurements (TL, fork length FL and standard length SL) are particularly suited
276 to estimation by R-CNN based networks because the longitudinal dimension of an ideal
277 detection corresponds with the distal extremes of the morphological features which delineate
278 these lengths. In this manuscript, TL was used to demonstrate the methodology, but other
279 measurements (including FL and SL) may be estimated by changing the ROIs defined in the
280 training and test images or using previously determined morphometric relationships (e.g.
281 Needle et al., 2015). To date, rectangular ROIs have no history of providing length data in
282 fisheries assessments because R-CNNs are a recent development in MV. However, our results
283 demonstrate the accuracy which can be achieved where body distortion can be limited. Where
284 curvature cannot be controlled, lengths can be estimated by identifying depth midpoints and
285 calculating the line bisecting these midpoints (Strachan, 1993; White et al., 2006) or line fitting
286 to subject contours (Miranda & Romero, 2017), which requires segmentation of the subject
287 from the background. Tensorflow supports this (He et al., 2017; Google, 2018) but further work
288 would be required to validate.

289 The fiducial marker deployed was particularly easy to identify in fully automated MV
290 processing pipelines and performed well as evidenced by the low bias and high detection rates.
291 Length was more accurately estimated on afloat platforms than on the shore, because a flat
292 surface was available to measure and photograph the subject. Across both platforms and all
293 camera models there was a small but consistent overestimate of size (mean bias error, 1.6%; 6
294 mm). Possible explanations include an underestimate of lens-subject distance during camera

295 calibration which did not account for the internal distance between the lens and the glass cover
296 of the cameras, or incorrect estimation of the parameters (e.g. mean profile height) used in the
297 length correction calculation.

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

304 Despite bias being largely eliminated, outliers in TL estimates were observed (minimised
305 under NASNet). Without rotation, this error was largely attributable to errors arising from the
306 subject pose in the image. Parallax errors arising through depth differences across the fiducial
307 marker and the subject will be a major source of error which are typically dealt with by
308 excluding images following manual review (e.g. Deakos, 2010, Rogers et al., 2017). Correction
309 for tangential deflection of MV designed fiducial markers is generally supported (Garrido-
310 Jurado et al., 2014), but this is unlikely to be a consistent correction for foreground fiducial
311 markers because the tangential displacement of the marker can differ from that of the subject.

312 **4.3 Transformations**

313 Detections and length estimations were robust to flipping and downsampling. Under
314 decreasing image size the fiducial marker was found to be the limiting factor for the automatic
315 extraction of TL. This is an intrinsic limitation of using a foreground fiducial marker where
316 increasing marker size could obscure salient features. The lowest IoU was observed on the
317 smallest fish sampled, where the marker occluded a comparatively large proportion of the
318 subject (Fig. 1d). The effectiveness of the CNN under substantial downsampling indicates that

319 image sizes can be significantly reduced prior to inference to improve speed and reduce
320 memory requirements.

321 Length estimates were unbiased and acceptably precise at small degrees of rotation. The
322 bounding box under rotation predicted the x-coordinates of the snout and caudal vertices
323 reasonably well, particularly under the NASNet network (see Supporting Information S4).
324 However, the geometric model (Appendix S1 Supporting Information, 1.4.3) largely failed to
325 improve length estimates under rotation. This failure is attributable to the divergence of the
326 geometric model (detailed in Appendix S1 Additional Methods) from the bounding features of
327 the subject. The CNN detections cannot be represented by the geometry of a rotating rectangle
328 (Appendix S4 Supporting Information). Development of a more accurate geometric correction
329 model would be possible should the use case demand it.

330 Failure to generalise through all rotations poses a serious limitation in some deployment
331 scenarios. Under volunteer image collection, a significant proportion of subject rotations could
332 exceed the experimental rotation limits. A trivially implemented approach to achieve rotation
333 invariance is the brute force repetition of detection through incremental rotations. The optimal
334 detection among all rotations is then determined by some combination of metrics, e.g. height
335 to width maxima. In this article accurate detections were achieved at absolute rotations to $\sim 15^\circ$
336 which suggests that 15° steps could be used to reduce the search space. However, it may be
337 more efficient to train the network on incrementally rotated images. This training is relatively
338 trivial and is supported in most CNN APIs. Nonetheless, data on rotation invariance under
339 rotated training images was not published by Zoph and Le, (2017) and R-CNNs are not
340 intrinsically rotation invariant.

341 **4.4 Applications**

342 A foreground marker is cheap and portable, and volunteers cannot inflate size estimates by
343 moving the marker further away from the subject as possible with a background marker. The

344 methodology applies to many visual markers and to multicamera systems, and to any organism
345 for which morphological estimates are made. Difficulties will arise in unconstrained camera
346 systems where the scale indicator is difficult to distinguish in the image, (e.g. lasers in intense
347 light). None specialist markers can be segmented and length estimated using machine vision,
348 such as a standard ruler (Konovalov et al., 2017). Opportunistic fiducial markers could also be
349 segmented (e.g. human face) and used to produce estimates of fish size from historical images
350 as has been done manually to provide ecological data on some species (McClenachan, 2009;
351 Rizgalla et al., 2017).

352 Correction for lens distortion is critical for accurate photogrammetry as show in this article,
353 particularly with increased use of robust and waterproof action cameras (Struthers et al., 2015;
354 Schmid et al., 2017) which have significant radial distortion. In small scale projects or where
355 the camera model can be restricted then it may be practical for images to be undistorted on an
356 ad hoc basis. However, to deploy large scale volunteer based metrological data gathering it will
357 be necessary to build a repository of lens correction profiles for each camera model. If a camera
358 supports multiple focal lengths and field of views then each unique combination requires a
359 separate profile. Fortunately cameras typically embed state data (e.g. focal length) and camera
360 model in image metadata which can be used to retrieve the correct profile to remove radial
361 distortion. Profile creation can be embedded in an application and requires the capture of
362 multiple images of a regular pattern (e.g. a chessboard). OpenCV (OpenCV team, 2018)
363 provides the open-source code to undistort images.

364 This article presents a closed problem with *a priori* knowledge that only a single class would
365 occur in the image, this may not be unusual where interest is in a single species. CNNs are
366 adept at discriminating between object classes (e.g. IMAGENET, 2018) and improved
367 predictive models are frequently released (Google, 2018). The task of generalizing to additional
368 species using R-CNN detectors and the combination of approaches outlined is eminently

369 achievable for many species and CNNs have been used in fine grained species classification
370 (e.g. Sun et al., 2016).

371 Good results were obtained with fewer than 1000 training images and this may be sufficient
372 for fine grained species classification. CNNs have performed well in classifying images
373 according to bird species with fewer than 100 examples per class (Lin et al., 2015).
374 Nonetheless, data augmentation can be employed to improve the models (Perez & Wang,
375 2017). Augmentation transforms training images as part of the training pipeline to artificially
376 boost the number of training images. Common transformations include rotation, blurring and
377 elastic transformations, and CNN APIs usually have native support for augmentation.
378 Alternatively augmentation can be managed prior to use in a preferred image processing API
379 (e.g. Jung, 2018). It will be extremely difficult to use MV to discriminate between some species
380 without large numbers of high resolution images. For example, identifying the flatfishes
381 *Pleuronectes platessa*, *Limanda limanda* and *Platichthys flesus* is challenging even for
382 postgraduate marine biologists (pers. observ.).

383 It will be impossible to obtain perfect object detections and length estimations, particularly
384 in diary like volunteer applications. Pragmatically, users could be prompted to provide “hints”
385 to any application to improve detection. For example, the IGFA fish catch log smartphone
386 application (International Game Fish Association, 2018) prompts users to identify the snout
387 and tail of the fish in an image to improve detection. This process could be used to determine
388 subject rotation. Users could also be prompted to identify species where there may be
389 uncertainty and these images can contribute to the training image set. Another smartphone
390 application has used user contributed images to train a species classifier from submitted images
391 (Fishbrain, 2018). Uncertain classifications and length estimations could be clarified by the
392 general public by crowd sourcing as in other successful citizen science projects (e.g. Joly et al.,

393 2014, Silvertown et al., 2015, Zooniverse, 2017) or by using paid-for crowdsourcing services
394 (e.g. Amazon, 2017).

395 **4.5 Conclusion**

396 Automatically extracting metrological data from images provides opportunities to greatly
397 increase the volume and type of data that can be collected in citizen science programmes,
398 directed surveys, remote electronic monitoring, virtual observers and other applications.
399 Further research is needed to reduce the potential bias and increase precision in extracted data
400 in machine vision (MV) systems to achieve mainstream adoption, but continued technological
401 advances will make automated data processing using machine vision in ecology an increasingly
402 viable option without needing a computer science expert to develop bespoke MV solutions.

403 **5 Funding and Acknowledgements**

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

406 **6 Data Accessibility**

407 Tensorflow configuration files, data and images are available at [https://github.com/seabass-](https://github.com/seabass-detection/seabass-detection)
408 [detection/seabass-detection](https://github.com/seabass-detection/seabass-detection). The Tensorflow API is available at
409 https://github.com/tensorflow/models/tree/master/research/object_detection.

410 **7 Author Contribution Statement**

411 GM designed the methodology, collected and analysed all data and authored all software
412 routines for the analysis (excepting 3rd party APIs as noted). FV provided guidance on the
413 methodological approaches. All authors contributed to conception and critically appraised the
414 drafts and gave final approval for publication.

415 **8 References**

416 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X. (2015).
417 TensorFlow: Large-scale machine learning on heterogeneous systems [Web Page].
418 Retrieved 10 March 2018, from <https://www.tensorflow.org/>

419 Amazon. (2017). Amazon Mechanical Turk: Artificial Intelligence [Web Page]. Retrieved 2
420 March 2017, from <https://www.mturk.com/mturk/welcome>

421 Bartholomew, D. C., Mangel, C., Alfaro-shigueto, J., Pingo, S., Jimenez, A., Godley, B. J.,
422 ... Godley, B. J. (2018). Remote electronic monitoring as a potential alternative to on-
423 board observers in small-scale fisheries. *Biological Conservation*, 219(May 2017), 35–
424 45. doi:10.1016/j.biocon.2018.01.003

425 Bicknell, A. W. J., Godley, B. J., Sheehan, E. V., Votier, S. C., & Witt, M. J. (2016). Camera
426 technology for monitoring marine biodiversity and human impact. *Frontiers in Ecology
427 and the Environment*, 14(8), 424–432. doi:10.1002/fee.1322

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

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

434 Deakos, M. H. (2010). Paired-laser photogrammetry as a simple and accurate system for
435 measuring the body size of free-ranging manta rays *Manta alfredi*. *Aquatic Biology*,
436 10(1), 1–10. doi:10.3354/ab00258

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

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

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

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

446 Garrido-Jurado, S., Muñoz-Salinas, R., Madrid-Cuevas, F. J., & Marín-Jiménez, M. J. (2014).
447 Automatic generation and detection of highly reliable fiducial markers under occlusion.

448 *Pattern Recognition*, 47(6), 2280–2292. doi:10.1016/j.patcog.2014.01.005

449 Google. (2018). Tensorflow detection model zoo [Web Page]. Retrieved 1 May 2018, from
450 [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)
451 [ction_model_zoo.md](https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/detection_model_zoo.md)

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

455 He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition.
456 In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp.
457 770–778). Retrieved from <http://arxiv.org/abs/1512.03385>

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

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

465 Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., ... Adam, H.
466 (2017). MobileNets: Efficient convolutional neural networks for mobile vision
467 applications. *ArXiv Preprint, 1704.04861*. Retrieved from
468 <http://arxiv.org/abs/1704.04861>

469 Hsiao, Y. H., Chen, C. C., Lin, S. I., & Lin, F. P. (2014). Real-world underwater fish
470 recognition and identification, using sparse representation. *Ecological Informatics*, 23,
471 13–21. doi:10.1016/j.ecoinf.2013.10.002

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

476 ICES. (2012). *Report on the Classification of Stock Assessment Methods developed by*

477 SISAM. ICES CM 2012/ACOM/SCICOM:01 (Report). Retrieved from
478 <http://www.ices.dk/community/Documents/SISAM/Report on the Classification of>
479 [Stock Assessment Methods developed by SISAM.pdf](http://www.ices.dk/community/Documents/SISAM/Report on the Classification of)

480 ICES. (2017). *Report of the Working Group on Recreational Fisheries Surveys (WGRFS)*, 6–
481 *10 June 2016. ICES CM 2016/SSGIEOM:10* (Report). Nea Peramos, Greece. Retrieved
482 from <https://www.ices.dk/sites/pub/Publication Reports/Expert Group>
483 [Report/SSGIEOM/2016/WGRFS/WGRFS_2016.pdf](https://www.ices.dk/sites/pub/Publication Reports/Expert Group)

484 IMAGENET. (2018). IMAGENET Large Scale Visual Recognition Challenge (ILSVRC)
485 [Web Page]. Retrieved 6 June 2018, from <http://www.image-net.org/challenges/LSVRC/>

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

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

492 Joly, A., Goëau, H., Bonnet, P., Bakić, V., Barbe, J., Selmi, S., ... Barthélémy, D. (2014).
493 Interactive plant identification based on social image data. *Ecological Informatics*, 23,
494 22–34. doi:10.1016/j.ecoinf.2013.07.006

495 Jung, A. (2018). imgaug: Image augmentation for machine learning experiments. Computer
496 Program. Retrieved from <https://github.com/aleju/imgaug>

497 Konovalov, D. A., Domingos, J. A., Bajema, C., White, R. D., & Jerry, D. R. (2017). Ruler
498 detection for automatic scaling of fish images. In *Proceedings of the International*
499 *Conference on Advances in Image Processing* (pp. 90–95). New York, NY, USA: ACM.
500 doi:10.1145/3133264.3133271

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

504 Lin, T., RoyChowdhury, A., & Maji, S. (2015). Bilinear CNN Models for Fine-grained
505 Visual Recognition. In *IEEE International Conference on Computer Vision* (pp. 1–14).

506 Santiago: IEEE. doi:<https://doi.org/10.1109/ICCV.2015.170>

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

510 McClenachan, L. (2009). Historical declines of goliath grouper populations in South Florida,
511 USA. *Endangered Species Research*, 7(3), 175–181. doi:10.3354/esr00167

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

515 National Oceanic and Atmospheric Administration. (2015). *A Cost Comparison of At-Sea*
516 *Observers and Electronic Monitoring for a Hypothetical Midwater Trawl Herring /*
517 *Mackerel Fishery*. (Report). Retrieved from
518 [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)
519 [_150904_v6.pdf](https://www.greateratlantic.fisheries.noaa.gov/fish/em_cost_assessment_for_gar_herring_150904_v6.pdf)

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

523 Needle, C. L., Dinsdale, R., Buch, T. B., Catarino, R. M. D., Drewery, J., & Butler, N.
524 (2015). Scottish science applications of Remote Electronic Monitoring. *ICES Journal of*
525 *Marine Science*, 72(4), 1214–1229. doi:10.1093/icesjms/fsu225

526 Neuswanger, J. R., Wipfli, M. S., & Rosenberger, A. E. (2016). Measuring fish and their
527 physical habitats : Versatile 2-D and 3-D video techniques with user-friendly software.
528 *Canadian Journal of Fisheries and Aquatic Sciences*, 13(June), 1–48. doi:10.1139/cjfas-
529 2016-0010

530 OpenCV team. (2018). OpenCV: Camera Calibration and 3D Reconstruction [Web Page].
531 Retrieved 23 April 2018, from
532 https://docs.opencv.org/master/d9/d0c/group__calib3d.html

533 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ...
534 Duchesnay, E. (2011). Scikit-learn: Machine learning in python. *Journal of Machine*

535 *Learning Research*, 12, 2825–2830.

536 Perez, L., & Wang, J. (2017). The effectiveness of data augmentation in image classification
537 using deep learning. *ArXiv Preprint*, 8. Retrieved from <http://arxiv.org/abs/1712.04621>

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

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

544 Ricard, D., Minto, C., Jensen, O. P., & Baum, J. K. (2012). Examining the knowledge base
545 and status of commercially exploited marine species with the RAM Legacy Stock
546 Assessment Database. *Fish and Fisheries*, 13(4), 380–398. doi:10.1111/j.1467-
547 2979.2011.00435.x

548 Rizgalla, J., Shinn, A. P., Ferguson, H. W., Paladini, G., Jayasuriya, N. S., & Bron, J. E.
549 (2017). A novel use of social media to evaluate the occurrence of skin lesions affecting
550 wild dusky grouper, *Epinephelus marginatus* (Lowe, 1834), in Libyan coastal waters.
551 *Journal of Fish Diseases*, 40(5), 609–620. doi:10.1111/jfd.12540

552 Rogers, T. D., Cambiè, G., & Kaiser, M. J. (2017). Determination of size, sex and maturity
553 stage of free swimming catsharks using laser photogrammetry. *Marine Biology*, 164(11),
554 1–11. doi:10.1007/s00227-017-3241-7

555 Rosen, S., Jørgensen, T., Hammersland-White, D., & Holst, J. C. (2013). DeepVision: a
556 stereo camera system provides highly accurate counts and lengths of fish passing inside
557 a trawl. *Canadian Journal of Fisheries and Aquatic Sciences*, 70(10), 1456–1467.
558 doi:10.1139/cjfas-2013-0124

559 Schmid, K., Reis-Filho, J. A., Harvey, E. S., & Giarrizzo, T. (2017). Baited remote
560 underwater video as a promising nondestructive tool to assess fish assemblages in
561 clearwater Amazonian rivers: testing the effect of bait and habitat type. *Hydrobiologia*,
562 784(1), 93–109. doi:10.1007/s10750-016-2860-1

563 Schneider, C. A., Rasband, W. S., & Eliceiri, K. W. (2012). NIH Image to ImageJ: 25 years

564 of image analysis. *Nature Methods*, 9(7), 671–5. Retrieved from
565 <http://www.ncbi.nlm.nih.gov/pubmed/22930834>

566 Silvertown, J., Harvey, M., Greenwood, R., Dodd, M., Rosewell, J., Rebelo, T., ...
567 McConway, K. (2015). Crowdsourcing the identification of organisms: A case-study of
568 iSpot. *ZooKeys*, (480), 125–146. doi:10.3897/zookeys.480.8803

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

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

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

579 Sun, X., Shi, J., Dong, J., & Wang, X. (2016). Fish Recognition from Low-resolution
580 Underwater Images. In *2016 9th International Congress on Image and Signal*
581 *Processing, BioMedical Engineering and Informatics* (pp. 471–476). Datong, China.
582 doi:10.1109/CISP-BMEI.2016.7852757

583 van Helmond, A. T. M., Chen, C., & Poos, J. J. (2017). Using electronic monitoring to record
584 catches of sole (*Solea solea*) in a bottom trawl fishery. *ICES Journal of Marine Science*,
585 74(5), 1421–1427. doi:10.1093/icesjms/fsw241

586 Weinstein, B. G. (2015). MotionMeerkat: Integrating motion video detection and ecological
587 monitoring. *Methods in Ecology and Evolution*, 6(3), 357–362. doi:10.1111/2041-
588 210X.12320

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

592 Zion, B., Alchanatis, V., Ostrovsky, V., Barki, A., & Karplus, I. (2007). Real-time

593 underwater sorting of edible fish species. *Computers and Electronics in Agriculture*,
594 56(1), 34–45. doi:10.1016/j.compag.2006.12.007

595 Zooniverse. (2017). Zooniverse: The list of active projects [Web Page]. Retrieved 10
596 February 2017, from <https://www.zooniverse.org/projects?status=live>

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

600

601 **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.

602

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

603

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

604

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

605