

# 1 Optimising a Simple Fully Convolutional Network (SFCN) for accurate 2 brain age prediction in the PAC 2019 challenge

3 **Weikang Gong<sup>1</sup>, Christian F. Beckmann<sup>1,3</sup>, Andrea Vedaldi<sup>2</sup>, Stephen M. Smith<sup>1</sup>, Han Peng<sup>1,2,3†</sup>**

4 <sup>1</sup>Wellcome Centre for Integrative Neuroimaging (WIN FMRIB), University of Oxford, Oxford, OX3  
5 9DU, United Kingdom

6 <sup>2</sup>Visual Geometry Group (VGG), University of Oxford, Oxford, OX2 6NN, United Kingdom

7 <sup>3</sup>Donders Institute for Brain, Cognition and Behaviour, Radboud University Nijmegen, Nijmegen,  
8 6525 EN, The Netherlands

9 † **Correspondence:**

10 Han Peng

11 [han.peng@ndcn.ox.ac.uk](mailto:han.peng@ndcn.ox.ac.uk)

12

13 **Keywords: predictive analysis; big data; deep learning; convolution neural network; brain age**  
14 **prediction; brain imaging (Min.5-Max. 8)**

15 **Abstract**

16 Brain age prediction from brain MRI scans not only helps improve brain ageing modelling generally,  
17 but also provides benchmarks for predictive analysis methods. Brain-age delta, which is the difference  
18 between a subject's predicted age and true age, has become a meaningful biomarker for the health of  
19 the brain. Here, we report the details of our brain age prediction models and results in the Predictive  
20 Analysis Challenge 2019. The aim of the challenge was to use T1-weighted brain MRIs to predict a  
21 subject's age in multicentre datasets. We apply a lightweight deep convolutional neural network  
22 architecture, Simple Fully Convolutional Neural Network (SFCN), and combined several techniques  
23 including data augmentation, transfer learning, model ensemble, and bias correction for brain age  
24 prediction. The model achieved first places in both of the two objectives in the PAC 2019 brain age  
25 prediction challenge: Mean absolute error (MAE) = 2.90 years without bias removal, and MAE = 2.95  
26 years with bias removal.

27 **1 Introduction**

28 Predictive analysis with data-driven machine learning algorithms brings huge promise in neuroimaging  
29 and neuroscience research. Predictive analysis can not only help disease diagnosis, such as Alzheimer's  
30 (Liu et al., 2018), Autism (Thomas et al., 2020), ADHD (Zou et al., 2017) and schizophrenia (Zeng et  
31 al., 2018), but also helps in formulating new hypotheses (Shmueli, 2010) and identifying new  
32 biomarkers (Rosenberg et al., 2018). Yet, the predictive analysis paradigm brings new challenges. First,  
33 a fair way to compare predictive analysis models is needed. In predictive analysis, it is common  
34 practice to build models in a training set, and then apply the models to a test set (Bzdok et al., 2020;  
35 Scheinost et al., 2019). It is important that no test data is used for model training or hyperparameter  
36 tuning (e.g. learning rate for gradient decent optimisations, number of layers in convnets) and to report  
37 the result objectively (LeCun et al., 2015) and avoid accidental data leakage (Lanka et al., 2019).  
38 Second, data is usually scarce for many diseases so that training a large deep learning model in such  
39 modest datasets is still hard (Raghu et al., 2019).

40 Brain ageing study is a recent example of the predictive analysis paradigm (Brown et al., 2012; Cole  
41 et al., 2018, 2017; Cole and Franke, 2017; Dosenbach et al., 2010; Franke et al., 2010; Levakov et al.,  
42 2020; Neeb et al., 2006). Studies showed that individuals' chronological age can be predicted  
43 accurately from brain MRI scans (Cole et al., 2017). Brain age delta, the difference of a subject's  
44 predicted (brain) age and chronological age, is linked with a variety of biological factors within the  
45 healthy population (Smith et al., 2020b), and group differences can be found in disease populations  
46 (Cole et al., 2019; Kaufmann et al., 2019). Yet, accurate prediction of a subject's age in healthy  
47 population is still a challenging task.

48 To tackle these challenges, a benchmarking platform is needed to objectively evaluate the models and  
49 strategies. Competitions have been seen in the field of computer vision (e.g. ImageNet (Russakovsky  
50 et al., 2015)) and proved to be a valuable vehicle for pushing AI technology (LeCun et al., 2015). In  
51 the field of neuroimaging, the Predictive Analysis Challenge (PAC) 2019 for brain age prediction<sup>1</sup>  
52 provides such opportunities for participants to train machine learning methods, and then objectively  
53 evaluate the models in a test dataset whose labels are hidden from the participants. PAC 2019 sets two  
54 objectives for brain age predictions: (1) to achieve the most accurate age prediction from brain  
55 structural MRI scans, and (2) to achieve the best accuracy while keeping the correlation between the  
56 prediction error and the ground truth age sufficiently small.

57 Our team 'BrainAgeDifference' achieved the first places in both two objectives among 79 participating  
58 teams. Our method is largely based on our previous work (Peng et al., 2019), with adaptations made  
59 for the challenge. In this report, we will provide a detailed description of our methods for PAC 2019,  
60 including the lightweight deep convnet architecture - Simple Fully Convolutional Neural Network  
61 (SFCN), and the combined techniques including data augmentation, transfer learning, model ensemble,  
62 and bias correction. We find that the lightweight model, which has achieved the state-of-the-art results  
63 in UK Biobank, works well in the multi-centre PAC 2019 dataset with a slightly adaptation in  
64 hyperparameters. SFCN pretrained on UK Biobank data achieves better single model performance than  
65 random initialized models in the PAC 2019 dataset. In addition, model ensemble with different T1-  
66 image derived maps, and different initializations, and training/validation data splits are important to  
67 achieve the best performance for the competition.

## 68 **2 Datasets and Preprocessing**

### 69 **2.1 PAC 2019**

70 The Predictive Analytic Challenge (PAC) 2019 was to predict age from brain MRI scans. The goal of  
71 the challenge includes two parts: (1) to achieve the most accurate age prediction, as measured by mean  
72 absolute error (MAE), and (2) to achieve the best MAE while keeping the Spearman correlation r-value  
73 between the prediction error (brain age delta) and the actual age below 0.1 ( $|r| < 0.1$ ). The dataset  
74 consists of both label-known training/validation dataset (2638 subjects in total) and a 'true' test set of  
75 660 subjects whose labels are unknown to the competition participants. The participants had a one-  
76 time opportunity to upload their predictions in the test set to the competition server for each objective,  
77 and the MAE and the Spearman's r-value were evaluated automatically. The subjects are from 17  
78 different sites. Most of the data is based on (Cole et al., 2017) and a few new sites were added by the  
79 organisers. The training set and the test set have the same age and site distribution.

---

<sup>1</sup> <https://web.archive.org/web/20200214101600/https://www.photon-ai.com/pac2019>

80 PAC 2019 organizers provide three version of MRI data: (a) raw T1 brain MRI scans, (b) white matter  
81 volume segmentation (WM) and (c) grey matter volume (GM) segmentation derived from T1 data. We  
82 use all three versions to develop deep learning models. We further preprocess the raw T1 images using  
83 FSL (Smith et al., 2004) (command `fsl_anat`) to derive two different pseudo-modalities: one is brain  
84 linearly registered to standard 1mm MNI space (by FLIRT), and the other is brain non-linearly  
85 registered to standard 1mm MNI space (by FNIRT). We use all the four pseudo-modalities to develop  
86 the convnet models. WM and GM segmentations are in 1.5mm MNI space as provided by the PAC  
87 2019 organisers, and the preprocessing pipeline is described in (Cole and Franke, 2017).

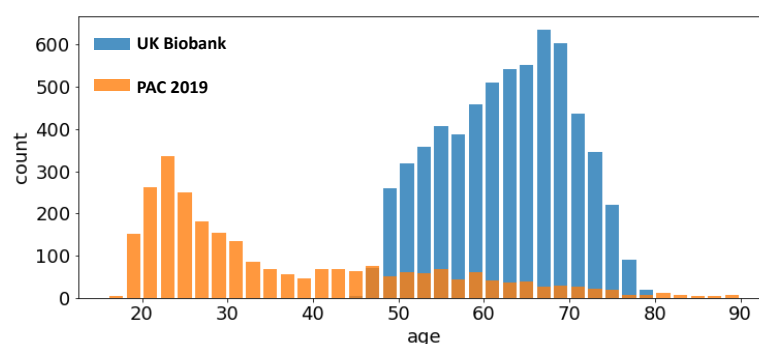
88 For linearly and non-linearly registered modalities, the input images are cropped to retain the central  
89 160x192x160 voxels, which is the same as what we had done with UK Biobank data. The WM and  
90 GM modalities are cropped in the central 96x128x96 voxels.

## 91 2.2 UK Biobank

92 UK Biobank brain imaging data consists of multimodal brain scans from a predominantly healthy  
93 cohort (Miller et al., 2016). Currently (year 2020) there are about 40,000 subjects released for research,  
94 and the number will eventually reach 100,000 (Smith et al., 2020a). In our previous study, we reported  
95 SFCN trained and tested on the initial 14,503 structural MRI brain images (Peng et al., 2019), and  
96 released the pretrained model in a GitHub repository ([https://github.com/ha-ha-ha-han/UKBiobank\\_deep\\_pretrain](https://github.com/ha-ha-ha-han/UKBiobank_deep_pretrain)). In this study, we mainly focus on optimising pipelines and models  
98 for PAC 2019, and most of the models are initialised randomly and then trained with the PAC 2019  
99 data unless otherwise stated. To apply transfer learning, we also use 5698 UK Biobank T1 images to  
100 pretrain a model, and then use the trained weights as initialisations for finetuning five models in the  
101 PAC 2019 dataset (see details in the section Experiments and Results – Transfer Learning).

102 The UK Biobank preprocessing pipeline can be found in (Alfaro-Almagro et al., 2018), and the UKB  
103 data release includes preprocessed data, so that researchers do not need to re-run the preprocessing  
104 pipeline. Models are trained/validated/tested separately. The inputs are in 1mm MNI space, cropped  
105 for the central 160x192x160 voxels to reduce GPU memory required.

106



**Figure 1. Age distribution of different datasets.** The UK Biobank (blue bars) and the PAC 2019 (orange bars) differ in age range and number of subjects.

## 107 2.3 Difference between UK Biobank and PAC 2019

108 UK Biobank and PAC 2019 datasets differ in age distribution and number of subjects. A summary of  
109 the statistics of both datasets (mean and standard deviation of age distribution, and number of subjects)  
110 is shown in Table 1 and visualised in Figure 1. The PAC 2019 dataset has a significantly smaller  
111 number of subjects and larger age range. Moreover, PAC 2019 contains multisite data with different  
112 data quality and scanner configurations. All these factors make the prediction task more difficult in  
113 PAC 2019 than UK Biobank.

114 Note that the test set labels are not available to the participants in the PAC 2019 challenge. This setup  
115 of a ‘true’ test set prevents the competition participants from the risk of accidental data leakage. During  
116 the competition, the prediction results were allowed to be uploaded only once, and then the  
117 performance metric was evaluated automatically. Therefore, no hyperparameter adjustment could be  
118 made for the testing process to elaboratively overfit the test set. In summary, we believe the results in  
119 the test set are an objective measurement of model performance in an unknown dataset with a similar  
120 age and site distribution.

121

Dataset	Age Range (yrs)	Age (yrs) Mean±STD	Number of Subject		Number of Site
			Training/Validation/Test	Total	
UK Biobank	44 - 80	62.7±7.5	5698 / 518 / -	6216	2
PAC 2019	17 - 90	35.9±16.2	2198 / 440 / 660	2638 with label + 660 without label	17

**Table 1. Difference in age distribution between PAC 2019 used in this study and UK Biobank dataset used in (Peng et al., 2019).**

## 122 3 Method

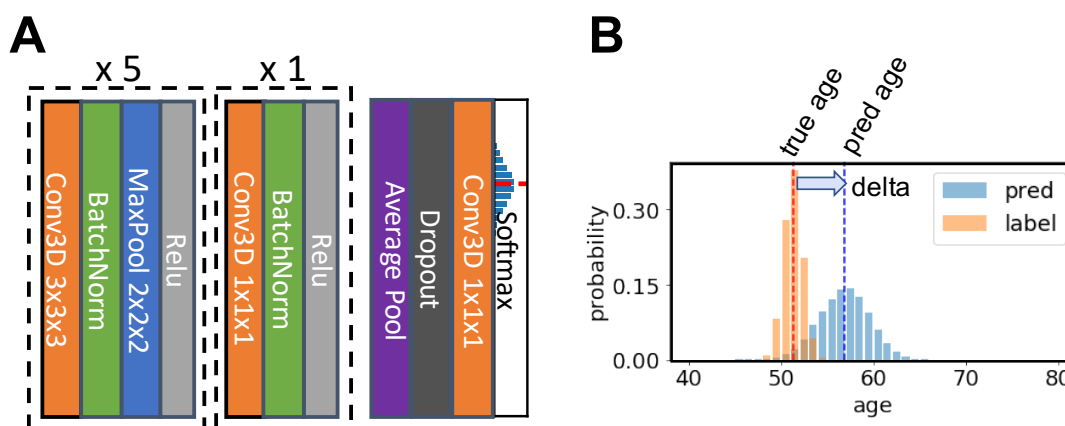
### 123 3.1 Model

124 The backbone of our method is the lightweight fully convolutional neural network architecture, Simple  
125 Fully Convolutional Neural Network (SFCN), that we proposed in (Peng et al., 2019). We briefly  
126 summarise the key aspects of the model and the adjustment for PAC 2019 here.

127 The SFCN model architecture is shown in Figure 2 (reproduced from the original work by (Peng et al.,  
128 2019)). The model consists of seven convolution blocks. Each of the first five blocks consist of a 3x3x3  
129 3D convolution layer, a batch normalisation layer, a max pooling layer, and a ReLU activation layer.  
130 The key facet of this architecture is that the model downsamples the input every time after a  
131 convolution layer. As a result, the spatial dimension is reduced quickly as the layer goes deeper, and it  
132 takes only five blocks to reduce the input data size from 160x192x160 to 5x6x5 (voxels). This simple

133 design saves GPU memory and reduced the number trainable weights. The sixth block is similar but  
 134 without a max pooling layer and uses a 1x1x1 3D convolution layer to increase non-linearity without  
 135 changing feature map spatial dimensions. The resulting 5x6x5 feature map is pooled by an average  
 136 pooling layer and then projected to the output layer with a linear transformation (i.e. fully connected  
 137 layer). For convenience of implementation, the fully connected layer is also treated as an 1x1x1  
 138 Conv3D in a 1x1x1 input ‘feature map’.

139 The input size is 160x192x160 voxels for both T1 non-linearly registered brains and linearly registered  
 140 brains, and 96x128x96 voxels for both WM and GM for PAC 2019. Note that the model is fully  
 141 convolutional; therefore it can take different input sizes without modifying the architecture. The feature  
 142 map size before the average pooling layer in the final block is 5x6x5 for the input size 160x192x160,  
 143 and 3x4x3 for the input size 96x128x96.



**Figure 2. Illustration of the core network for the Simple Fully Convolutional Neural Network (SFCN) model. A) SFCN model architecture. B) An example of soft labels and output probabilities. The figure is reproduced from (Peng et al., 2019) under CC-BY-NC-ND 4.0.**

144

### 145 3.2 Model Output and Loss function

146 We treat the regression as a soft classification problem. In this set-up, the label of the age is not treated  
 147 as a single number, but a discretized Gaussian probability distribution centred at the true age. The  
 148 output of the model is also a probability distribution. Kullback-Leibler divergence is used to measure  
 149 the similarity between the two probabilities.

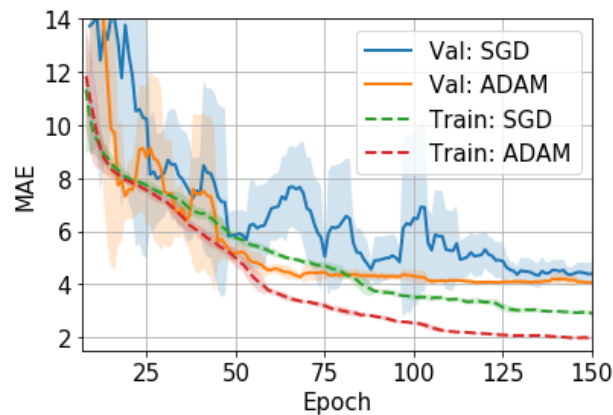
150 The output is 40 digits standing for 40 age bins for the UK Biobank data. Each age bin covers a 1-year  
 151 range. The number of age bins is 38 for trained-from-scratch models for PAC 2019, each of which  
 152 covers a 2-year range. The sigma of the Gaussian distribution for the labels is set to be the size of one  
 153 age bin (i.e. 1 year for UK Biobank and 2 years for PAC 2019). The final age prediction is the average  
 154 of all the age bins weighted by the output probability.

155 For models pretrained in UK Biobank and finetuned in PAC 2019, the number of output age bins is set  
156 to 40 to reduce coding effort (although the bins stand for different age ranges).

### 157 3.3 Hyper parameter, optimiser choice and training

158 Hyper parameters are tuned with the validation set. We also evaluate different optimizers, namely,  
159 ADAM and SGD. In UK Biobank we find ADAM easily overfits the model and thus performs worse  
160 than SGD (Peng et al., 2019). However, in PAC 2019, we find that ADAM, although it overfits more  
161 than SGD (as measured by the val-train gap in Figure 3), performs slightly better than SGD in the  
162 validation set. Also, ADAM is observed to be more stable during the training process for the PAC 2019  
163 dataset (as shown in Figure 3), so that we use ADAM for PAC 2019 for the rest of our experiments.

164 The validation set is used to evaluate model performance after every epoch (i.e. one iteration through  
165 the full dataset) in the training set, and the model weights for the best validation performance within  
166 150 epochs are chosen for testing.



**Figure 3. Training curves for the SGD and ADAM optimisers in PAC 2019 data.** The curves are smoothed with a 7-step averaging window. The shading areas show the standard deviation within the window.

167 Data augmentation and weight regularization are important to achieve the best prediction accuracy and  
168 to reduce overfitting. We use the same augmentation and regularization strategy as specified in detail  
169 in (Peng et al., 2019) for all experiments reported in this work: voxel shifting, mirroring and dropout.

## 170 4 Experiments and Results

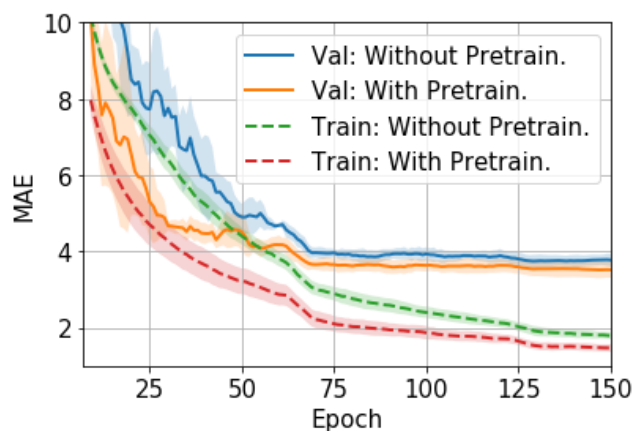
171 To achieve accurate brain age prediction, we use several techniques in the competition setup besides  
172 the lightweight SFCN model, the regularization and the data augmentations. For a single model, we  
173 applied transfer learning to boost the single model prediction accuracy. We also train multiple models  
174 using different (pseudo-)modalities to form an ensemble for better performance. As summarised in  
175 Table 2, we find that the best ensemble uses all the modalities. While transfer learning stably achieves  
176 better single-model performance, only 5 out of 45 models in the final ensemble are transferred from  
177 UK Biobank, due to the limit of time and computational power. The details of the experiments and the  
178 results are described below.

### 179 4.1 Transfer learning

180 To test how pretraining in the large UK Biobank dataset can help smaller datasets such as PAC 2019,  
181 we compare the performance of models that are pretrained-and-finetuned and those trained-from-  
182 scratch using the PAC 2019 data only.

183 The finetuning process and all the hyperparameters are the same with the trained-from-scratch ones  
184 except for the initialisation of model weights. For the pretraining, an SFCN model is trained with 5698  
185 UK Biobank subjects using the methods specified in (Peng et al., 2019) and achieving validation MAE  
186 = 2.20 yrs in UK Biobank dataset. This MAE is slightly worse than the reported value due to the smaller  
187 training dataset size we use. The trained weights are then used to initialise models that are finetuned  
188 with the PAC 2019 dataset. There are five models initialised with the same weights, and then trained  
189 with different train-validation split under a five-fold cross validation scheme using the PAC 2019  
190 training data. as shown in Figure 4, the five finetuned models achieve a mean MAE of  $3.69 \pm 0.19$  yrs  
191 (mean $\pm$ STD), which is 0.22 years better than the randomly initialised models (MAE =  $3.91 \pm 0.13$  yrs,  
192 mean $\pm$ STD). The pretrained models also converge faster. This result shows that initialising models  
193 with pretrained weights from UK Biobank can help achieve better performance in small datasets, even  
194 using a naïve finetuning protocol.

195



**Figure 4. Training curves for transfer learning.** The curves are averaged by five models trained with five-fold cross-validation splitting, and then smoothed with a 7-step averaging window. The shading areas show the standard deviation within the window.

## 196 4.2 Performance of different (pseudo-)modalities and model ensembles

197 Different T1-derived data contain distinct information regarding brain ageing. We find that averaging  
198 predictions with different pseudo-modalities (outputs from distinct pre-processing approaches applied  
199 to the same original input data modality, here T1) is an effective method to utilise the independent  
200 information to achieve the overall best ensemble performance. We train and test 10 models (from  
201 scratch, no pretraining) in each pseudo-modality, namely, T1 data linearly registered to the MNI space  
202 (Lin), raw T1 data nonlinearly registered to the MNI space (NonLin), segmented grey matter (GM)  
203 and white matter (WM) volumes. Lin and NonLin modalities are preprocessed by us, and GM and WM  
204 are provided by the organiser. Models are randomly initialized (with different random seeds). As shown  
205 in Table 2, models trained with Lin, NonLin and GM achieve comparable MAEs ranging from 3.89 to

206 3.93 years, which are all better than the MAE for WM (4.19 years), and is in accordance with our  
 207 previous findings (Peng et al., 2019).

208 We show in our previous work (Peng et al., 2019) that, even though with comparable MAEs, brain-  
 209 PADs contain different information from different pseudo-modalities. This result is consolidated in the  
 210 PAC 2019 dataset using the left-out validation set (not used in cross-validation) in Figure 5. Models  
 211 with the same modalities show higher correlation for the brain-PAD prediction.

212 To achieve the best performance in the challenge, we use all four pseudo-modalities to form an  
 213 ensemble. For every pseudo-modality, there are 10 models initialised randomly and trained separately  
 214 with different train/validation splits. For the Lin modality, 5 additional models are pretrained in UK  
 215 Biobank and finetuned in PAC 2019, as previously mentioned, adding up to 45 models in total. All

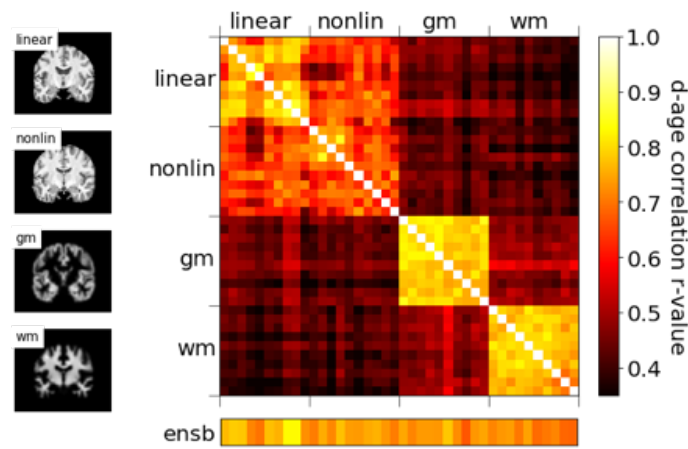
Modality	Performance			
	Single Model		Ensemble	
	MAE (yrs)	r value	MAE (yrs)	r value
Raw, linearly registered, Pretrained with UK Biobank x 5	3.69±0.08	0.946±0.006	3.22	0.960
Raw, linearly registered x 10	3.91±0.13	0.935±0.007	3.48	0.951
Raw, non-linearly registered x 10	3.89±0.16	0.937±0.006	3.40	0.957
Grey matter x 10	3.93±0.13	0.948±0.003	3.54	0.957
White matter x 10	4.19±0.09	0.937±0.003	3.74	0.951
All 45 models	3.95±0.19	0.940±0.007	2.98	0.971

**Table 2. Performance of model ensembles with different pseudo modalities in PAC 2019.** 5 models are initialized with pretrained weights and then finetuned with linearly registered brains. For all other experiments, 10 models are trained from scratch for each modality and used to predict brain age individually. The mean and the standard deviation of the single model performances are computed within each modality.



216 models are trained separately, and make predictions independently. For every subject, mean and  
 217 standard deviation (STD) are computed for the 45 age predictions, and the predictions deviating more  
 218 than  $\lambda$ -STD from the mean are treated as outliers ( $\lambda$  is a coefficient of our choice), and the final  
 219 prediction is the new average of the rest predictions.  $\lambda$  is set to be 1.1 to optimise the performance in  
 220 the left-out validation set, which makes the ensemble performance slightly biased towards this  
 221 ‘validation’ set. This strategy achieves MAE = 2.98 yrs in the left-out validation set and MAE = 2.90  
 222 yrs in the test set, as shown in Table 3. Our result in the test set ranks the first for the first goal of PAC  
 223 2019 (best MAE), and is 0.18/0.42 years better than the second/third place (MAE: Ours = 2.904 yrs;  
 224 Second Place = 3.086 yrs; Third Place = 3.328 yrs).

225

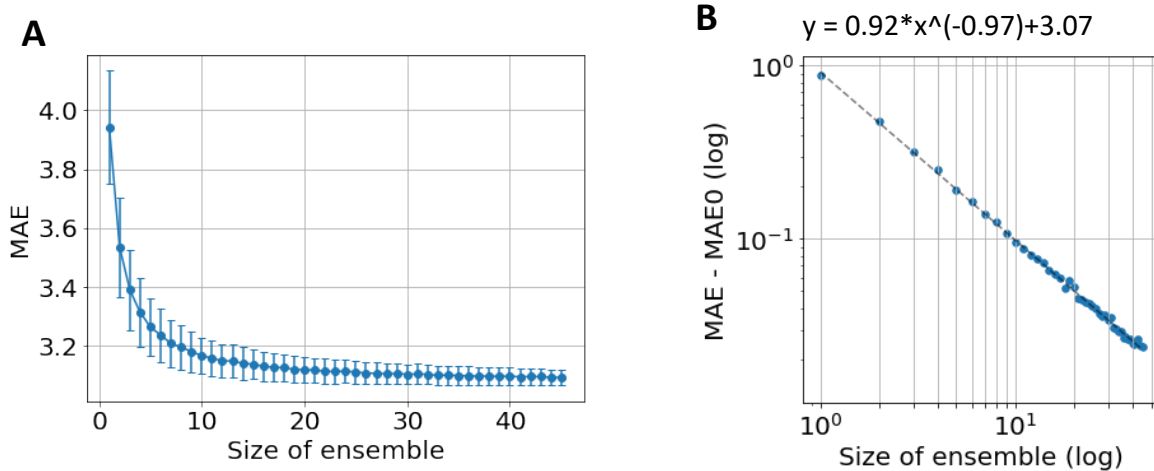


**Figure 5. Correlations of predicted brain age difference (d-age) between different models, showing similar results as (Peng et al., 2019).**

226 In our previous work (Peng et al., 2019), we showed that independent predictions are important to form  
 227 a good ensemble. Here, we further show that a sufficiently large number of models is also important  
 228 for good ensemble performance. To demonstrate this, we explore the ensemble performance with  
 229 different number of models, as summarised in Figure 6. Ensembles are randomly formed using some  
 230 of the 45 trained models (replacement allowed) and predictions are made using the mean without  
 231 excluding outliers. As the number of models increases, the MAE decreases and finally saturate. A  
 232 power law can be fitted to empirically describe the quantitative relationship between the size of  
 233 ensemble and the MAE, as shown in Figure 6B. A ‘critical point’ of MAE of 3.07 yrs is estimated, and  
 234 can be interpreted as the ideal MAE if we can increase the number of models to infinity. This empirical  
 235 observation suggests that simply increasing ensemble size will result in only limited performance gain.

236 The ‘critical’ MAE is worse than the actual MAE we get from the all the models. This is because the  
237 bootstrap process allows replacement, i.e. the same model is allowed to be selected more than once,  
238 which reduces the independent information gathered from the ensemble.

### 239 4.3 Bias correction



**Figure 6. Ensemble performance with different number of models.** **A)** Average performance in MAE with different number of models used by ensemble. The mean and standard deviation come from 1000-time bootstraps. **B)** The fitted line of a power law.  $MAE_0$  is the critical point if an infinite number of models are used to form the ensemble.

240 We follow (Smith et al., 2019) and (Peng et al., 2019) to fit a straight line between the predicted brain-  
241 PAD and the ground truth age in the left-out validation set, and then apply the fitted parameters (slope  
242 and intercept) to bias-correct predictions in the test set whose labels are unknown. We correct the bias  
243 for the ensemble predictions rather than for every single model.

244 For the validation set, this linear regression method reduces the Spearman’s r-value (between delta and  
245 age) from -0.44 to -0.06 with a small increase (0.03 years) in the MAE. The generalization to the test  
246 set reduces the Spearman’s r-value from -0.39 to 0.03, with a small increase of 0.05 years in the MAE  
247 (from MAE = 2.90 to MAE = 2.95). This result is summarised in Table 3.

248 The result in the test set achieves the first place for the second goal of the competition (smallest MAE  
249 with sufficiently small Spearman’s r-value between brain-PAD and the true age), and it leads by a large  
250 margin (MAE: Ours = 2.950 yrs; Second Place = 3.799 yrs; Third Place = 3.924 yrs).

251

Model	Performance		Performance with Bias Correction	
	MAE (years)	Spearman Correlation d-age vs age	MAE (years)	Spearman Correlation d-age vs age
45 Model Ensemble (Left-out validation set)	2.98	-0.44	3.01	-0.06
45 Model Ensemble (PAC Test Set)	2.90	-0.39	2.95	-0.03

**Table 3. Bias correction results.**

## 253 5 Discussion and conclusion

254 We note that different datasets may require distinct hyperparameters and optimisers for optimal  
255 performance for a deep learning algorithm. For example, we showed in our previous study that ADAM  
256 easily overfits the model and thus performs worse than SGD in UK Biobank data (Peng et al., 2019).  
257 In this study, we find ADAM works comparable or even slightly better than SGD in PAC 2019  
258 validation data. We have not fully explored the mechanism behind this empirical difference. One can  
259 assume that PAC 2019 is a more difficult dataset for deep learning models to optimize, due to the  
260 multi-site origin and inhomogeneous data quality, and this may be the reason why ADAM performs  
261 better in PAC 2019; it has been shown to be a more powerful optimizer for other problems (Kingma  
262 and Ba, 2014). For future studies, it may be beneficial to explore and choose different optimisers for  
263 different datasets even for similar tasks.

264 Despite additional hyperparameter tuning, we have shown that the SFCN method together with the  
265 data augmentation and model regularisation methods are generalisable outside the UK Biobank dataset.  
266 However, this ‘generalisability’ requires retraining or finetuning in the targeting dataset, and may not  
267 be feasible for smaller datasets (e.g. a dataset with 100-subject). Also, although PAC 2019 provides a  
268 true measurement for generalisability of models to unseen data (because the test set labels are hidden  
269 from the participants), this does not guarantee the generalisability to unseen scanning site (because the  
270 test set follows the same site and age distribution as the training set). For applications requiring site  
271 generalisability, see recent work aiming to address this specific issue (Dinsdale et al., 2020).

272 Finally, we need to point out that our choice of hyperparameters, transfer learning and the naïve  
273 ensemble strategy may not be optimal, due to the limit of time and computation power in the  
274 competition setup.

275 To conclude, we have applied the lightweight convnet - SFCN model, data augmentation,  
276 regularisation, and bias correction techniques proposed in (Peng et al., 2019) to PAC 2019 challenge  
277 and achieved leading results. Besides initialising models randomly, we have shown that initialising  
278 weights pretrained in UK Biobank achieve better single-model results for the PAC 2019 dataset (after  
279 retraining/finetuning). For ensembles with multiple models, we have shown that the best ensemble  
280 comes from a large number of models taking the input of different pseudo-modalities.

## 281 **6 Conflict of Interest**

282 *The authors declare that the research was conducted in the absence of any commercial or financial*  
283 *relationships that could be construed as a potential conflict of interest.*

## 284 **7 Author Contributions**

285 **Weikang Gong:** Conceptualization, Methodology, Software, Writing - Review & Editing. **Christian**  
286 **F. Beckmann:** Conceptualization, Writing - Review & Editing, Methodology, Funding acquisition,  
287 Supervision. **Andrea Vedaldi:** Conceptualization, Writing - Review & Editing, Methodology,  
288 Funding acquisition, Supervision. **Stephen M. Smith:** Conceptualization, Writing - Review & Editing,  
289 Methodology, Funding acquisition, Supervision. **Han Peng:** Conceptualization, Methodology,  
290 Software, Writing- Original draft preparation, Writing - Review & Editing, (Co-)Supervision.

## 291 **8 Funding**

292 This project is supported by the DeepMedicine project in the Oxford Martin School and the Innovative  
293 Medicines Initiative 2 Joint Undertaking under grant agreement No 777394 (for AIMS-2-TRIALS)  
294 which receives support from the European Union's Horizon 2020 research and innovation programme  
295 and EFPIA and AUTISM SPEAKS, Autistica, SFARI. We are also grateful for funding from the  
296 Wellcome Trust (215573/Z/19/Z, 203139/Z/16/Z).

## 297 **9 Acknowledgments**

298 This research has been conducted in part using the UK Biobank Resource under Application 8107. We  
299 are grateful to UK Biobank for making the data available, and to all UK Biobank study participants,  
300 who generously donated their time to make this resource possible. Computation used the Oxford  
301 Biomedical Research Computing (BMRC) facility, a joint development between the Wellcome Centre  
302 for Human Genetics and the Big Data Institute supported by Health Data Research UK and the NIHR  
303 Oxford Biomedical Research Centre.

## 304 **10 Data Availability Statement**

305 The PAC 2019 dataset consists of several public available datasets and a few datasets provided by the  
306 organiser. Interested researchers can apply for the access to the public available datasets as specified  
307 in (Cole et al., 2017) and need to contact the PAC 2019 organisers for the rest of the sites. The UK  
308 Biobank dataset is accessible upon applications via the website: <https://www.ukbiobank.ac.uk/>

## 309 **11 References**

- 310 Alfaro-Almagro, F., Jenkinson, M., Bangerter, N.K., Andersson, J.L.R., Griffanti, L., Douaud, G.,  
311 Sotiropoulos, S.N., Jbabdi, S., Hernandez-Fernandez, M., Vallee, E., Vidaurre, D., Webster, M.,  
312 McCarthy, P., Rorden, C., Daducci, A., Alexander, D.C., Zhang, H., Dragonu, I., Matthews,  
313 P.M., Miller, K.L., Smith, S.M., 2018. Image processing and Quality Control for the first 10,000  
314 brain imaging datasets from UK Biobank. *Neuroimage* 166, 400–424.  
315 <https://doi.org/10.1016/j.neuroimage.2017.10.034>
- 316 Brown, T.T., Kuperman, J.M., Chung, Y., Erhart, M., McCabe, C., Hagler, D.J., Venkatraman, V.K.,  
317 Akshoomoff, N., Amaral, D.G., Bloss, C.S., Casey, B.J., Chang, L., Ernst, T.M., Frazier, J.A.,  
318 Gruen, J.R., Kaufmann, W.E., Kenet, T., Kennedy, D.N., Murray, S.S., Sowell, E.R., Jernigan,  
319 T.L., Dale, A.M., 2012. Neuroanatomical assessment of biological maturity. *Curr. Biol.* 22,  
320 1693–1698. <https://doi.org/10.1016/j.cub.2012.07.002>
- 321 Bzdok, D., Varoquaux, G., Steyerberg, E.W., 2020. Prediction, not association, paves the road to  
322 precision medicine. *JAMA Psychiatry*. <https://doi.org/10.1001/jamapsychiatry.2020.2549>
- 323 Cole, J., Raffel, J., Friede, T., Eshaghi, A., Brownlee, W., Chard, D., De Stefano, N., Enzinger, C.,  
324 Pirpamer, L., Filippi, M., Gasperini, C., Rocca, M., Rovira, A., Ruggieri, S., Sastre-Garriga, J.,  
325 Stromillo, M., Uitdehaag, B., Vrenken, H., Barkhof, F., Nicholas, R., Ciccarelli, O., 2019.  
326 Accelerated brain ageing and disability in multiple sclerosis. *bioRxiv* 584888.  
327 <https://doi.org/10.1101/584888>
- 328 Cole, J.H., Franke, K., 2017. Predicting Age Using Neuroimaging: Innovative Brain Ageing  
329 Biomarkers. *Trends Neurosci.* 40, 681–690. <https://doi.org/10.1016/j.tins.2017.10.001>
- 330 Cole, J.H., Poudel, R.P.K., Tsagkrasoulis, D., Caan, M.W.A., Steves, C., Spector, T.D., Montana, G.,  
331 2017. NeuroImage Predicting brain age with deep learning from raw imaging data results in a  
332 reliable and heritable biomarker. *Neuroimage* 163, 115–124.  
333 <https://doi.org/10.1016/j.neuroimage.2017.07.059>
- 334 Cole, J.H., Ritchie, S.J., Bastin, M.E., Valdés Hernández, M.C., Muñoz Maniega, S., Royle, N.,  
335 Corley, J., Pattie, A., Harris, S.E., Zhang, Q., Wray, N.R., Redmond, P., Marioni, R.E., Starr,  
336 J.M., Cox, S.R., Wardlaw, J.M., Sharp, D.J., Deary, I.J., 2018. Brain age predicts mortality.  
337 *Mol. Psychiatry* 23, 1385–1392. <https://doi.org/10.1038/mp.2017.62>
- 338 Dinsdale, N.K., Jenkinson, M., Namburete, A.I.L., 2020. Deep Learning-Based Unlearning of  
339 Dataset Bias for MRI Harmonisation and Confound Removal. *bioRxiv* 2020.10.09.332973.  
340 <https://doi.org/10.1101/2020.10.09.332973>
- 341 Dosenbach, N.U.F., Nardos, B., Cohen, A.L., Fair, D.A., Power, J.D., Church, J. a, Nelson, S.M.,  
342 Wig, G.S., Vogel, A.C., Coalson, R.S., Jr, J.R.P., Barch, D.M., 2010. Prediction of Individual  
343 Brain Maturity Using fMRI. *Science* (80-. ). 329, 1358–1361.  
344 <https://doi.org/http://dx.doi.org/10.1016/j.jrurstud.2015.06.009>
- 345 Franke, K., Ziegler, G., Klöppel, S., Gaser, C., 2010. Estimating the age of healthy subjects from T1-  
346 weighted MRI scans using kernel methods: Exploring the influence of various parameters.  
347 *Neuroimage* 50, 883–892. <https://doi.org/10.1016/j.neuroimage.2010.01.005>
- 348 Kaufmann, T., van der Meer, D., Doan, N.T., Schwarz, E., Lund, M.J., Agartz, I., Alnæs, D., Barch,  
349 D.M., Baur-Streubel, R., Bertolino, A., Bettella, F., Beyer, M.K., Bøen, E., Borgwardt, S.,

- 350 Brandt, C.L., Buitelaar, J., Celius, E.G., Cervenka, S., Conzelmann, A., Córdova-Palomera, A.,  
351 Dale, A.M., de Quervain, D.J.F., Carlo, P., Djurovic, S., Dørum, E.S., Eisenacher, S.,  
352 Elvsåshagen, T., Espeseth, T., Fatouros-Bergman, H., Flyckt, L., Franke, B., Frei, O., Haatveit,  
353 B., Håberg, A.K., Harbo, H.F., Hartman, C.A., Heslenfeld, D., Hoekstra, P.J., Høgestøl, E.A.,  
354 Jernigan, T.L., Jonassen, R., Jönsson, E.G., Kirsch, P., Kłoszewska, I., Kolskår, K.K., Landrø,  
355 N.I., Hellard, S., Lesch, K.-P., Lovestone, S., Lundervold, A., Lundervold, A.J., Maglanoc,  
356 L.A., Malt, U.F., Mecocci, P., Melle, I., Meyer-Lindenberg, A., Moberget, T., Norbom, L.B.,  
357 Nordvik, J.E., Nyberg, L., Oosterlaan, J., Papalino, M., Papassotiropoulos, A., Pauli, P., Pergola,  
358 G., Persson, K., Richard, G., Rokicki, J., Sanders, A.-M., Selbæk, G., Shadrin, A.A., Smeland,  
359 O.B., Soininen, H., Sowa, P., Steen, V.M., Tsolaki, M., Ulrichsen, K.M., Vellas, B., Wang, L.,  
360 Westman, E., Ziegler, G.C., Zink, M., Andreassen, O.A., Westlye, L.T., 2019. Common brain  
361 disorders are associated with heritable patterns of apparent aging of the brain. *Nat. Neurosci.* 22,  
362 1617–1623. <https://doi.org/10.1038/s41593-019-0471-7>
- 363 Kingma, D.P., Ba, J.L., 2014. Adam: A method for stochastic optimization. *arXiv Prepr.*  
364 [arXiv1412.6980](https://arxiv.org/abs/1412.6980).
- 365 Lanka, P., Rangaprakash, D., Dretsch, M.N., Katz, J.S., Denney, T.S., Deshpande, G., 2019.  
366 Supervised machine learning for diagnostic classification from large-scale neuroimaging  
367 datasets. *Brain Imaging Behav.* 1–39. <https://doi.org/10.1007/s11682-019-00191-8>
- 368 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *Nature* 521, 436–444.  
369 <https://doi.org/10.1038/nature14539>
- 370 Levakov, G., Rosenthal, G., Shelef, I., Raviv, T.R., Avidan, G., 2020. From a deep learning model  
371 back to the brain—Identifying regional predictors and their relation to aging. *Hum. Brain Mapp.*  
372 *hbm.25011*. <https://doi.org/10.1002/hbm.25011>
- 373 Liu, M., Zhang, J., Adeli, E., Shen, D., 2018. Landmark-based deep multi-instance learning for brain  
374 disease diagnosis. *Med. Image Anal.* 43, 157–168. <https://doi.org/10.1016/j.media.2017.10.005>
- 375 Miller, K.L., Alfaro-Almagro, F., Bangerter, N.K., Thomas, D.L., Yacoub, E., Xu, J., Bartsch, A.J.,  
376 Jbabdi, S., Sotiropoulos, S.N., Andersson, J.L.R., Griffanti, L., Douaud, G., Okell, T.W., Weale,  
377 P., Dragonu, I., Garratt, S., Hudson, S., Collins, R., Jenkinson, M., Matthews, P.M., Smith,  
378 S.M., 2016. Multimodal population brain imaging in the UK Biobank prospective  
379 epidemiological study. *Nat. Neurosci.* 19, 1523–1536. <https://doi.org/10.1038/nn.4393>
- 380 Neeb, H., Zilles, K., Shah, N.J., 2006. Fully-automated detection of cerebral water content changes:  
381 Study of age- and gender-related H2O patterns with quantitative MRI. *Neuroimage* 29, 910–  
382 922. <https://doi.org/10.1016/j.neuroimage.2005.08.062>
- 383 Peng, H., Gong, W., Beckmann, C.F., Vedaldi, A., Smith, S.M., 2019. Accurate brain age prediction  
384 with lightweight deep neural networks. *bioRxiv* 879346.
- 385 Raghu, M., Zhang, C., Kleinberg, J., Bengio, S., 2019. Transfusion: Understanding Transfer  
386 Learning for Medical Imaging, in: *Advances in Neural Information Processing Systems* 32.  
387 Curran Associates, Inc., pp. 3347–3357.
- 388 Rosenberg, M.D., Casey, B.J., Holmes, A.J., 2018. Prediction complements explanation in  
389 understanding the developing brain. *Nat. Commun.* 9, 1–13. <https://doi.org/10.1038/s41467->

390 018-02887-9

391 Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla,  
392 A., Bernstein, M., Berg, A.C., Fei-Fei, L., 2015. ImageNet Large Scale Visual Recognition  
393 Challenge. *Int. J. Comput. Vis.* 115, 211–252. <https://doi.org/10.1007/s11263-015-0816-y>

394 Scheinost, D., Noble, S., Horien, C., Greene, A.S., Lake, E.M., Salehi, M., Gao, S., Shen, X.,  
395 O’Connor, D., Barron, D.S., Yip, S.W., Rosenberg, M.D., Constable, R.T., 2019. Ten simple  
396 rules for predictive modeling of individual differences in neuroimaging. *Neuroimage*.  
397 <https://doi.org/10.1016/j.neuroimage.2019.02.057>

398 Shmueli, G., 2010. To explain or to predict? *Stat. Sci.* 25, 289–310. <https://doi.org/10.1214/10->  
399 STS330

400 Smith, S.M., Douaud, G., Chen, W., Hanayik, T., Alfaro-Almagro, F., Sharp, K., Elliott, L.T., 2020a.  
401 Enhanced Brain Imaging Genetics in UK Biobank. *bioRxiv* 2020.07.27.223545.  
402 <https://doi.org/10.1101/2020.07.27.223545>

403 Smith, S.M., Elliott, L.T., Alfaro-Almagro, F., McCarthy, P., Nichols, T.E., Douaud, G., Miller,  
404 K.L., 2020b. Brain aging comprises many modes of structural and functional change with  
405 distinct genetic and biophysical associations. *Elife* 9, 1–28. <https://doi.org/10.7554/eLife.52677>

406 Smith, S.M., Jenkinson, M., Woolrich, M.W., Beckmann, C.F., Behrens, T.E.J., Johansen-Berg, H.,  
407 Bannister, P.R., De Luca, M., Drobnjak, I., Flitney, D.E., Niazy, R.K., Saunders, J., Vickers, J.,  
408 Zhang, Y., De Stefano, N., Brady, J.M., Matthews, P.M., 2004. Advances in functional and  
409 structural MR image analysis and implementation as FSL, in: *NeuroImage*. Academic Press, pp.  
410 S208–S219. <https://doi.org/10.1016/j.neuroimage.2004.07.051>

411 Smith, S.M., Vidaurre, D., Alfaro-Almagro, F., Nichols, T.E., Miller, K.L., 2019. Estimation of brain  
412 age delta from brain imaging. *Neuroimage* 200, 528–539.  
413 <https://doi.org/10.1016/j.neuroimage.2019.06.017>

414 Thomas, R.M., Gallo, S., Cerliani, L., Zhutovsky, P., El-Gazzar, A., van Wingen, G., 2020.  
415 Classifying Autism Spectrum Disorder Using the Temporal Statistics of Resting-State  
416 Functional MRI Data With 3D Convolutional Neural Networks. *Front. Psychiatry* 11, 1.  
417 <https://doi.org/10.3389/fpsy.2020.00440>

418 Zeng, L.L., Wang, H., Hu, P., Yang, B., Pu, W., Shen, H., Chen, X., Liu, Z., Yin, H., Tan, Q., Wang,  
419 K., Hu, D., 2018. Multi-Site Diagnostic Classification of Schizophrenia Using Discriminant  
420 Deep Learning with Functional Connectivity MRI. *EBioMedicine* 30, 74–85.  
421 <https://doi.org/10.1016/j.ebiom.2018.03.017>

422 Zou, L., Zheng, J., Miao, C., McKeown, M.J., Wang, Z.J., 2017. 3D CNN Based Automatic  
423 Diagnosis of Attention Deficit Hyperactivity Disorder Using Functional and Structural MRI.  
424 *IEEE Access* 5, 23626–23636. <https://doi.org/10.1109/ACCESS.2017.2762703>

425