

Automatic interpretation of cod otoliths using deep learning

Endre Moen^{1*}, Rune Vabø¹, Szymon Smoliński¹, Come Denechaud¹, Ketil Malde^{1,2},

1 Institute of Marine Research, Bergen, Norway

2 Department of Informatics, University of Bergen, Norway

* endre.moen@hi.no

Abstract

Introduction

Knowledge of fish age structure is central to the study of fish and stock dynamics. It informs on population growth and mortality and, with size distribution, is one of the main criteria used for determining the health of exploited populations and monitoring the effects of selective fishing (4; 14). Changes in the age distribution can track significant changes in population structure, such as a particularly strong year-class skewing the distribution (23), or the gradual truncation of older age classes as selective fishing mortality removes larger individuals (24). Hard structures such as scales and otoliths are used worldwide as one of the primary sources of fish age estimates, due to their ability as natural physiological and environmental recorders to form regular, temporally resolved growth increments at the daily and annual levels (2; 5; 12). While age is inferred from the “simple” counting of annual increments, the interpretation of this zonation pattern is species or even population-specific (15) and is based on precise knowledge of the timing of zone formation and of the correct identification of true and false zones (18). This process therefore requires specific expertise and is subject to uncertainties in both between-reader precision and “true” age accuracy (12). Because those estimates are central to stock assessment, ageing errors or wrong interpretation of

otolith zonation can have dramatic effects on the evaluation of fish biology and
consequently stock size and structure (3; 22; 26).
19
20

Otolith reading is time and resource consuming. Training of expert readers can take
several years depending on the species, and otoliths often undergo a long processing
phase before the final age estimates can be produced (6). This is particularly true for
demersal fish species, like Atlantic cod (*Gadus morhua*), that have large opaque otoliths
that can't be read whole and need to be prepared. These routines vary between
populations and institutes and range from direct reading of broken otoliths under a
magnifying glass, to embedding, thin sectioning and finally imaging of the sections
under a microscope. There has been a variety of methods proposed to automatically
interpret otoliths, which range from one-dimensional data analysis like intensity
transects (17) to the more recent effort toward developing machine learning (ML)
frameworks (9; 20). Despite fast progress the results remain mixed and often yield lower
precision and consistency than those obtained by trained human readers, which limits
the application of automated methods in real conditions. However, one aspect that is
often under considered by such studies are the practical time and cost benefits that
implementing a functional ML framework would provide. As noted by (11) in their
review of digital techniques for otolith analysis, “costs for human and machine ageing
systems are broadly similar since a large part of the cost is associated with preparing
the otolith sections”. As such, the net benefit of automated ageing routines is directly
dependent on the ability to scale performance using a comparatively smaller number of
samples than human readers or, alternatively, to train them on “rougher” data that can
be produced faster and at a more efficient cost.
41
42
43
44
45
46
47
48
49
50

In this study, we develop a deep learning network for estimating Atlantic cod age
using multi-exposure images of broken otoliths set in place using simple plasticine. Our
results are positive and show the potential for developing automated pipelines that
require minimum processing and could be able to produce near at-sea age estimates.
42
43
44
45

There are two families of models used, EfficientNet with CNNs B0-B7 (25) and
EfficientNetV2 with convolutional neural-networks (CNNs) small, medium, Large, and
Xtra-Large (25). The EfficientNet family of models, was introduced in 2019 and the
largest model B7 achieved state-of-the-art result on the ImageNet (8) benchmark. It
uses neural architecture search to scale image-size and the network. The EfficientNetV2
46
47
48
49
50

family of models was introduced in 2021 and Xtra-large achieved state-of-the-art result
on the ImageNet benchmark again. It extends on the previous work and introduces new
ideas, like scaling up test-set image-size. In this work we investigate EfficientNet B4-B6,
and EfficientNetV2 medium and large which shows the best compromise between
training-time and accuracy.

Method and materials

Data collection structure should be:

1. Data collection (cruises and archives) and sampling
2. photographic protocol
3. resulting images (size, exposures, number, method)
4. split into datasets and configuration

Data Collection

”1. Data collection and sampling”

We used a dataset sampled from 5150 cod otoliths which has been collected on
surveys in the period 2012-2018 conducted by Institute of Marine Research (IMR) and
aged by otolith experts. On each of the surveys, the otoliths are sampled using a
random-stratified sampling based on fish length for each trawl station, and the otoliths
from individual fish are randomly sampled.

”2. Photographic protocol”

The otolith was broken and placed on a mount, before it was captured by six images
with three light exposures and one rotation of 180°. We used the first 3 images, which
positioned the otolith so the ventral side of the otolith was near the bottom of the
camera.

”3. resulting images (size, exposures, number, method)”

The images are 3744×5616 pixels which are re-scaled for training to between
380×380 to 512×512. The image light exposure varies depending on light condition
outside, and are stored in the metadata of the JPG file. Typically the exposure order is
middle-dark-light then the rotation, and then middle-light-dark again. Sometimes the
order is changed, so the order is recovered by reading the metadata property.

Figure 1. Otolith from 2016 with read age 6 years and light exposure medium, low, and high, then rotated 180° and three new images.



The details of how the data-set is collected and sampled from surveys, camera and
mount setup, and how the otolith was processed before imaging, the resulting exposures,
naming and folders organization can be found in (10) as well as where the data-set is
available.

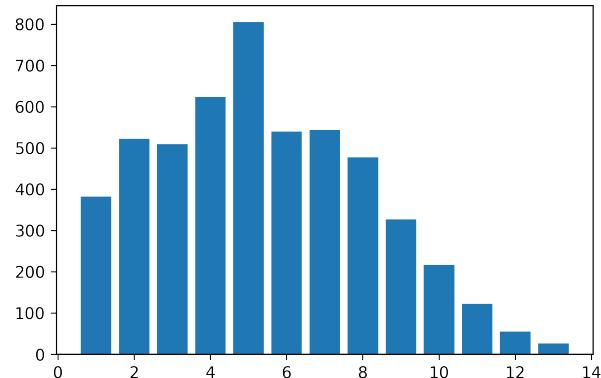


Figure 2. Age distribution of all 5150 images

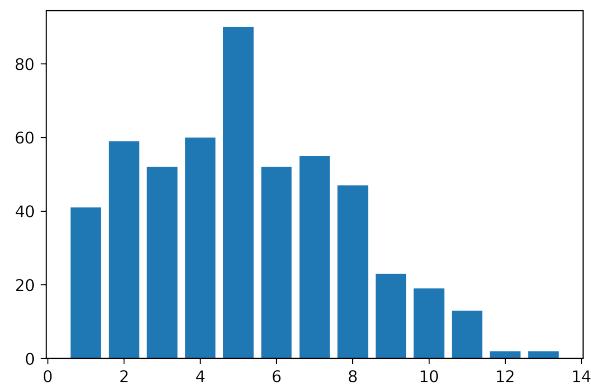


Figure 3. Age distribution of 515 images from the test set

Convolutional neural network architecture

84

Table 1. EfficientNet and EfficientNetV2 models trained with image exposure. The models are numbered for reference of ensembles later

CNN family / Image exposure	EfficientNet			EfficientNetV2	
	B4	B5	B6	Medium	Large
Minimum	1	2	3	4	5
Medium	6	7	8	9	10
Maximum	11	12	13	14	15
All (3 images)	x	x	x	16	17

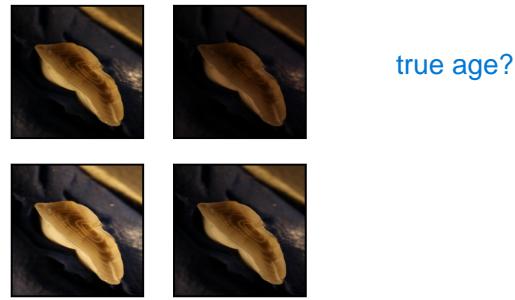
extra-large?

Each CNN was trained using transfer learning by loading ImageNet weights. The image size varies between 380×380 and 528×528 . While test-set size prediction has been done on 380×380 and 384×384 . To investigate the image-taking protocol described in (10) we have also training on 9-channel images. Three RGB-images are stacked to produce a 9-channel image. Using Timm(27) the imageNet weights were duplicated on the input layer to accommodate 9 channels. The 3 images used are of dark, medium and light exposure of the first orientation.

CNNs was selected based on performance on the ImageNet benchmark and availability of open-source implementations with imageNet weights. The imageNet benchmark is for classification while we treat aging as a regression problem (9) (R. et al.). The last layer of the CNNs has been modified to output a linear output. In the EfficientNetV2 family we have done this by applying three multi-layer perceptron layers going from 1280 output of last hidden layer to dense 256-layer, then a leakyRelu (28) layer, and then dense 32-layer, then a leakyRelu layer, and finally a linear output layer. For EfficientNet we only change the last layer from softmax output to a linear output output.

To each fold we normalize the age on the training-set by removing the mean and scaling to unit variance. The normalization is then applied to validation and test-set using sklearns StandardScalar. Test-set predictions are obtained by applying the inverse transform.

Figure 4. Otolith from 2013, read age: 6. With light exposure: medium, low, high, and expectation per channel of the three exposures.



Implementation and training

106

EfficientNetV1 B4, B5, and B6 was implemented with TensorFlow (1) and Keras (7) software packages in Python. Computation was done using CUDA 11.1 and CuDNN with Nvidia(Nvidia Corp., Santa Clara, California) A6000 accelerator card with 48 GB of GPU memory, EfficientNetV2 medium, and large was implemented with the PyTorch (19) and timm (27) software package. Computation was done on P100 cards with 12 GB of GPU memory and RTX 3090 with 24 GB of GPU memory. Pretrained weights for EfficientNet was available from Keras, and pretrained weights for EfficientNetV2 was available from timm.

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

Explaining 10 fold split for the reader:

For each model architecture the training is performed on a 10-fold split of the training dataset. This means that 10 different folds are created such that in each fold 9 parts of the dataset is used as the actual training data and one part is used as the validation data. Between different folds different part of the total training set is used for validation such that each data point participates in the validation set once. This produces 10 different trained models. The results recorded per model architecture are therefore the average of these 10 different fold models. Each fold model will give different predictions per data point and this is expected to result in an ensemble average closer to the true value than any of the single fold models. Each model architecture is also trained on three different image exposures. (By taking the ensemble across all three images exposures we can also get predictions as averages across 30 models.)
fold 1: |v|t|t|t|t|t|t|t|t|t|
fold 2: |t|v|t|t|t|t|t|t|t|t|

Augmentation was applied to the training-set. The images were augmented using rotation between 0 and 360 degrees, and reflection by the vertical axis. The pixel values has a range between 0 and 255 which was normalized to between 0 and 1.

The augmented data set can produce $360 \times 2 \times 5150 = 3.708.000$ possible images. Depending on the augmentation factor and the number of images in a training cycle, the model will likely never see the same image twice.

The cost-function is mean squared error (MSE) while the primary metric used for evaluating the models and comparing it to expert readers is accuracy. Accuracy is obtained by rounding the floating point number predictions to nearest integer and comparing the age classification against the true labels.

To get the most out of a small data-set we applied 10-fold cross-validation on 90% of the data-set, 4635 otoliths. Each fold of the 10 folds consists of 90% of the cross-validation set and 81% of the whole data-set, 4172 otoliths for training. Each fold had then 463 otoliths for validation which is 10% of the cross-validation set, and 9% of

I thought each fold did its own prediction and the average of the 10 was reported ?

does this mean that some images may appear several times in the validation set in different folds? Not True What I Wrote on Page 6?

the whole data-set. Each model is training on the 4172 otoliths and the model with the best MSE on the 463 otoliths in the validation set is chosen. The best model on the validation set was then used to predict the age on the test-set, and the metric for accuracy and MSE was recorded. The test-set is chosen at random, while the 10-fold split is chosen using stratified-kfold split which preserves a similar distribution of the whole cross-validation set in each validation set. That means the 463 images in the validation-set will have similar age distribution to that of the 4635 images in the cross-validation set. Both the test-set and the whole data-set follows a normal distribution with largest age-class being 5-year-olds.

Hyper-parameters

The CNN hyper-parameters configurations varies a little between the two families of networks, but are kept the same within the families. Some hyper-parameters that has been tuned are batch size, learning rate, k-fold size, weight decay, step size, number of epochs, early stopping, and patience. Some parameters are constrained by the GPU memory, like batch-size which is kept at 8 except for the B6 model which was run on the large A6000 GPU.

EfficientNet uses learning-rate with no scheduler while EfficientNetV2 uses Cosine Annealing scheduler (16). The training- and validation image size is as described in the papers except for large which uses smaller validation image size. The exact configuration of each network is available with each network result in the github page of the project (<https://github.com/emoen/Deep-learning-for-regression-of-cod-otoliths>).

Table 2. Hyper-parameters on each model

Param/CNN	B4	B5	B6	Medium	Large
<code>train_batch_size</code>	8	8	16	8	8
<code>img_size</code>	380	456	528	384	384
<code>val_img_size</code>	380	456	528	384	384
<code>steps_per_epoch</code>	1600	1600	1600	1600	1600
<code>epochs</code>	150	150	(150,250x2)	450	(450,250,-,450)
<code>early_stopping</code>	-	-	-	40	40
<code>early_stopping_patience</code>	14	14	(14,22,22)	-	-
<code>reduceLROnPlateau_patience</code>	7	7	(7,11,11)	-	-

Medium, all exposure was run with `steps_per_epoch=160`

Table 3. Hyper-parameters on all models

learning_rate	1e-05
n_fold	10
test_size	0.1
in_chans	3 or 9

151

in_chans is the number of channels as input for the model. It is either 3 for an
RGB image or 9 channels for 3 images.

152

153

Table 4. Hyper-parameters on TensorFlow models (B4,B5, B6)

reduceLROnPlateau_factor	0.2
which_exposure	min, medium, max

154

Table 5. Hyper-parameters on PyTorch models (medium, large)

scheduler	CosineAnnealingLR
T_max	10
min_lr	1e-06
weight_decay	1e-06
which_exposure	min, medium, max, all

155

We trained 10 models using 10-fold cross-validation which produced an ensemble prediction based on the test-set prediction on the test-set. Typically the ensemble prediction is better than any single fold prediction. Ensembles are better because they improve performance. An ensemble can make better predictions and achieve better performance than any single contributing model, just as more experts will produce higher accuracy in predicting a single otolith. Robustness; An ensemble reduces the spread or dispersion of the predictions and model performance. This result can be improved further by taking ensemble predictions of ensembles. We look at all ensembles from tuple-ensembles, consisting of 2 models, which produces an ensemble of 20 models, and triplet-ensembles consisting of 3 models, to ensemble of all models which produces an ensemble consisting of 180 models.

156

157

158

159

160

161

162

163

164

165

166

Ensemble predictions tend to make better predictions than single predictions whenever single predictions are not 100% precise. Given a true value, x , some predictions are above x and some are below x , hence the ensemble average tends to regress towards x . With several ensembles, each ensemble can be treated as a single prediction thus the same mechanism applies. The ensemble of ensembles should give even better predictions.

By choosing the best model we are over fitting to the test-set, but selecting a subset of the best of these ensembles should produce a candidate ensemble of ensemble which will produce the best prediction on a hold-out test-set.

167

168

169

Simple ensemble learning with averaging

Ensemble averaging is a simple form of committee machines (13). We investigate both simple mean and weighted mean of the 10-fold ensemble models. Simple mean gives each model equal importance and weighted mean is represented by a set of weights that sum to 1.0.

Why does ensembles work? Assume we measure a random variable (x), with a normal distribution, which is denoted as $\mathcal{N}(\mu, \sigma^2)$ with μ, σ the mean and standard deviation.

Measuring only one variable once, we know $\mathbb{E}[x_1] = \mu$ and $Var(x_1) = \sigma^2$ for any $x_1 \in (x)$

Suppose we measure the random variable (x), P times (x_1, x_2, \dots, x_p). That is, measurement in the form of $(x_1, x_2, \dots, x_p)/P$. Then the mean will still be μ . However, the variance will be smaller:

$$Var\left(\frac{x_1 + \dots + x_p}{P}\right) = \frac{Var(x_1) + \dots + Var(x_p)}{P^2} = \frac{P\sigma^2}{P^2} = \frac{\sigma^2}{P}$$

So the mean stays the same, while the variance is averaged. Hence the variance is reduced.

Results

In table 1 and table 2 are the accuracy and MSE metrics for ensemble predictions on the 10-fold training. It can be observed that in the EfficientNet family, larger networks has better MSE, while accuracy is not as correlated. A similar pattern can be observed for the EfficientNetV2 networks. However it seems like EfficientNet is better than EfficientNetV2 in both metrics unlike the results observed on the ImageNet benchmark.

Table 6. Accuracy by light exposure and CNN architectures

MSE:light/CNN	B4	B5	B6	Medium	Large
min	72.8*	74.4	73.4	74.0	-
medium	71.5	73.4	74.4	72.4	71.8*
max	70.9	-	71.5	71.1	-
9 channels	x	x	x	74.0	71.7

Table 7. MSE by light exposure and CNN architectures

ACC:light/CNN	B4	B5	B6	Medium	Large
min	.277	.277	.272	.273	-
medium	.285	.273	.262	.292	.280
max	.291	-	.305	.290	-
9 channels	x	x	x	.273	.281

192

We compare the 10-fold prediction accuracy, and MSE of all the models in a box plot in figure 5, and 9. The red line is the ensemble accuracy or MSE. The orange line is the mean accuracy or MSE. The ensemble metric is either better than or in the upper quantile for all the folds. 193
194
195
196

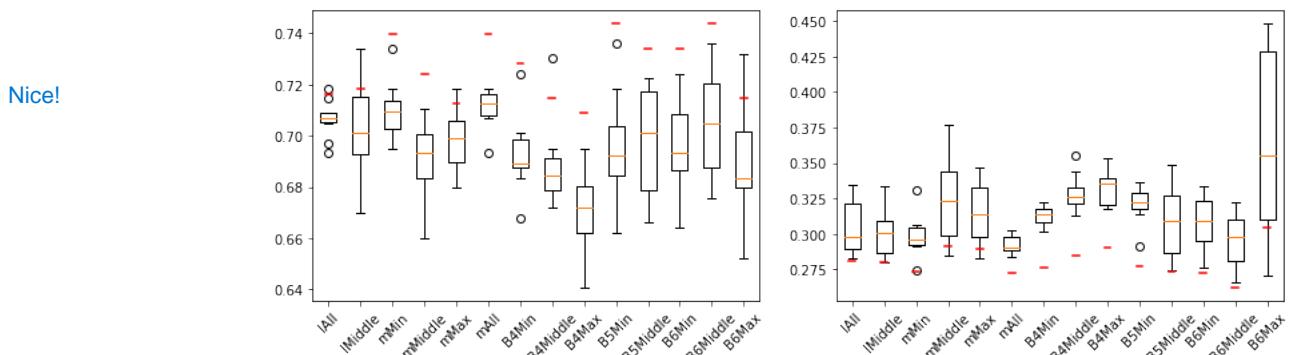


Figure 5. Accuracy score(left) and MSE of all the 17 models and the red line is simple ensemble-average prediction accuracy

By comparing the models on MSE we can see that larger models are better, e.g B6 has higher mean than B5 and B4, and large is better than medium. We also see that the EfficientNetV2 networks has higher mean than the first generation EfficientNet. However, this is not true for the ensemble predictions (red line) nor for the fold-mean or ensemble of the accuracy. We can also see that the effect of adding 3 images, creating 9 channels, on the model is that the variance is reduced, the fold mean metric increases, but the ensemble metric is reduced. 197
198
199
200
201
202
203

The box plots are produced from the folds given in table 12 and 13 in appendix C 204

Prediction by age class and residuals

Figure 6 shows the residual error as the average across all models. It looks like the residuals follow a normal-like distribution. Assuming all age groups for all models are 205
206
207

normal distributed, a table with mean and standard deviation can be found in
208
Appendix B for each model and age-group
209

The figures below shows the predictions per age group on the test-set. We can see
210
that the prediction follows a linear trend $y = x$ except for the 2-3 last years, when the
211
mean drops below $y = x$. This is even more obvious in the residual plots where the
212
prediction drops below $y = 0$ for the last 2-3 age groups. The models has a bias towards
213
lower age which is a sign of under-fitting. This correlates with the limited number of
214
otolihts in the oldest age groups.
215

Figure 8 shows scatter plots of all predictions that results in a miss-classification.
216
That is predictions that has an error greater than 0.5 in magnitude. Predictions that
217
miss by more than 1.5 in magnitude are shown with red dots.
218

Ensembles with averaging

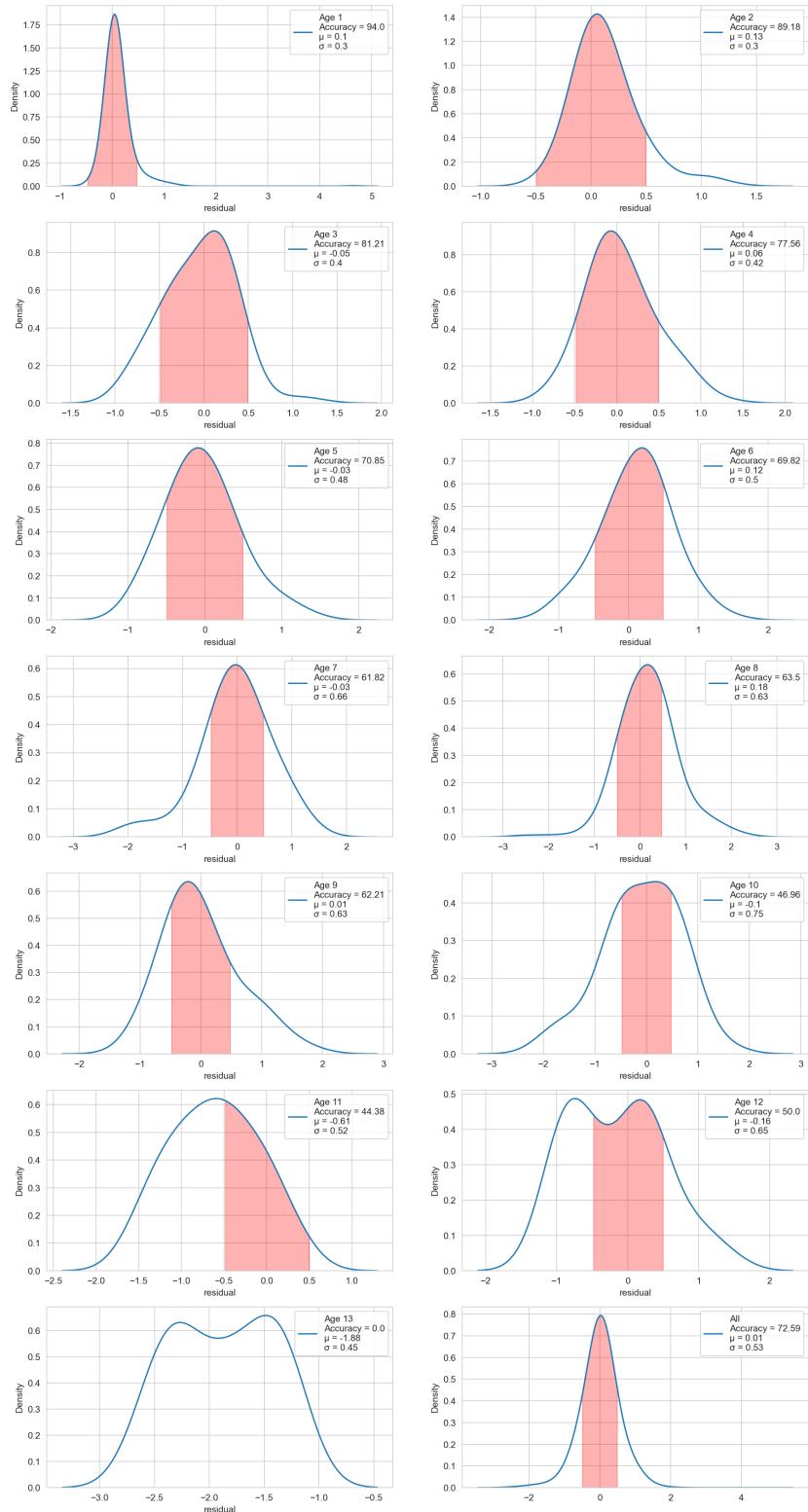
We search the space of ensembles with simple-average and weighted-average predictions
220
which are given by $\sum_{k=1}^N \binom{N}{k}$ where $N = 17$ and $k \in 1..N$ and find three ensembles
221
with simple-average which produce the best results overall with accuracy of 75.9%,
222
76.1%, and 76.9% and MSE 0.247, 0.248, and 0.248 from (1) average of all networks, (2)
223
average of B4, B5 and B6 with min exposure, and (3)average of B4, B5, B6 and
224
Medium with min exposure.
225

Table 8. Binomial combinations of simple average of ensembles accuracy

comb.	#comb	best	model
2	136	75.9	(3, 7)
3	680	77.5	(1, 3, 4)
4	2380	77.9	(1, 2, 3, 4)
5	6188	77.9	(1, 2, 3, 4, 10)
6	19448	78.6*	(1, 2, 3, 4, 7, 10)
7	24310	78.1	(1, 2, 3, 4, 6, 7, 10)
8	24310	77.5	(1, 2, 3, 4, 6, 7, 8, 11)
9	19448	77.3	(1, 2, 3, 4, 6, 7, 9, 10, 11)
10	12376	77.1	(1, 2, 3, 4, 5, 6, 7, 8, 10, 11)
11	6188	76.5	(1, 2, 3, 4, 5, 6, 7, 8, 10, 13, 14)
13	2380	76.5	(1, 2, 3, 4, 5, 7, 8, 9, 10, 12, 13, 14)
14	680	76.1	(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 12, 13, 14)
15	136	75.7	(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14)
16	17		
17	1		

Figure 6. Residuals per age class over all models. Red; correctly classified

Nice!



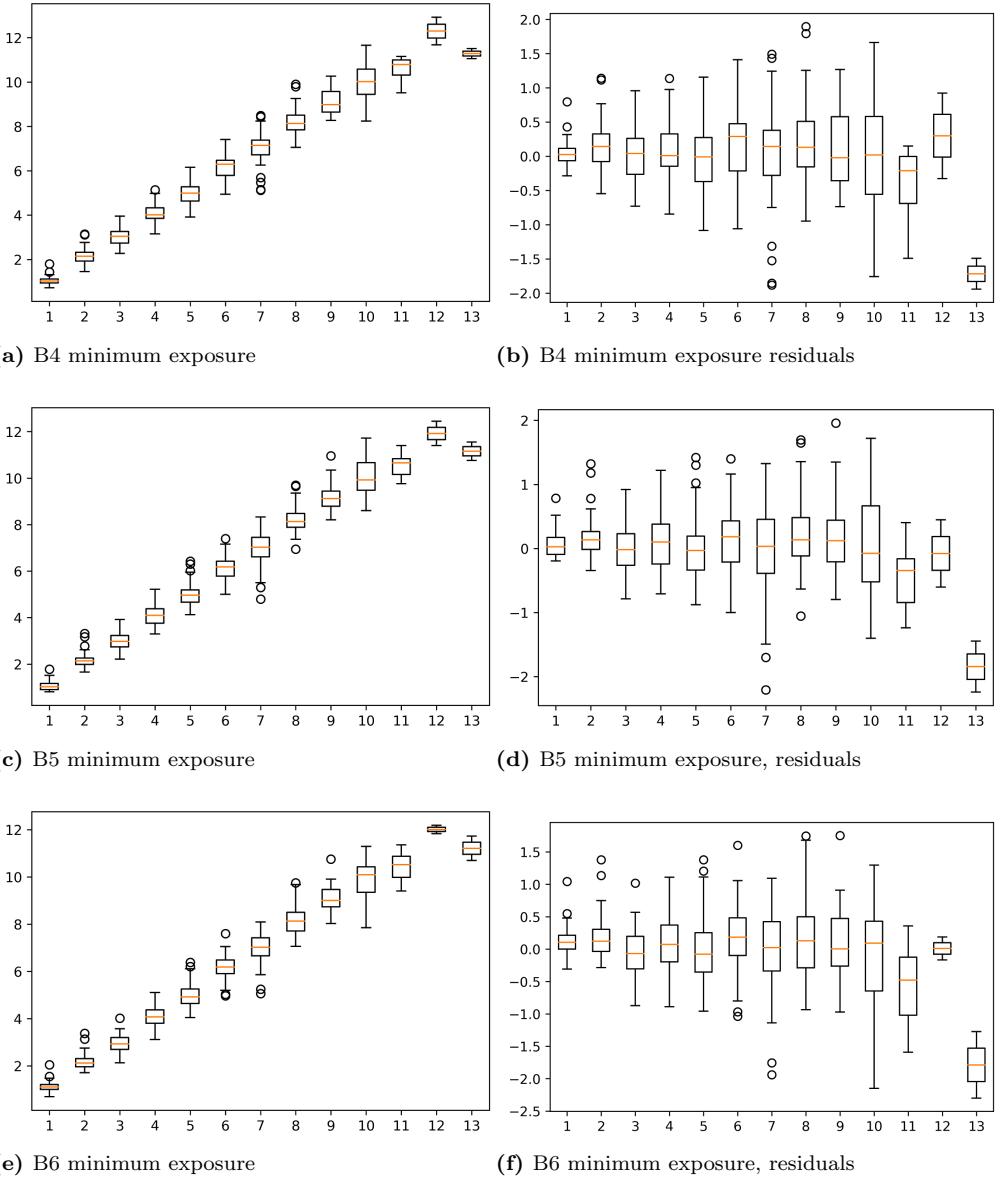
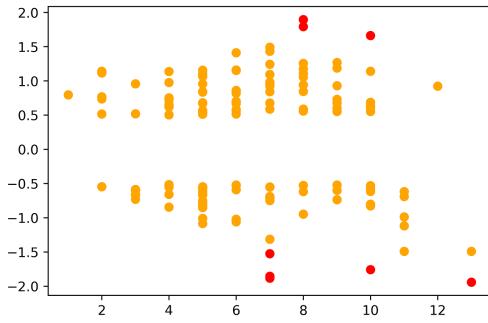
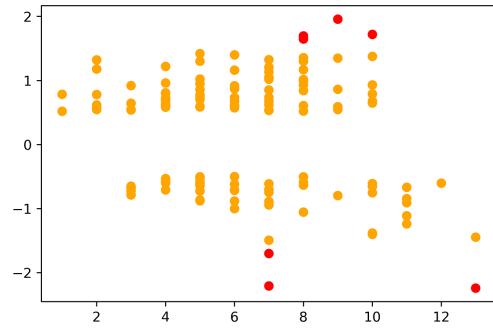


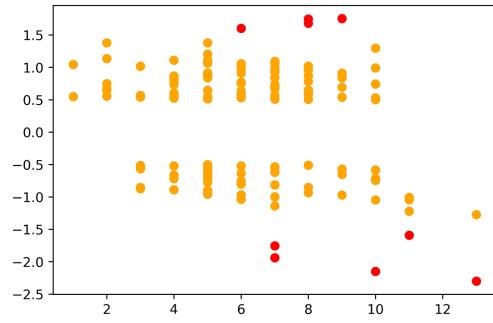
Figure 7. Comparing the models, looking at age per age class, and the residuals per prediction



(a) B4 minimum exposure



(b) B5 minimum exposure



(c) B6 minimum exposure

Figure 8. Comparing the models, looking at **age per age class**, and the reciduals per prediction

?

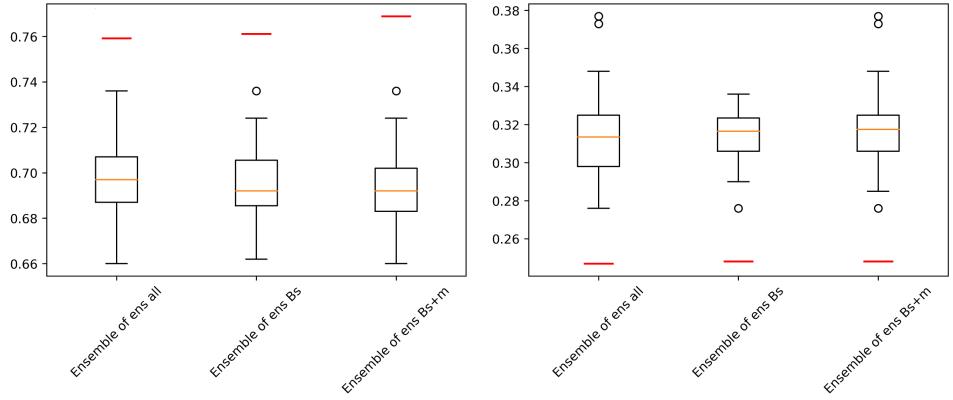


Figure 9. Ensemble of ensemble: accuracy(left) and MSE of the 3 best models

Table 9. Binomial combinations of simple average of ensembles MSE

comb.	#comb	best	model
2	136	.285	(8, 12)
3	680	.279	(8, 12, 13)
4	2380	.275	(8, 9, 12, 14)
5	6188	.273	(8, 9, 12, 13, 14)
6	19448	.271	(4, 8, 9, 12, 13, 14)
7	24310	.265	(4, 8, 9, 11, 12, 13, 14)
8	24310	.259	(4, 8, 9, 10, 11, 12, 13, 14)
9	19448	.256	(4, 5, 8, 9, 10, 11, 12, 13, 14)
10	12376	.253	(4, 5, 7, 8, 9, 10, 11, 12, 13, 14)
11	6188	.251	(4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14)
12	2380	.250	(1, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14)
13	680	.248	(1, 2, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14)
14	136	.247	(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14)
15	1		
16	1		
17	1		

227

Table 10. Accuracy/MSE pr ensemble of average. Eoe1 is ensemble of ensemble of all models, Eoe2 is for B4, B5 and B6, and Eoe3 is Eoe2 plus efficientNetV2 medium.

score/ensemble	eoe1	eoe2	eoe3
Accuracy	75.9	76.1	76.9
MSE	.247	.248	.248

228

Outliers

229

Looking at figure 7 we can see that the model under-predicts the age of older otoliths. 230
This pattern is especially observable for individuals read as 13 years. To better 231
understand the bias, figure 10 shows 6 images which has an error of more than 1 year. 232
The index of the images in the test-set is (13, 71, 270, 342, 360 and 369). Which 233
networks made the miss-prediction and by how much as well as other images that had a 234
prediction error of more than 1 year can be found in table 11 in appendix A. "More 235
here on images and outliers". 236

Figure 10. Some of the most common images with miss-predicted of more than 1.5 years



Discussion-The effect of data size:

A crucial issue in machine learning projects is to determine how much training data is needed to achieve a specific performance goal. Adapted from the number of images and classes in the ImageNet dataset a common rule of thumb for computer vision is to have a thousand images for each class. This number can be reduced if transfer learning is applied for images within the same domain. In our case of cod otolith images the domain is very different from images in ImageNet. In addition, regression towards 13 age groups is the task instead of classification into 1000 classes. The optimal number of images for our problem is difficult to estimate precisely but the computer vision rule of thumb might suggest around 13,000 images as optimal. Despite the different image domain for our problem we do see a significant performance boost in using transfer learning, suggesting that fewer images are needed than if trained from scratch. Excessive use of augmentation also reduces the number of images required. On the other hand, a general insight from deep learning is also that more training data is always better. The number of images used in this study, ~5000, might be close to the optimal but we still think that a larger training set would improve performance. During initial training we trained a B4 network on about 2000 images and obtained an accuracy of around 60%. Later another 3000 images was added and the same network was trained on ca 5000 images which resulted in accuracy of ca 70%. This suggests to us that if our training set were even larger, say 10,000 images, this would boost performance further, maybe even approaching human level accuracy of 85%.

Correlation of predictions across models

237

From the outliers we can see there is a correlation of predicting outliers across models. 238
Lets look at the correlation of models on the test-set predictions. 239
We can see that EfficientNetV2 models are most correlated to each other, and that 240
B4 models are correlated. Also B6 on middle exposure is correlated to all models. All 241
the results are highly correlations. 242

We can also look at the correlation between models pr age-class. 243

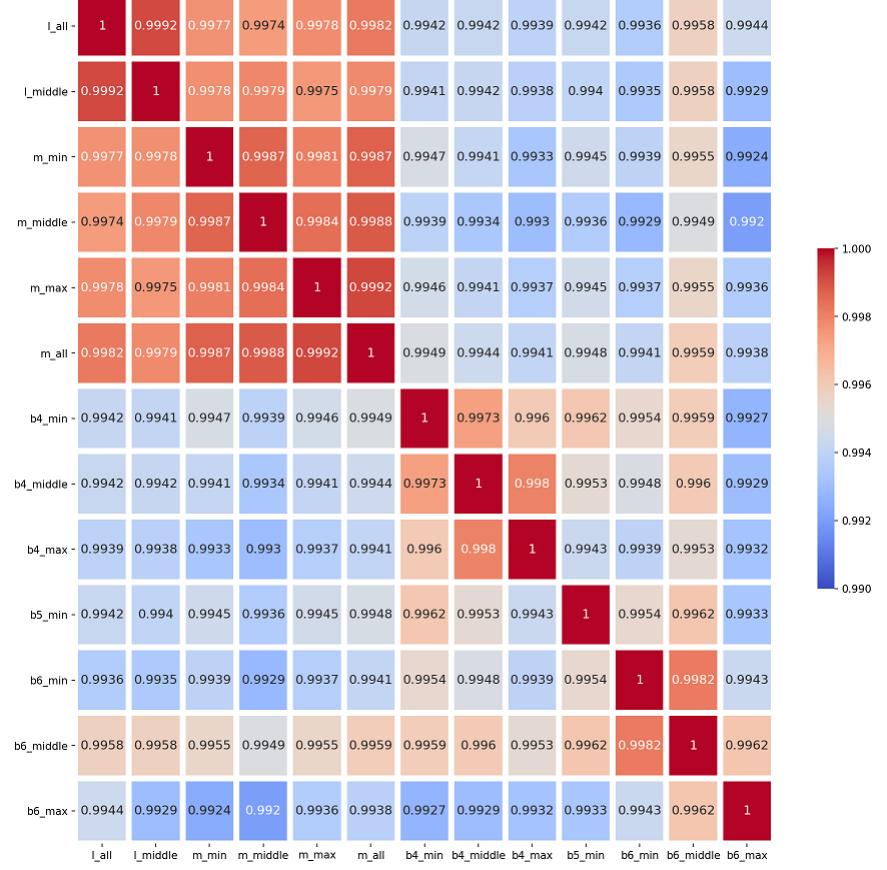
The same figure as 13 but with all the age-groups in the same scatter-plot found by 244
looking at the residuals. 245

Discussion

246

During initial training we trained a B4 network on ca 2000 images and obtained an 247
accuracy of ca 60%, later another 3000 images was added and the same network was 248

Figure 11. Pearson correlation of each model prediction on the test-set



trained on ca 5000 images which resulted in accuracy of ca 70%. It could be interesting to investigating if adding another 3-5000 images would increase the accuracy to 80%.

To reach human level accuracy a score of 85% or higher is required (?), and a score of 90% is considered good.

References

1. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., et al. (2016). Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*.
2. Albuquerque, C. Q., Lopes, L. C. S., Jaureguizar, A. J., and Condini, M. V.

Figure 12. Scatter plot of each age-class by Large-all \times Large-medium

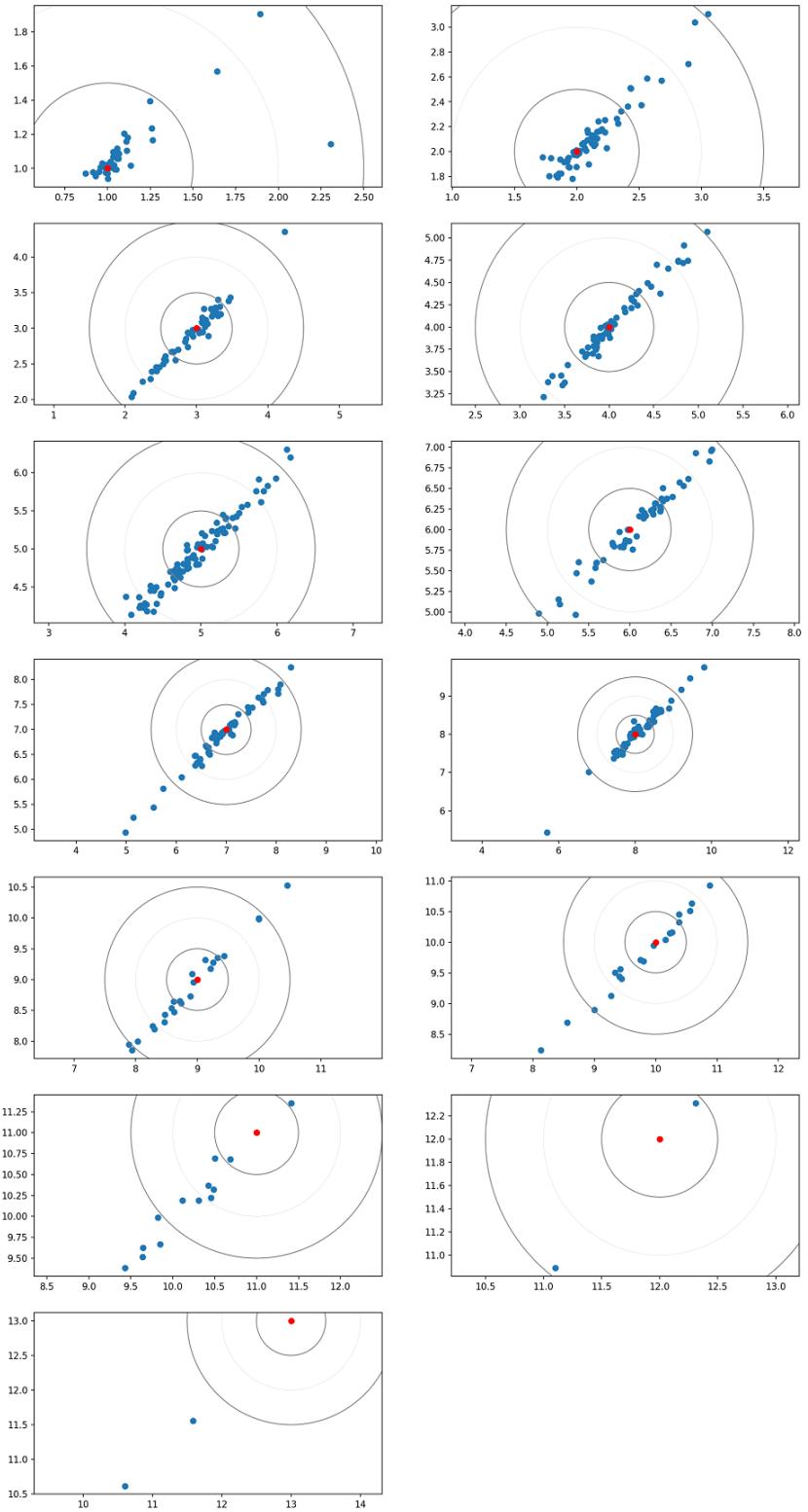


Figure 13. Scatter plot of residuals of age-class by Large-all \times Large-medium

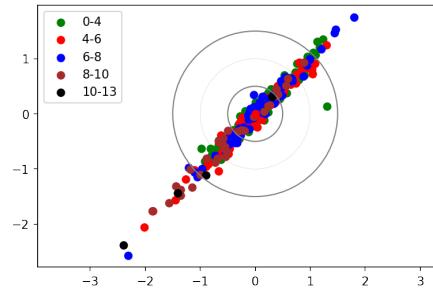


Figure 14. Sample of 25 predictions on a model of training on EfficientNetV2 size medium with minimum light exposure, left number is prediction, and right number is age read

(2019). The visual quality of annual growth increments in fish otoliths increases with latitude. *Fisheries Research*, 220: 105351.

3. Beamish, R. J. and McFarlane, G. A. (1995). A discussion of the importance of aging errors, and an application to walleye pollock: the world's largest fishery. in recent developments in fish otolith research. *University of South Carolina Press, Columbia, S.C.*, pp. 545–565.
4. Brunel, T. and Piet, G. J. (2013). Is age structure a relevant criterion for the

- health of fish stocks? *ICES Journal of Marine Science*, 70: 270–283. 266
5. Campana, S. (2001). Accuracy, precision and quality control in age determination, 267
including a review of the use and abuse of age validation methods. *Journal of fish 268
biology*, 59(2):197–242. 269
6. Carbonara, P. and Follesa, M. C. (2019). Handbook on fish age determination: a 270
mediterranean experience. *General Fisheries Commission for the Mediterranean. 271
Studies and Reviews*, 98: 1–179. 272
7. Chollet, F. and others (2018). Keras 2.1.3. <https://github.com/fchollet/keras>. 273
8. Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: 274
A large-scale hierarchical image database. In *Proceedings of IEEE Conference on 275
Computer Vision and Pattern Recognition*, pages 248–255. IEEE. 276
9. E., M., Handegard, N. O., Allken, V., Albert, O. T., Harbitz, A., and Malde, K. 277
(2018). Automatic interpretation of otoliths using deep learning. 278
10. et al., M., Thorsen, A., Godiksen, J., Malde, K., and Handegard, N. (2019). An 279
efficient protocol and data set for automated otolith image analysis. *GeoScience 280
Data Journal*. 281
11. Fisher, M. and Hunter, E. (2018). Digital imaging techniques in otolith data 282
capture, analysis and interpretation. *Marine Ecology Progress Series*, 598: 213–231. 283
12. Francis, R. C. and Campana, S. E. (2011). Inferring age from otolith 284
measurements: a review and a new approach. *Canadian Journal of Fisheries and 285
Aquatic Sciences. NRC Research Press Ottawa, Canada.* 286
<https://cdnsciencepub.com/doi/abs/10.1139/f04-063> (Accessed 3 February 2022). 287
13. HAYKIN, S. (1999). Neural networks - a comprehensive foundation. *Second 288
edition. Pearson Prentice Hall*. 289
14. Hidalgo, M., Rouyer, T., Molinero, J. C., Massutí, E., Moranta, J., Guijarro, B., , 290
and Stenseth, N. C. (2011). Synergistic effects of fishing-induced demographic 291
changes and climate variation on fish population dynamics. *Marine Ecology 292
Progress Series*, 426: 1–12. 293

15. Høie, H., Millner, R. S., McCully, S., Nedreaas, K. H., Pilling, G. M., and Skadal, J. (2009). Latitudinal differences in the timing of otolith growth: A comparison between the barents sea and southern north sea. *Fisheries Research*, *96*: 319–322.
16. Loshchilov, I. and Hutter, F. (2016). Sgdr: Stochastic gradient descent with warm restarts.
17. Mahé, K. (2009). Project no. 044132. *Automated FISH Ageing (AFISA): final activity report*.
18. Panfili, J., de Pontual, H., Troadec, H., and Wrigg, P. J. (2002). Manual of fish sclerochronology. <https://archimer.ifremer.fr/doc/00017/12801/> (Accessed 3 February 2022).
19. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
20. Politikos, D. V., Petasis, G., Chatzispyrou, A., Mytilineou, C., and Anastasopoulou, A. (2021). Automating fish age estimation combining otolith images and deep learning: The role of multitask learning. *Fisheries Research*, *242*: 106033.
- R. et al.. R., V., Moen, E., Smoliński, S., Åse Husebø, Handegard, N. O., and Malde, K. Automatic interpretation of salmon scales using deep learning. *Ecol. Inform.* *63*, 101322 (2021).
22. Ragonese, S. (2018). Methuselah or butterfly? when fish age estimates and validations tell different stories. the case of the european hake (*merluccius merluccius* l. 1758) in the mediterranean sea. ””.
23. Reglero, P. and Mosegaard, H. (2006). Onset of maturity and cohort composition

- at spawning of baltic sprat sprattus sprattus on the basis of otolith macrostructure 322
analysis. *Journal of Fish Biology*, 68: 1091–1106. 323
24. Siskey, M. R., Wilberg, M. J., Allman, R. J., Barnett, B. K., and Secor, D. H. 324
(2016). Forty years of fishing: changes in age structure and stock mixing in 325
northwestern atlantic bluefin tuna (*thunnus thynnus*) associated with size-selective 326
and long-term exploitation. *ICES Journal of Marine Science*, 73: 2518–2528. 327
25. Tan, M. and Le, Q. V. (2019). Efficientnet: Rethinking model scaling for 328
convolutional neural networks. *CoRR*, abs/1905.11946. 329
26. Tyler, A. V., Beamish, R. J., and McFarlane, G. A. (1989). Implications of age 330
determination errors to yield estimates. *ICES Journal of Marine Science*, 108: 331
27–35. 332
27. Wightman, R. (2019). Pytorch image models. 333
<https://github.com/rwightman/pytorch-image-models>. 334
28. Xu, B., Wang, N., Chen, T., and Li, M. (2015). Empirical evaluation of rectified 335
activations in convolutional network. *CoRR*, abs/1505.00853. 336

A Common outliers of more than 1.5 years, 'm' is 337

Medium, and 'l' Large network 338

B Mean and standard deviation per model x per 339

Age group 340

C Accuracy and MSE per model and per fold 341

Table 11. Outliers with more than 1.5 year error. Prediction and true age, per model

Idx	13	48	71	92	270	279	312	320	362	342	369	393	423	444	502
m, middle		4.96	10.95		9.93		5.11	10.35	8.17		5.39				
l, all		4.98			9.79	9.42		5.14	10.6	8.13		5.69			
l, middle		4.94			9.75	9.38	5.44	5.23	10.61	8.23	10.53		5.43		
B4, min	9.79	5.14		11.66	9.89		5.47	5.11	11.05	8.24					
B5, min	9.64	4.79		11.71	9.69			5.29	10.75				10.95		
B6, min	9.74	7.6	5.06			9.67		5.24	10.69	7.85	10.75			9.4	
B6, middle	9.58	5.12		11.53	9.7			5.15	10.84	8.29	10.83		9.43		
Age	8	6	7	13	10	8	11	7	7	13	10	9	8	9	11

Table 12. MSE per CNN and per fold, 'm' is Medium, and 'l' Large network

CNN/fold	1	2	3	4	5	6	7	8	9	10	ens.
B4,min	.320	.318	.306	.313	.322	.314	.315	.316	.306	.302	.277
B4,middle	.344	.328	.316	.334	.326	.320	.355	.326	.313	.325	.285
B4,max	.340	.317	.318	.347	.336	.336	.336	.320	.354	.336	.291
B5,min	.324	.322	.325	.336	.291	.314	.320	.331	.33	.317	.277
B5,middle											
B5,max											
B6,min	.325	.329	.334	.293	.312	.290	.320	.300	.276	.306	.272
B6,middle	.323	.301	.312	.268	.294	.266	.309	.311	.278	.289	.262
B6,max	.435	.306	.306	.270	.390	.321	.411	.321	.294	.448	.305
m,min	.292	.292	.294	.275	.298	.304	.304	.331	.307	.295	.273
m,middle	.321	.377	.332	.285	.285	.325	.311	.348	.295	.373	.292
m,max	.337	.297	.302	.291	.315	.347	.338	.321	.313	.283	.289
m,all	.292	.289	.289	.326	.307	.327	.283	.300	.335	.295	.281
l,min											
l,middle	.301	.281	.299	.318	.282	.305	.280	.334	.3	.310	.280
l,max											
l,all	.292	.289	.289	.326	.307	.327	.283	.30	.335	.295	.281

Table 13. Accuracy per CNN and per fold, 'm' is Medium, and 'l' Large network

CNN/fold	1	2	3	4	5	6	7	8	9	10	ens.
B4, min	69.9	68.9	68.7	68.3	68.9	70.1	69.7	66.8	68.9	72.4	72.8
B4, middle	68.5	69.3	73.0	68.5	67.8	68.2	67.2	67.2	68.3	69.5	71.5
B4, max	64.1	68.2	67.2	66.2	67.8	69.5	67.2	69.3	66.2	65.2	70.9
B5, min	71.8	69.1	69.3	66.8	73.6	70.7	66.2	68.3	69.5	68.7	74.4
B5, middle											
B5, max											
B6, min	68.3	68.5	66.4	72.4	70.7	70.9	69.3	69.3	72.0	68.9	73.4
B6, middle	68.5	69.9	67.6	73.6	72.8	72	68	69.3	72	71.1	74.4
B6, max	70.5	68.2	65.2	73.2	69.1	67.8	68.0	68.0	72.8	68.5	71.5
m, min	71.1	71.1	69.5	73.4	71.8	70.9	70.9	69.7	70.1	71.5	74.0
m, middle	68.7	67.6	68.3	71.1	70.1	70.5	69.9	68.3	69.9	66	72.4
m, max	68.9	70.1	70.3	71.3	70.7	68.5	69.7	68.0	69.1	71.8	71.3
m, all	71.7	70.7	69.3	71.3	71.8	71.8	71.3	71.7	71.1	70.7	74.0
l, min											
l, middle	69.7	73.4	69.1	67	71.8	69.9	72.6	68.2	70.5	70.3	71.8
l, max											
l, all	70.9	70.7	70.5	70.7	71.5	69.3	70.7	71.8	69.7	70.9	71.7

Figure 15. Mean of residuals per age-group

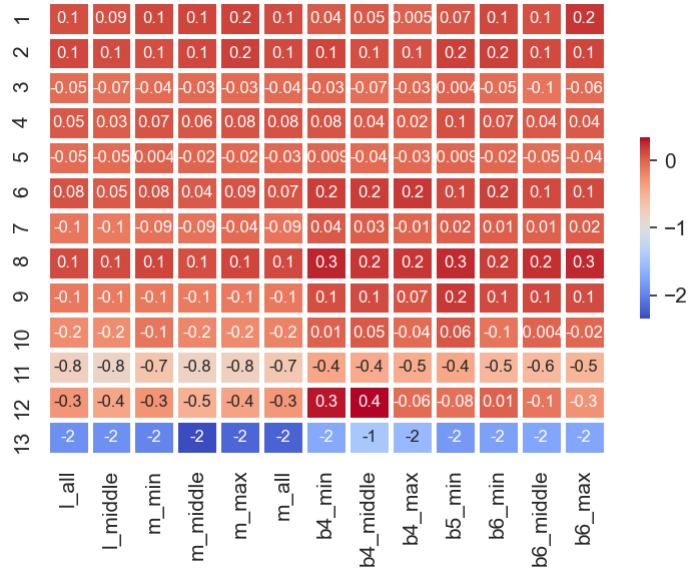


Figure 16. Standard deviation of residuals per age-group

