

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; 13). Changes in the age distribution can track significant changes in population structure, such as a particularly strong year-class skewing the distribution (22), or the gradual truncation of older age classes as selective fishing mortality removes larger individuals (23). 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 (14) and is based on precise knowledge of the timing of zone formation and of the correct identification of true and false zones (17). 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; 21; 25).
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 (16) to the more recent effort toward developing machine learning (ML)
frameworks (9; 19). 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 (24) and
EfficientNetV2 with convolutional neural-networks (CNNs) small, medium, Large, and
Xtra-Large (24). 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.

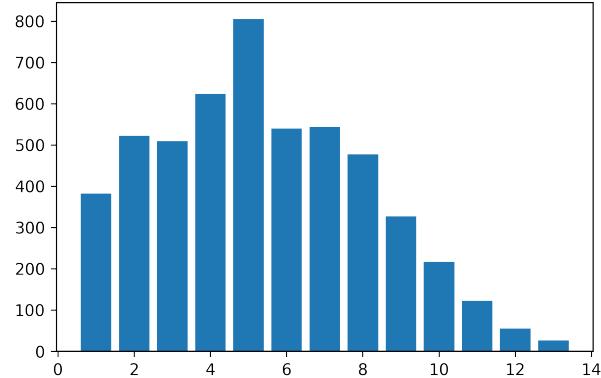


Figure 2. Age distribution of all 5150 images

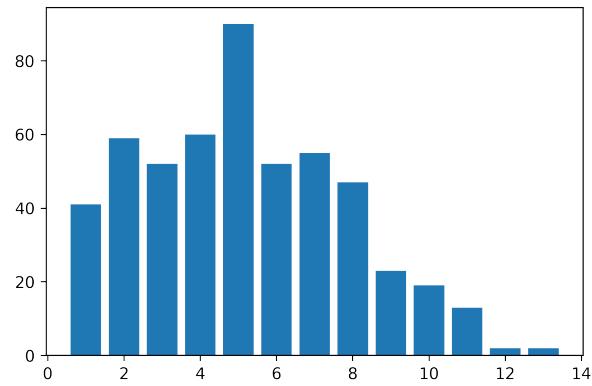


Figure 3. Age distribution of 515 images from the test set

Convolutional neural network architecture

84

Table 1. EfficientNet and EfficientNetV2 models trained with image exposure

CNN family / Image exposure	EfficientNet			EfficientNetV2	
	B4	B5	B6	medium	Large
Minimum	v	v	v	v	v
Medium	v	v	v	v	v
Maximum	v	v	v	v	v
All (3 images)	x	x	x	v	v

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(26) 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.

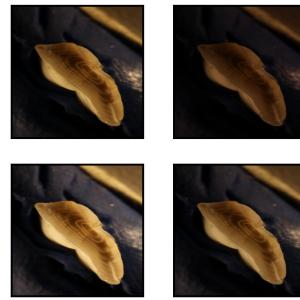
CNNs was selected based on performance on the ImageNet benchmark and availability of open-source implementations with imageNet weights. The imageNet benchmark is for classification while we treat aging as a regression problem (9) (R. et al.). The last layer of the CNNs has been modified to output a linear 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 (27) 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.

Implementation and training

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

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



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 (18) and timm (26) 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.

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 reach human level accuracy a score of 85% or higher is required (?).

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 the whole data-set. Each model is trained 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 (15). 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	160	(1600,160,1600,160)
<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)	-	-

`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.

Table 3. Hyper-parameters on all models

<code>learning_rate</code>	1e-05
<code>n_fold</code>	10
<code>test_size</code>	0.1
<code>in_chans</code>	3 or 9

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

<code>reduceLROnPlateau_factor</code>	0.2
<code>which_exposure</code>	min, medium, max

155

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

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

156

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. s 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.

157

158

159

160

161

162

163

164

165

166

167

168

169

170

Results

171

In table 1 and table 2 are the accuracy and MSE metrics for ensembled 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.

172

173

174

175

176

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	-	74.4	72.4	71.8*
max	70.9	-	71.5	71.1	-
9 channels	x	x	x	74.0	71.7

177

Table 7. MSE by light exposure and CNN architectures

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

178

We compare the 10-fold prediction accuracy, and MSE of all the models in a box plot in figure 6, and 12. 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.

179

180

181

182

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.

183

184

185

186

187

188

189

The box plots are produced from the folds given in table 11 and 12 in appendix C

190

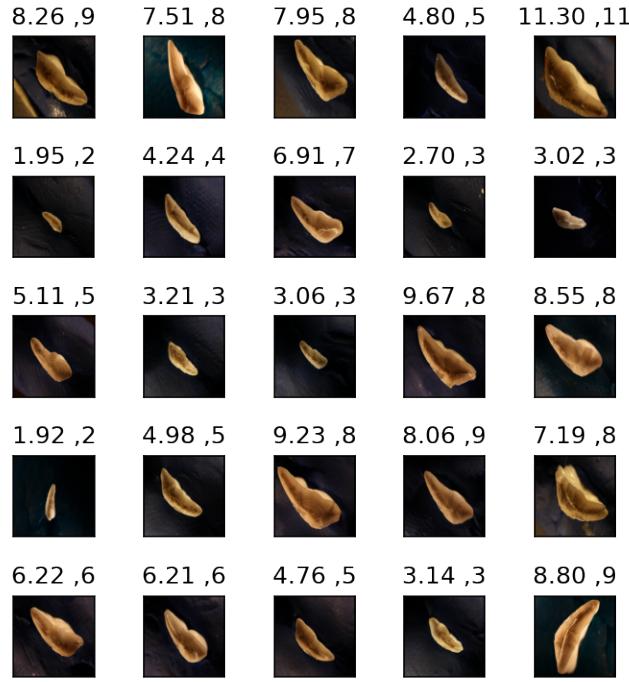


Figure 5. Sample of 25 predictions on a model of training on EfficientNetV2 size medium with minimum light exposure, left number is prediction, and right number is age read

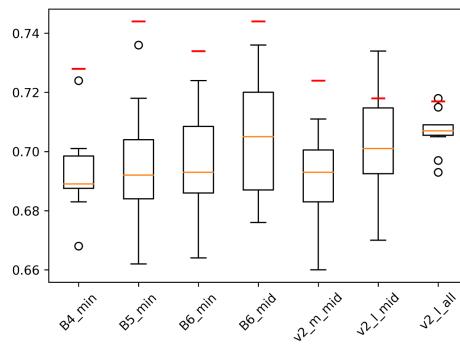


Figure 6. Accuracy score of all the 17 models and the red line is ensemble prediction accuracy

Prediction by age class and residuals

Figure 17 shows the residual error as the average across all models. It looks like a normal distribution. Assuming all age groups for all models are normally distributed, a table with mean and standard deviation can be found in Appendix B for each model and age-group

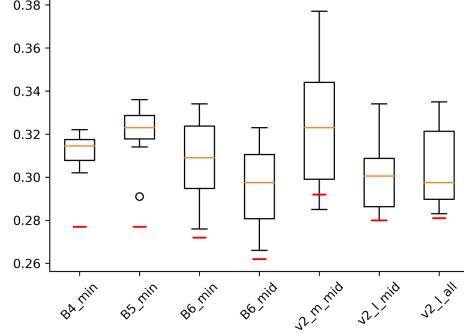


Figure 7. MSE score of all the 17 models
and the red line is ensemble prediction
MSE

The figures below shows the predictions per age group on the test-set. We can see
that the prediction follows a linear trend $y = x$ except for the 2-3 last years, when the
mean drops below $y = x$. This is even more obvious in the residual plots where the
prediction drops below $y = 0$ for the last 2-3 age groups.

Figure 10 shows scatter plots of all predictions that results in a miss-classification.
That is predictions that error greater than 0.5 in magnitude. Predictions that miss by
more than 1.5 in magnitude are shown with red dots.

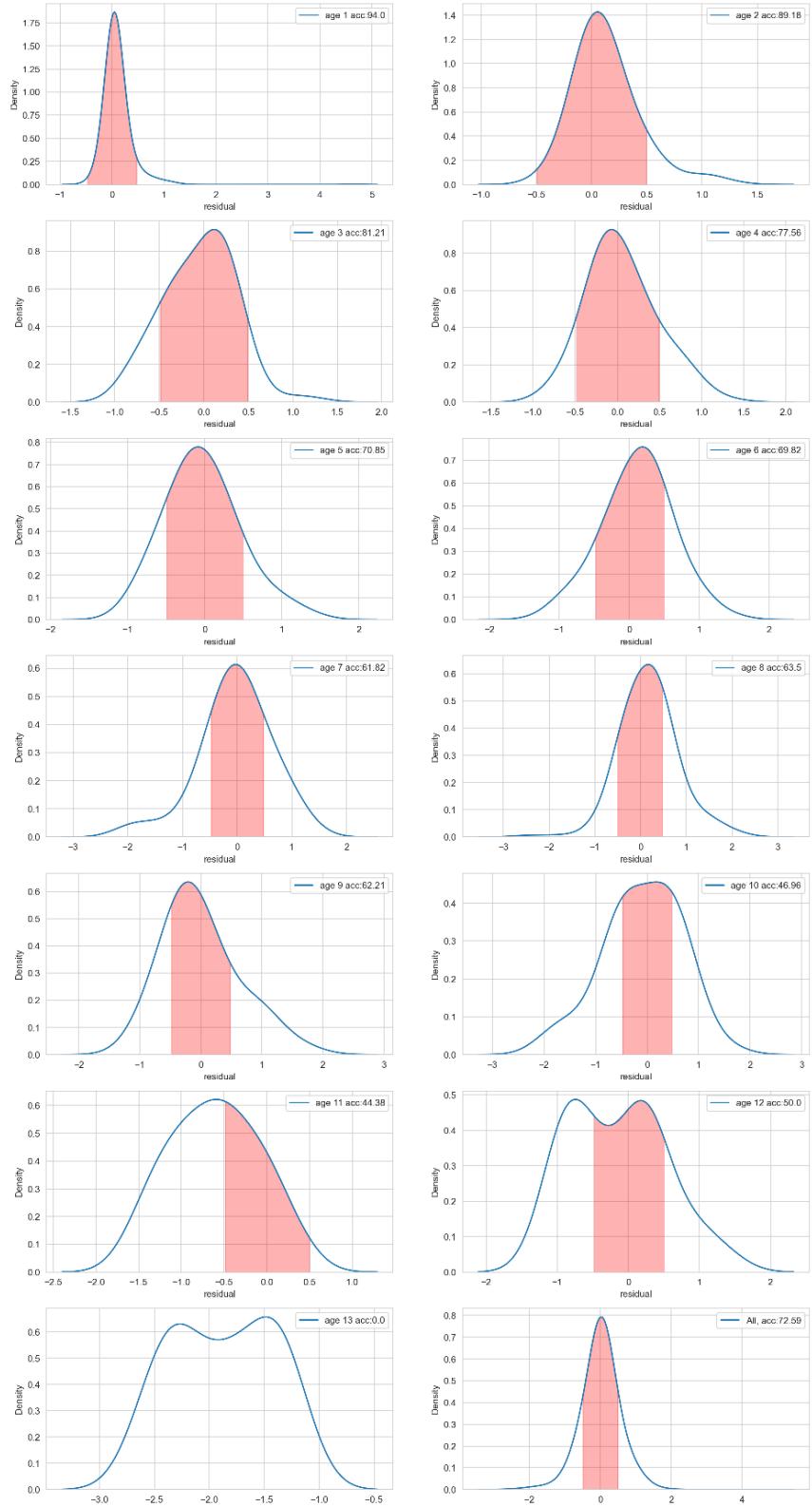
Ensemble of ensembles

We search the space of ensemble of ensemble predictions which are given by $\sum_{k=1}^N \binom{N}{k}$
where $N = 22$ and $k \in 1..N$ and find three ensemble of ensembles which produce the
best results overall with accuracy of 75.9%, 76.1%, and 76.9% and MSE 0.247, 0.248,
and 0.248 from ensemble of all networks, ensemble of B4, B5 and B6 with min exposure,
and ensemble of B4, B5, B6 and middle with min exposure.

Table 8. Accuracy/MSE pr ensemble of ensemble. 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

Figure 8. Residuals per age class over all models. Red; correctly classified



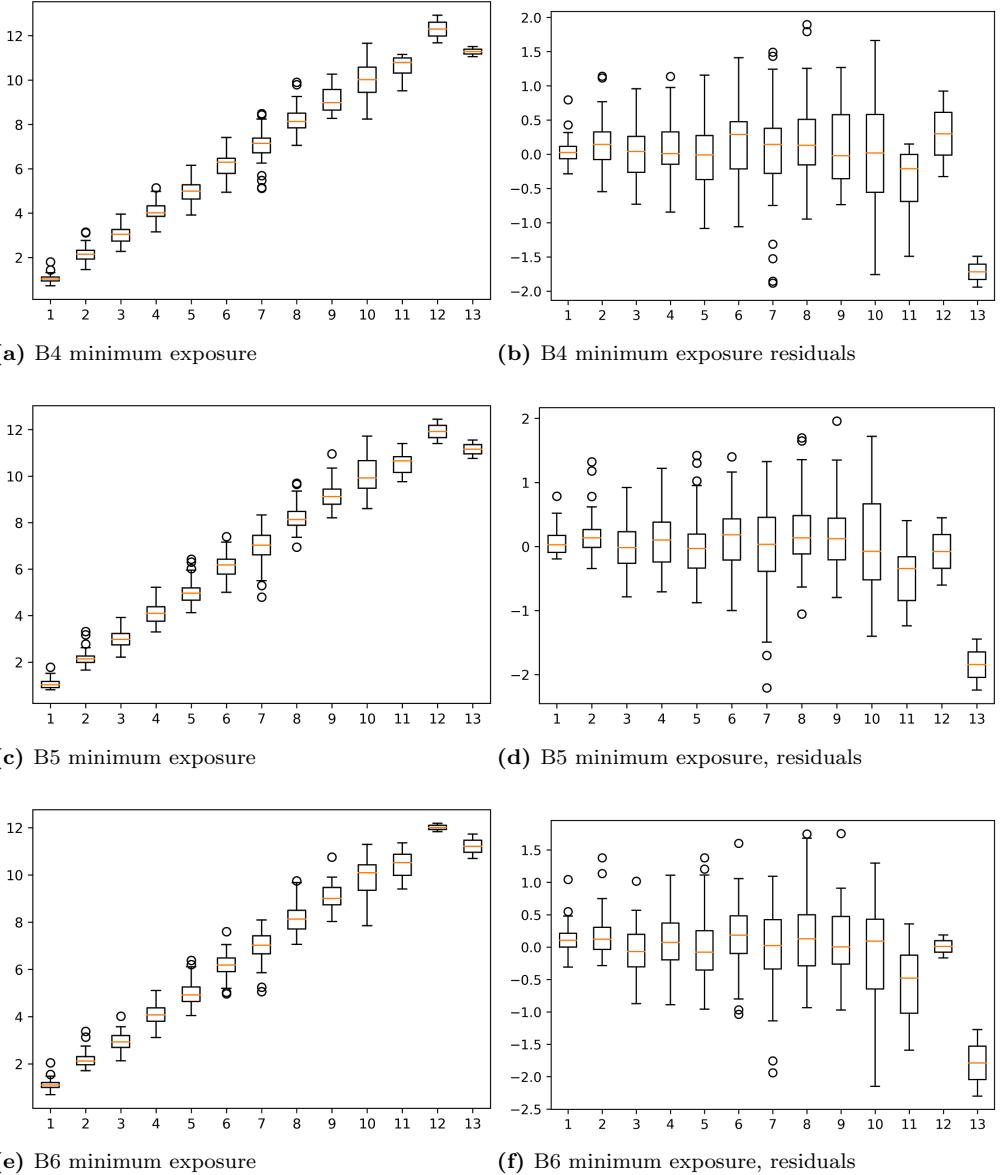
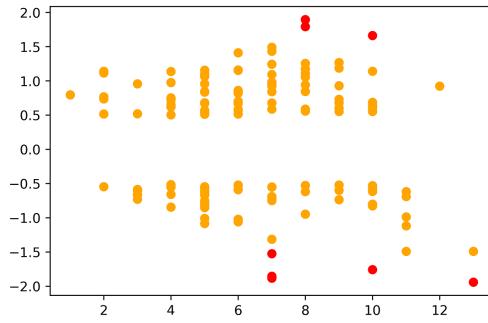
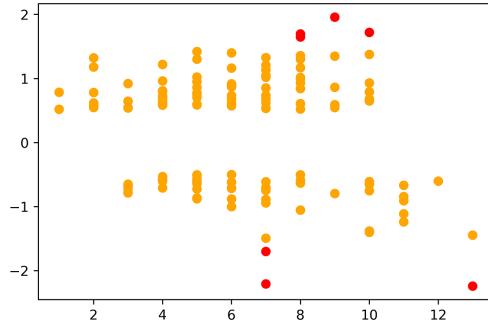


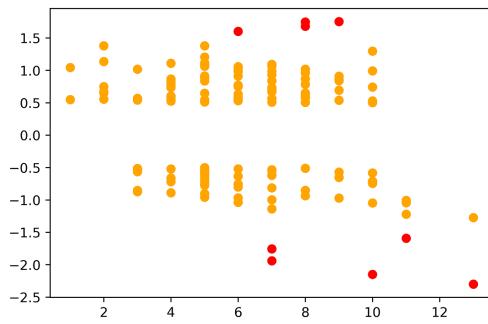
Figure 9. Comparing the models, looking at age per age class, and the residuals per prediction



(a) B4 minimum exposure



(b) B5 minimum exposure



(c) B6 minimum exposure

Figure 10. Comparing the models, looking at age per age class, and the reciduals per prediction

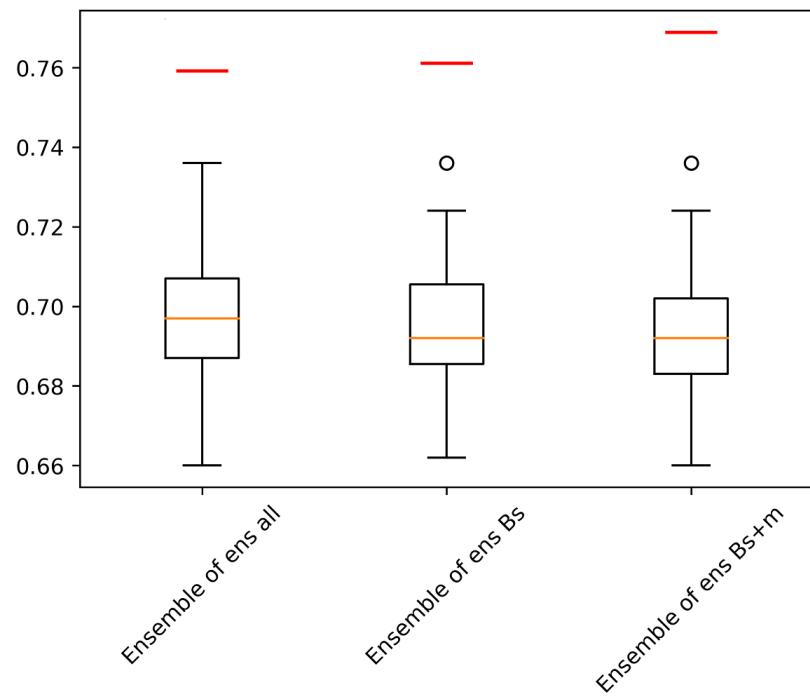


Figure 11. Ensemble of ensemble:
accuracy of the 3 best models

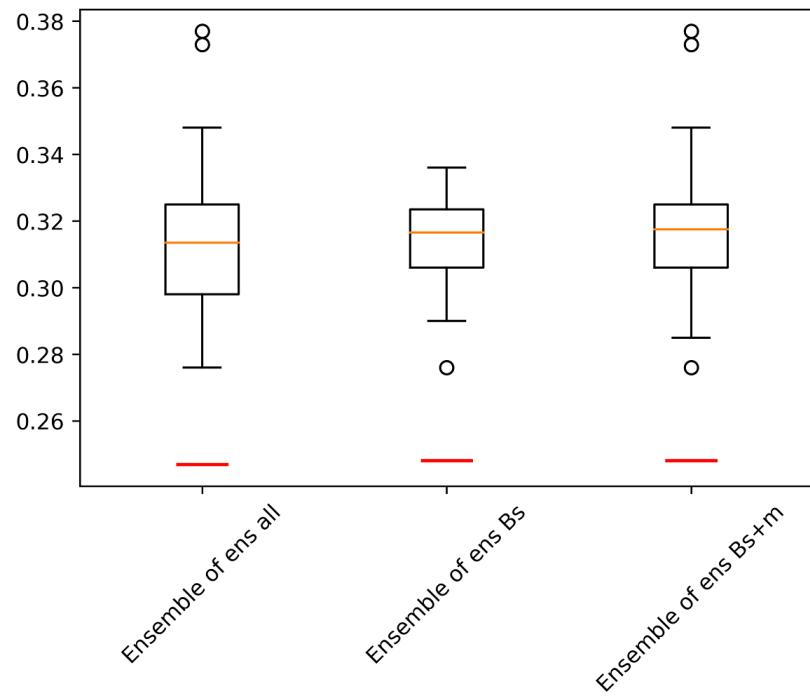


Figure 12. Ensemble of ensemble: mse of
the 3 best models

Outliers

Looking at figure 9 we can see that the model under-predicts the age of older otoliths.
This pattern is especially observable for individuals read as 15 years and older. The
oldest predication is 18 years while the test set contains individuals as old as 22 years.
To better understand the bias, figure 12 shows the 4 largest outliers from the test set
which come from two pairs

Figure 13. Some of the most common images with miss-predicted of more than 1.5 years



Figure 10 are the most commonly miss-classified images with greatest magnitude of error.

Correlation of predictions across models

From the outliers we can see there is a correlation of predicting outliers across models.
Lets look at the correlation of models on the test-set predictions.

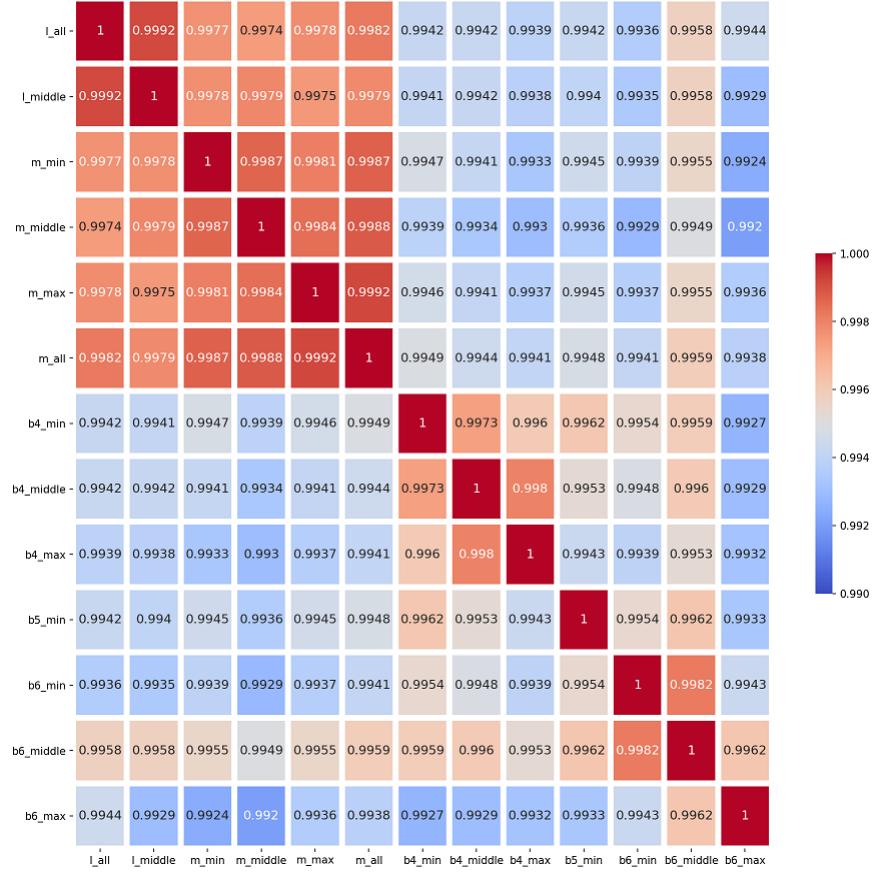
We can see that EfficientNetV2 models are most correlated to each other, and that B4 models are correlated. Also B6 on middle exposure is correlated to all models. All the results are highly correlations.

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

Discussion

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

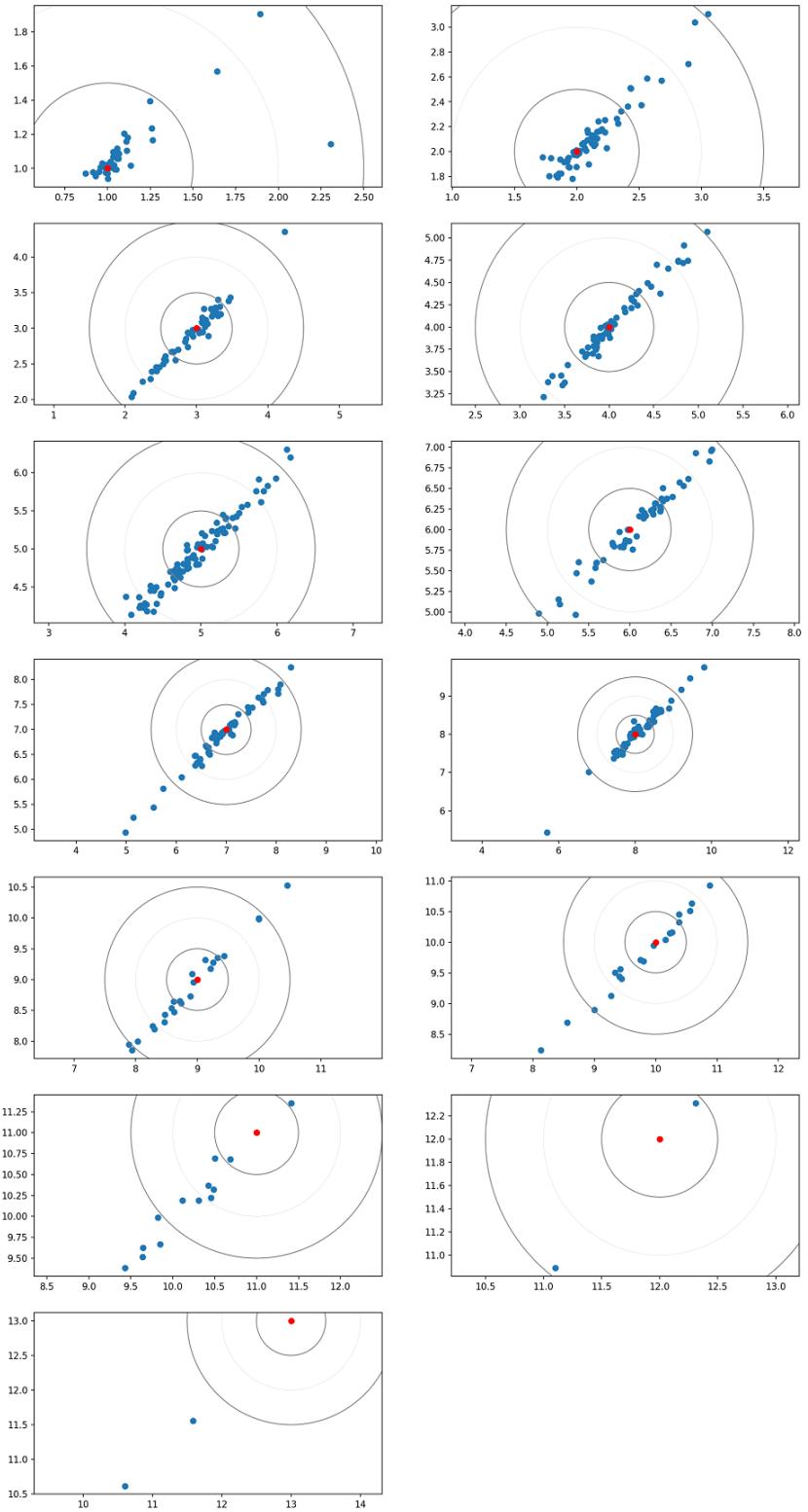
Figure 14. Pearson correlation of each model prediction on the test-set



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*. 232
2. Albuquerque, C. Q., Lopes, L. C. S., Jaureguizar, A. J., and Condini, M. V. (2019). The visual quality of annual growth increments in fish otoliths increases with latitude. *Fisheries Research*, 220: 105351. 235
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 238

Figure 15. Scatter plot of each age-class by Large-all \times Large-medium



- recent developments in fish otolith research. *University of South Carolina Press, Columbia, S.C.*, pp. 545–565. 240
241
4. Brunel, T. and Piet, G. J. (2013). Is age structure a relevant criterion for the 242
health of fish stocks? *ICES Journal of Marine Science*, 70: 270–283. 243
5. Campana, S. (2001). Accuracy, precision and quality control in age determination, 244
including a review of the use and abuse of age validation methods. *Journal of fish 245
biology*, 59(2):197–242. 246
6. Carbonara, P. and Follesa, M. C. (2019). Handbook on fish age determination: a 247
mediterranean experience. *General Fisheries Commission for the Mediterranean 248
Studies and Reviews*, 98: 1–179. 249
7. Chollet, F. and others (2018). Keras 2.1.3. <https://github.com/fchollet/keras>. 250
8. Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: 251
A large-scale hierarchical image database. In *Proceedings of IEEE Conference on 252
Computer Vision and Pattern Recognition*, pages 248–255. IEEE. 253
9. E., M., Handegard, N. O., Allken, V., Albert, O. T., Harbitz, A., and Malde, K. 254
(2018). Automatic interpretation of otoliths using deep learning. 255
10. et al., M., Thorsen, A., Godiksen, J., Malde, K., and Handegard, N. (2019). An 256
efficient protocol and data set for automated otolith image analysis. *GeoScience 257
Data Journal*. 258
11. Fisher, M. and Hunter, E. (2018). Digital imaging techniques in otolith data 259
capture, analysis and interpretation. *Marine Ecology Progress Series*, 598: 213–231. 260
12. Francis, R. C. and Campana, S. E. (2011). Inferring age from otolith 261
measurements: a review and a new approach. *Canadian Journal of Fisheries and 262
Aquatic Sciences. NRC Research Press Ottawa, Canada*. 263
<https://cdnsciencepub.com/doi/abs/10.1139/f04-063> (Accessed 3 February 2022). 264
13. Hidalgo, M., Rouyer, T., Molinero, J. C., Massutí, E., Moranta, J., Guijarro, B., , 265
and Stenseth, N. C. (2011). Synergistic effects of fishing-induced demographic 266

- changes and climate variation on fish population dynamics. *Marine Ecology Progress Series*, 426: 1–12. 267
268
14. Høie, H., Millner, R. S., McCully, S., Nedreaas, K. H., Pilling, G. M., and Skadal, 269
J. (2009). Latitudinal differences in the timing of otolith growth: A comparison 270
between the barents sea and southern north sea. *Fisheries Research*, 96: 319–322. 271
15. Loshchilov, I. and Hutter, F. (2016). Sgdr: Stochastic gradient descent with 272
warm restarts. 273
16. Mahé, K. (2009). Project no. 044132. *Automated FISH Ageing (AFISA): final 274
activity report.* 275
17. Panfili, J., de Pontual, H., Troadec, H., and Wrig, P. J. (2002). Manual of fish 276
sclerochronology. <https://archimer.ifremer.fr/doc/00017/12801/> (Accessed 3 277
February 2022). 278
18. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., 279
Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., 280
Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and 281
Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning 282
library. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., 283
and Garnett, R., editors, *Advances in Neural Information Processing Systems 32*, 284
pages 8024–8035. Curran Associates, Inc. 285
19. Politikos, D. V., Petasis, G., Chatzispyrou, A., Mytilineou, C., and 286
Anastasopoulou, A. (2021). Automating fish age estimation combining otolith 287
images and deep learning: The role of multitask learning. *Fisheries Research*, 242: 288
106033. 289
- R. et al.. R., V., Moen, E., Smoliński, S., Åse Husebø, Handegard, N. O., and Malde, 290
K. Automatic interpretation of salmon scales using deep learning. *Ecol. Inform.* 63, 291
101322 (2021). 292
21. Ragonese, S. (2018). Methuselah or butterfly? when fish age estimates and 293
validations tell different stories. the case of the european hake (*merluccius* 294
merluccius l. 1758) in the mediterranean sea. ””. 295

22. Reglero, P. and Mosegaard, H. (2006). Onset of maturity and cohort composition
296 at spawning of baltic sprat sprattus sprattus on the basis of otolith macrostructure
297 analysis. *Journal of Fish Biology*, 68: 1091–1106.
298
23. Siskey, M. R., Wilberg, M. J., Allman, R. J., Barnett, B. K., and Secor, D. H.
299 (2016). Forty years of fishing: changes in age structure and stock mixing in
300 northwestern atlantic bluefin tuna (*thunnus thynnus*) associated with size-selective
301 and long-term exploitation. *ICES Journal of Marine Science*, 73: 2518–2528.
302
24. Tan, M. and Le, Q. V. (2019). Efficientnet: Rethinking model scaling for
303 convolutional neural networks. *CoRR*, abs/1905.11946.
304
25. Tyler, A. V., Beamish, R. J., and McFarlane, G. A. (1989). Implications of age
305 determination errors to yield estimates. *ICES Journal of Marine Science*, 108:
306 27–35.
307
26. Wightman, R. (2019). Pytorch image models.
308
<https://github.com/rwightman/pytorch-image-models>.
309
27. Xu, B., Wang, N., Chen, T., and Li, M. (2015). Empirical evaluation of rectified
310 activations in convolutional network. *CoRR*, abs/1505.00853.
311

A

Common outliers of more than 1.5 years

B

Mean and standard deviation per model x per Age group

C

Accuracy and MSE per model and per fold

Table 9. Outliers with more than 1.5 year error. Index of image in test-set per model

V2-m,mid.	V2-m,mid.	V2-l,all	V2-l,mid.	B4,min	B5,min	B6,min	B6,mid.
				13	13	13	13
						48	
71	71	71	71	71	71	71	71
92	92						
				270	270		270
279	279	279	279	279	279	279	279
		312	312				
			320	320			
362	362	362	362	362	362	362	362
342	342	342	342	342	342	342	342
369	369	369	369	369		369	369
			393			393	393
423	423	423	423				
					444		
						502	502
7	7	7	9	8	7	9	9

Table 10. Outliers with more than 1.5 year error. Prediction and true age, per model

Idx	V2-m,mid.	V2-l,all	V2-l,mid.	B4,min	B5,min	B6,min	B6,mid.	Age
13				9.79	9.64	9.74	9.58	8
48						7.6		6
71	4.96	4.98	4.94	5.14	4.79	5.06	5.12	7
92	10.95							13
270				11.66	11.71		11.53	10
279	9.93	9.79	9.75	9.89	9.69	9.67	9.7	8
312		9.42	9.38					11
320			5.44	5.47				7
362	5.11	5.14	5.23	5.11	5.29	5.24	5.15	7
342	10.35	10.6	10.61	11.05	10.75	10.69	10.84	13
369	8.17	8.13	8.23	8.24		7.85	8.29	10
393			10.53			10.75	10.83	9
423	5.39	5.69	5.43					8
444					10.95			9
502						9.4	9.43	11

Figure 16. Mean of residuals per age-group

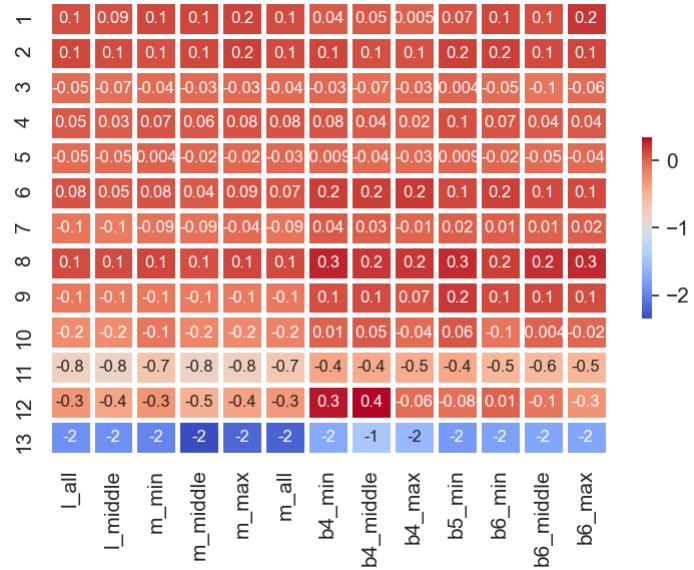


Figure 17. Standard deviation of residuals per age-group

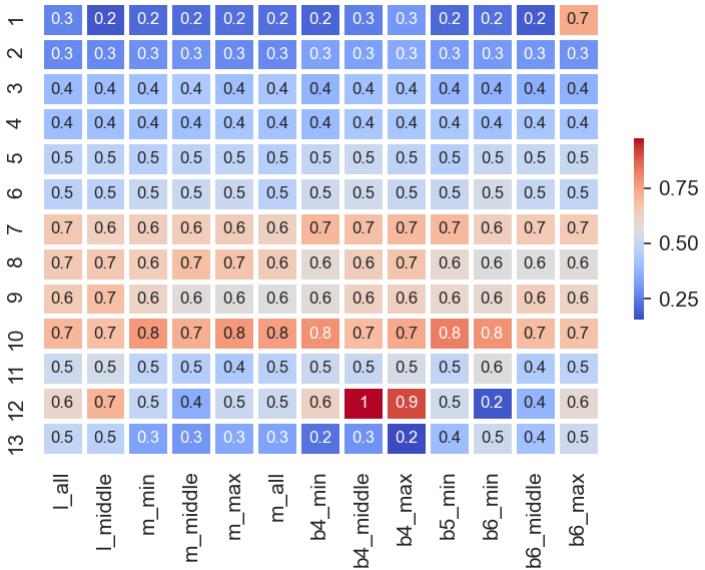


Table 11. MSE per CNN and per fold

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
medium, min	.292	.292	.294	.275	.298	.304	.304	.331	.307	.295	.273
med., mid.	.321	.377	.332	.285	.285	.325	.311	.348	.295	.373	.292
medium, max	.305	.413	.319	.327	.310	.284	.309	.315	.302	.287	.290
medium, all	.292	.289	.289	.326	.307	.327	.283	.300	.335	.295	.281
large, min											
large, middle	.301	.281	.299	.318	.282	.305	.280	.334	.3	.310	.280
large, max											
large, all	.292	.289	.289	.326	.307	.327	.283	.30	.335	.295	.281

Table 12. Accuracy per CNN and per fold

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
medium, min	71.1	71.1	69.5	73.4	71.8	70.9	70.9	69.7	70.1	71.5	74.0
med., mid.	68.7	67.6	68.3	71.1	70.1	70.5	69.9	68.3	69.9	66	72.4
medium, max	68.9	62.5	66.8	70.5	68.9	70.9	69.3	70.7	69.7	72.6	71.1
medium, all	71.7	70.7	69.3	71.3	71.8	71.8	71.3	71.7	71.1	70.7	74.0
large, min											
large, middle	69.7	73.4	69.1	67	71.8	69.9	72.6	68.2	70.5	70.3	71.8
large, max											
large, all	70.9	70.7	70.5	70.7	71.5	69.3	70.7	71.8	69.7	70.9	71.7