

PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a publisher's version.

For additional information about this publication click this link.

<https://repository.ubn.ru.nl/handle/2066/233683>

Please be advised that this information was generated on 2021-10-28 and may be subject to change.



Optimising a Simple Fully Convolutional Network for Accurate Brain Age Prediction in the PAC 2019 Challenge

Weikang Gong¹, Christian F. Beckmann^{1,2}, Andrea Vedaldi³, Stephen M. Smith¹ and Han Peng^{1,2,3*}

¹ Wellcome Centre for Integrative Neuroimaging (WIN Centre for Functional MRI of the Brain), University of Oxford, Oxford, United Kingdom, ² Donders Institute for Brain, Cognition and Behaviour, Radboud University Nijmegen, Nijmegen, Netherlands, ³ Visual Geometry Group, University of Oxford, Oxford, United Kingdom

OPEN ACCESS

Edited by:

Christian Gaser,
Friedrich Schiller University
Jena, Germany

Reviewed by:

Tao Liu,
Beihang University, China
Ryu-ichiro Hashimoto,
Showa University, Japan

*Correspondence:

Han Peng
han.peng@ndcn.ox.ac.uk

Specialty section:

This article was submitted to
Computational Psychiatry,
a section of the journal
Frontiers in Psychiatry

Received: 10 November 2020

Accepted: 12 April 2021

Published: 10 May 2021

Citation:

Gong W, Beckmann CF, Vedaldi A,
Smith SM and Peng H (2021)
Optimising a Simple Fully
Convolutional Network for Accurate
Brain Age Prediction in the PAC
2019 Challenge.
Front. Psychiatry 12:627996.
doi: 10.3389/fpsy.2021.627996

Brain age prediction from brain MRI scans not only helps improve brain ageing modelling generally, but also provides benchmarks for predictive analysis methods. Brain-age delta, which is the difference between a subject's predicted age and true age, has become a meaningful biomarker for the health of the brain. Here, we report the details of our brain age prediction models and results in the Predictive Analysis Challenge 2019. The aim of the challenge was to use T1-weighted brain MRIs to predict a subject's age in multicentre datasets. We apply a lightweight deep convolutional neural network architecture, Simple Fully Convolutional Neural Network (SFCN), and combined several techniques including data augmentation, transfer learning, model ensemble, and bias correction for brain age prediction. The model achieved first place in both of the two objectives in the PAC 2019 brain age prediction challenge: Mean absolute error (MAE) = 2.90 years without bias removal (Second Place = 3.09 yrs; Third Place = 3.33 yrs), and MAE = 2.95 years with bias removal, leading by a large margin (Second Place = 3.80 yrs; Third Place = 3.92 yrs).

Keywords: predictive analysis, big data, deep learning, convolution neural network, brain age prediction, brain imaging

INTRODUCTION

Predictive analysis with data-driven machine learning algorithms brings huge promise in neuroimaging and neuroscience research. Predictive analysis can not only help disease diagnosis, such as Alzheimer's (1), Autism (2), ADHD (3) and schizophrenia (4), but also helps in formulating new hypotheses (5) and identifying new biomarkers (6). Yet, the predictive analysis paradigm brings new challenges. First, a fair way to compare predictive analysis models is needed. In predictive analysis, it is common practise to build models in a training set, and then apply the models to a test set (7, 8). It is important that no test data is used for model training or hyperparameter tuning (e.g., learning rate for gradient decent optimisations, number of layers in convnets) and to report the result objectively (9) and avoid accidental data leakage (10). Second, data is usually scarce for many diseases so that training a large deep learning model in such modest datasets is still hard (11).

Brain ageing study is a recent example of the predictive analysis paradigm (12–19). Studies showed that individuals' chronological age can be predicted accurately from brain MRI scans (14).

Brain age delta, the difference of a subject's predicted (brain) age and chronological age, is linked with a variety of biological factors within the healthy population (20), and group differences can be found in disease populations (21, 22). Yet, accurate prediction of a subject's age in healthy population is still a challenging task.

To tackle these challenges, a benchmarking platform is needed to objectively evaluate the models and strategies. Competitions have been seen in the field of computer vision [e.g., ImageNet (23)] and proved to be a valuable vehicle for pushing AI technology (9). In the field of neuroimaging, the Predictive Analysis Challenge (PAC) 2019 for brain age prediction¹ provides such opportunities for participants to train machine learning methods, and then objectively evaluate the models in a test dataset whose labels are hidden from the participants. PAC 2019 sets two objectives for brain age predictions: (1) to achieve the most accurate age prediction from brain structural MRI scans, and (2) to achieve the best accuracy while keeping the correlation between the prediction error and the ground truth age sufficiently small.

Our team "BrainAgeDifference" achieved the first places in both two objectives among 79 participating teams. Our method is largely based on our previous work (24), with adaptations made for the challenge. In this report, we will provide a detailed description of our methods for PAC 2019, including the lightweight deep convnet architecture - Simple Fully Convolutional Neural Network (SFCN), and the combined techniques including data augmentation, transfer learning, model ensemble, and bias correction. We find that the lightweight model, which has achieved the state-of-the-art results in UK Biobank, works well in the multi-centre PAC 2019 dataset with a slightly adaptation in hyperparameters. SFCN pretrained on UK Biobank data achieves better single model performance than random initialised models in the PAC 2019 dataset. In addition, model ensemble with different T1-image derived maps, and different initializations, and training/validation data splits are important to achieve the best performance for the competition.

DATASETS AND PREPROCESSING

PAC 2019

The Predictive Analytic Challenge (PAC) 2019 was to predict age from brain MRI scans. The goal of the challenge includes two parts: (1) to achieve the most accurate age prediction, as measured by mean absolute error (MAE), and (2) to achieve the best MAE while keeping the Spearman correlation r -value between the prediction error (brain age delta) and the actual age below 0.1 ($|r| < 0.1$). The dataset consists of both label-known training/validation dataset (2,638 subjects in total) and a "true" test set of 660 subjects whose labels are unknown to the competition participants. The participants had a one-time opportunity to upload their predictions in the test set to the competition server for each objective, and the MAE and the Spearman's r -value were evaluated automatically. The subjects

are from 17 different sites. Most of the data is based on (14) and a few new sites were added by the organisers, including MRI data from both 1.5 and 3 T scanners. The training set and the test set have the same age and site distribution.

PAC 2019 organisers provide three version of MRI data: (a) raw T1 brain MRI scans, (b) white matter volume segmentation (WM), and (c) grey matter volume (GM) segmentation derived from T1 data. We use all three versions to develop deep learning models. We further preprocess the raw T1 images using FSL (25) (command `fsl_anat`) to derive two different pseudo-modalities: one is brain linearly registered to standard 1 mm MNI space (by FLIRT), and the other is brain non-linearly registered to standard 1 mm MNI space (by FNIRT). We use all the four pseudo-modalities to develop the convnet models. WM and GM segmentations are in 1.5 mm MNI space as provided by the PAC 2019 organisers, and the preprocessing pipeline is described in (15).

For linearly and non-linearly registered modalities, the input images are cropped to retain the central $160 \times 192 \times 160$ voxels, which is the same as what we had done with UK Biobank data. The WM and GM modalities are cropped in the central $96 \times 128 \times 96$ voxels.

UK Biobank

UK Biobank brain imaging data consists of multimodal brain scans from a predominantly healthy cohort (26). Currently (year 2020) there are about 40,000 subjects released for research, and the number will eventually reach 100,000 (27). In our previous study, we reported SFCN trained and tested on the initial 14,503 structural MRI brain images (24), and released the pretrained model in a GitHub repository (https://github.com/ha-ha-ha-ha/UKBiobank_deep_pretrain). In this study, we mainly focus on optimising pipelines and models for PAC 2019, and most of the models are initialised randomly and then trained with the PAC 2019 data unless otherwise stated. To apply transfer learning, we also use 5,698 UK Biobank T1 images to pretrain a model, and then use the trained weights as initialisations for finetuning five models in the PAC 2019 dataset (see details in the section Experiments and Results – Transfer Learning).

The UK Biobank preprocessing pipeline can be found in (28), and the UKB data release includes preprocessed data, so that researchers do not need to re-run the preprocessing pipeline. Models are trained/validated/tested separately. The inputs are in 1 mm MNI space, cropped for the central $160 \times 192 \times 160$ voxels to reduce GPU memory required.

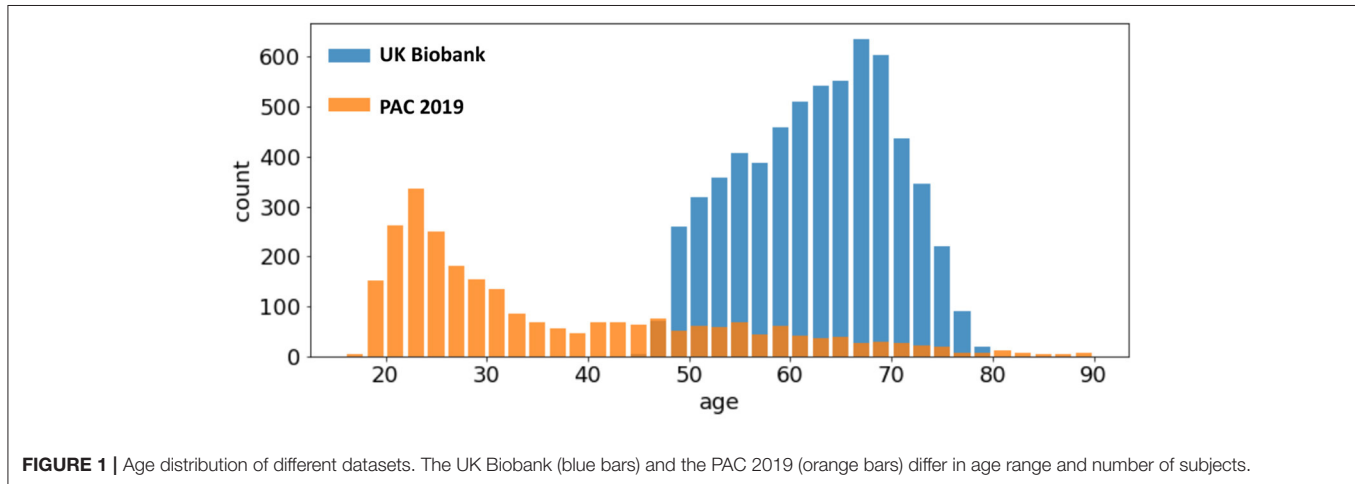
Difference Between UK Biobank and PAC 2019

UK Biobank and PAC 2019 datasets differ in age distribution and number of subjects. A summary of the statistics of both datasets (mean and standard deviation of age distribution, and number of subjects) is shown in **Table 1** and visualised in **Figure 1**. The PAC 2019 dataset has a significantly smaller number of subjects and larger age range. Moreover, PAC 2019 contains multisite data with different data quality and scanner configurations. All these

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

TABLE 1 | Difference in age distribution between PAC 2019 used in this study and UK Biobank dataset used in Peng et al. (24).

Dataset	Age range (yrs)	Age (yrs) mean \pm STD	Number of subject		Number of site
			Training/validation/test	Total	
UK Biobank	44–80	62.7 \pm 7.5	5,698/518/–	6,216	2
PAC 2019	17–90	35.9 \pm 16.2	2,198/440/660	2,638 with label + 660 without label	17



factors make the prediction task more difficult in PAC 2019 than UK Biobank.

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

METHOD

Model

The backbone of our method is the lightweight fully convolutional neural network architecture, Simple Fully Convolutional Neural Network (SFCN), that we proposed in (24). We briefly summarise the key aspects of the model and the adjustment for PAC 2019 here.

The SFCN model architecture is shown in **Figure 2** [reproduced from the original work by (24)]. The model consists of seven convolution blocks. Each of the first five blocks consist of a $3 \times 3 \times 3$ 3D convolution layer, a batch normalisation layer, a max pooling layer, and a ReLU activation layer. The key facet of this architecture is that the model downsamples the input every time after a convolution layer. As a result, the spatial dimension is reduced quickly as the layer goes deeper, and it takes only five blocks to reduce the input data size from $160 \times$

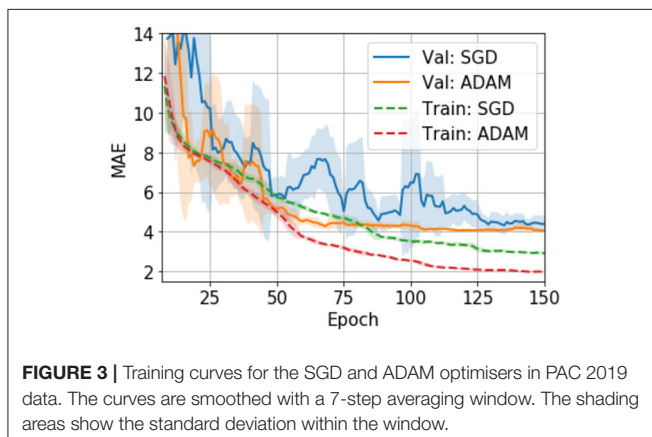
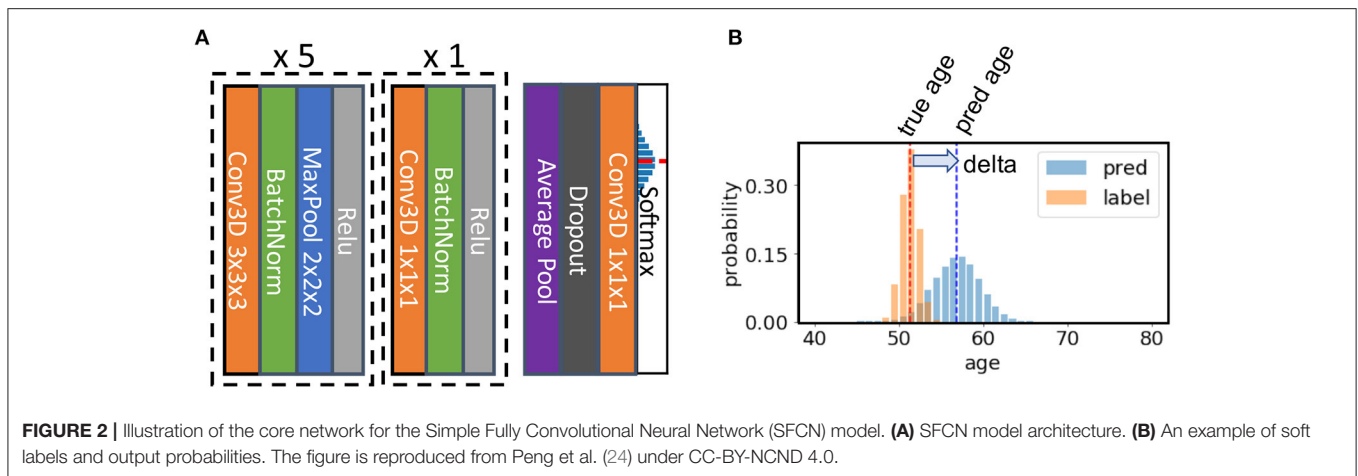
192×160 to $5 \times 6 \times 5$ (voxels). This simple design saves GPU memory and reduced the number trainable weights. The sixth block is similar but without a max pooling layer and uses a $1 \times 1 \times 1$ 3D convolution layer to increase non-linearity without changing feature map spatial dimensions. The resulting $5 \times 6 \times 5$ feature map is pooled by an average pooling layer and then projected to the output layer with a linear transformation (i.e., fully connected layer). For convenience of implementation, the fully connected layer is also treated as an $1 \times 1 \times 1$ Conv3D in a $1 \times 1 \times 1$ input “feature map.”

The input size is $160 \times 192 \times 160$ voxels for both T1 non-linearly registered brains and linearly registered brains, and $96 \times 128 \times 96$ voxels for both WM and GM for PAC 2019. Note that the model is fully convolutional; therefore it can take different input sizes without modifying the architecture. The feature map size before the average pooling layer in the final block is $5 \times 6 \times 5$ for the input size $160 \times 192 \times 160$, and $3 \times 4 \times 3$ for the input size $96 \times 128 \times 96$.

Model Output and Loss Function

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

The output is 40 digits standing for 40 age bins for the UK Biobank data. Each age bin covers a 1-year range. The number of age bins is 38 for trained-from-scratch models for PAC 2019, each of which covers a 2-year range. The sigma of the Gaussian distribution for the labels is set to be the size of one age bin (i.e.,



1 year for UK Biobank and 2 years for PAC 2019). The final age prediction is the average of all the age bins weighted by the output probability.

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

Hyper Parameter, Optimiser Choice and Training

Hyper parameters are tuned with the validation set. We also evaluate different optimizers, namely, an adaptive moment estimation optimizer (i.e., ADAM) (29) and a stochastic gradient descent optimizer (SGD) (30). In UK Biobank we find ADAM easily overfits the model and thus performs worse than SGD (24). However, in PAC 2019, we find that ADAM, although it overfits more than SGD (as measured by the val-train gap in **Figure 3**), performs slightly better than SGD in the validation set. Also, ADAM is observed to be more stable during the training process for the PAC 2019 dataset (as shown in **Figure 3**), so that we use ADAM for PAC 2019 for the rest of our experiments.

The validation set is used to evaluate model performance after every epoch (i.e., one iteration through the full dataset) in

the training set, and the model weights for the best validation performance within 150 epochs are chosen for testing.

Data augmentation and weight regularisation are important to achieve the best prediction accuracy and to reduce overfitting. We use the same augmentation and regularisation strategy as specified in detail in (24) for all experiments reported in this work: voxel shifting, mirroring and dropout.

EXPERIMENTS AND RESULTS

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

Transfer Learning

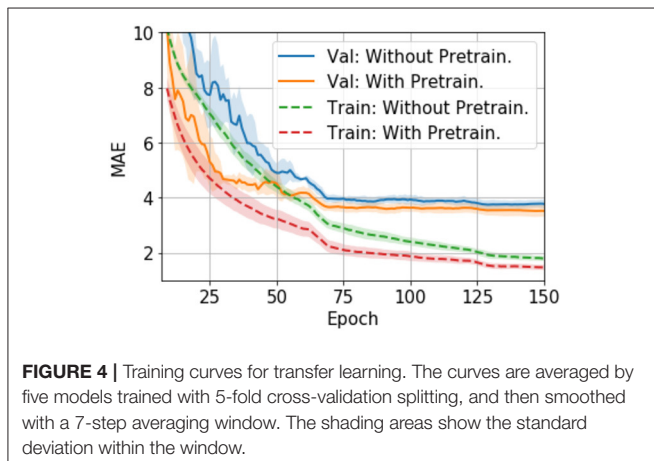
To test how pretraining in the large UK Biobank dataset can help smaller datasets such as PAC 2019, we compare the performance of models that are pretrained-and-finetuned and those trained-from-scratch using the PAC 2019 data only.

The finetuning process and all the hyperparameters are the same with the trained-from-scratch ones except for the initialisation of model weights. For the pretraining, an SFCN model is trained with 5,698 UK Biobank subjects using the methods specified in (24) and achieving validation MAE = 2.20 yrs in UK Biobank dataset. This MAE is slightly worse than the reported value due to the smaller training dataset size we use. The trained weights are then used to initialise models that are finetuned with the PAC 2019 dataset. There are five models initialised with the same weights, and then trained with different train-validation split under a 5-fold cross validation scheme

TABLE 2 | Performance of model ensembles with different pseudo modalities in PAC 2019.

Modality	Performance			
	Single model		Ensemble	
	MAE (yrs)	<i>r</i> value	MAE (yrs)	<i>r</i> value
Raw, linearly registered, Pretrained with UK Biobank × 5	3.69 ± 0.08	0.946 ± 0.006	3.22	0.960
Raw, linearly registered × 10	3.91 ± 0.13	0.935 ± 0.007	3.48	0.951
Raw, non-linearly registered × 10	3.89 ± 0.16	0.937 ± 0.006	3.40	0.957
Grey matter × 10	3.93 ± 0.13	0.948 ± 0.003	3.54	0.957
White matter × 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

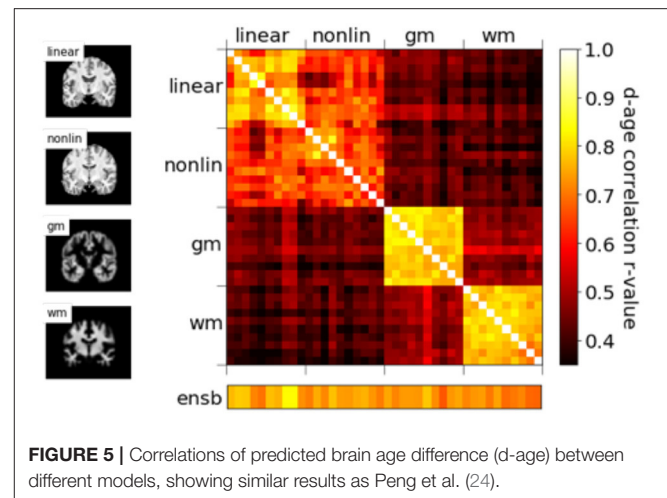
Five models are initialised with pretrained weights and then finetuned with linearly registered brains. For all other experiments, 10 models are trained from scratch for each modality and used to predict brain age individually. The mean and the standard deviation of the single model performances are computed within each modality.

**FIGURE 4** | Training curves for transfer learning. The curves are averaged by five models trained with 5-fold cross-validation splitting, and then smoothed with a 7-step averaging window. The shading areas show the standard deviation within the window.

using the PAC 2019 training data. as shown in **Figure 4**, the five finetuned models achieve a mean MAE of 3.69 ± 0.19 yrs (mean ± STD), which is 0.22 years better than the randomly initialised models (MAE = 3.91 ± 0.13 yrs, mean ± STD). The pretrained models also converge faster. This result shows that initialising models with pretrained weights from UK Biobank can help achieve better performance in small datasets, even using a naïve finetuning protocol. Finally, as is shown in **Table 2**, models initialised with UK Biobank pretrained weights result in a better-performing ensemble (MAE = 3.22 yrs) than randomly initialised models (MAE = 3.48 yrs). This result suggests that one could use UK Biobank pretrained models (as we released in the GitHub) and finetune them in a new smaller dataset, and achieve better prediction.

Performance of Different (Pseudo-)Modalities and Model Ensembles

Different T1-derived data contain distinct information regarding brain ageing. We find that averaging predictions with different

**FIGURE 5** | Correlations of predicted brain age difference (d-age) between different models, showing similar results as Peng et al. (24).

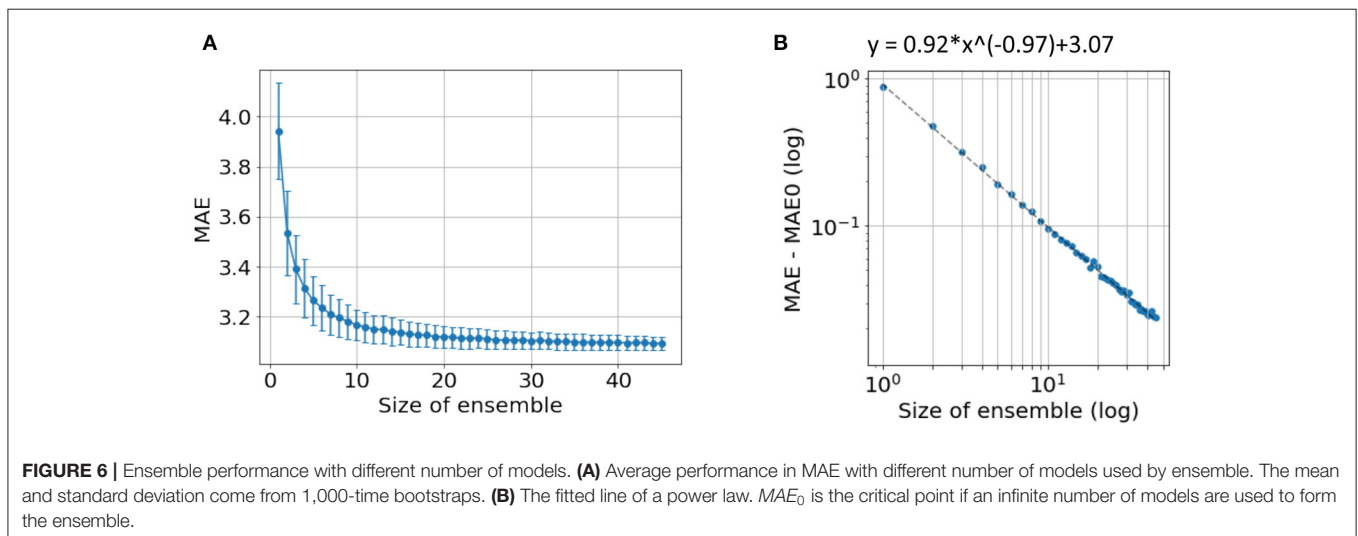
pseudo-modalities (outputs from distinct pre-processing approaches applied to the same original input data modality, here T1) is an effective method to utilise the independent information to achieve the overall best ensemble performance. We train and test 10 models (from scratch, no pretraining) in each pseudo-modality, namely, T1 data linearly registered to the MNI space (Lin), raw T1 data nonlinearly registered to the MNI space (NonLin), segmented grey matter (GM) and white matter (WM) volumes. Lin and NonLin modalities are preprocessed by us, and GM and WM are provided by the organiser. Models are randomly initialised (with different random seeds). As shown in **Table 2**, models trained with Lin, NonLin and GM achieve comparable MAEs ranging from 3.89 to 3.93 years, which are all better than the MAE for WM (4.19 years), and is in accordance with our previous findings (24).

We show in our previous work (24) that, even though with comparable MAEs, brain-PADs contain different information from different pseudo-modalities. This result is consolidated in the PAC 2019 dataset using the left-out validation set (not used in cross-validation) in **Figure 5**. Models with the same modalities show higher correlation for the brain-PAD prediction.

To achieve the best performance in the challenge, we use all four pseudo-modalities to form an ensemble. For every pseudo-modality, there are 10 models initialised randomly and trained separately with different train/validation splits. For the Lin modality, five additional models are pretrained in UK Biobank and finetuned in PAC 2019, as previously mentioned, adding up to 45 models in total. All models are trained separately, and make predictions independently. For every subject, mean and standard deviation (STD) are computed for the 45 age predictions, and the predictions deviating more than λ -STD from the mean are treated as outliers (λ is a coefficient of our choice), and the final prediction is the new average of the rest predictions. λ is set to be 1.1 to optimise the performance in the left-out validation set, which makes the ensemble performance slightly biased toward this “validation” set. This strategy achieves MAE = 2.98 yrs in the left-out validation set and MAE = 2.90 yrs in the test set, as shown in **Table 3**. Our result in the test set ranks the first for the first

TABLE 3 | Bias correction results.

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



goal of PAC 2019 (best MAE), and is 0.18/0.42 years better than the second/third place (MAE: Ours = 2.904 yrs; Second Place = 3.086 yrs; Third Place = 3.328 yrs).

We also found a sex difference in brain age prediction accuracy in the left-out validation set (females, MAE = 2.85 yrs; males, MAE = 3.09 yrs), which is in line with previous results (31).

In our previous work (24), we showed that independent predictions are important to form a good ensemble. Here, we further show that a sufficiently large number of models is also important for good ensemble performance. To demonstrate this, we explore the ensemble performance with different number of models, as summarised in **Figure 6**. Ensembles are randomly formed using some of the 45 trained models (replacement allowed) and predictions are made using the mean without excluding outliers. As the number of models increases, the MAE decreases and finally saturate. A power law can be fitted to empirically describe the quantitative relationship between the size of ensemble and the MAE, as shown in **Figure 6B**. A “critical point” of MAE of 3.07 yrs is estimated, and can be interpreted as the ideal MAE if we can increase the number of models to infinity. This empirical observation suggests that simply increasing ensemble size will result in only limited performance gain.

The “critical” MAE is worse than the actual MAE we get from all the models. This is because the bootstrap process allows replacement, i.e., the same model is allowed to be selected more

than once, which reduces the independent information gathered from the ensemble.

Bias Correction

We follow (24, 32) to fit a straight line between the predicted brain-PAD and the ground truth age in the left-out validation set, and then apply the fitted parameters (slope and intercept) to bias-correct predictions in the test set whose labels are unknown. We correct the bias for the ensemble predictions rather than for every single model.

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

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

DISCUSSION AND CONCLUSION

We note that different datasets may require distinct hyperparameters and optimisers for optimal performance

for a deep learning algorithm. For example, we showed in our previous study that ADAM easily overfits the model and thus performs worse than SGD in UK Biobank data (24). In this study, we find ADAM works comparable or even slightly better than SGD in PAC 2019 validation data. We have not fully explored the mechanism behind this empirical difference. One can assume that PAC 2019 is a more difficult dataset for deep learning models to optimise, due to the multi-site origin and inhomogeneous data quality, and this may be the reason why ADAM performs better in PAC 2019; it has been shown to be a more powerful optimizer for other problems (29). For future studies, it may be beneficial to explore and choose different optimisers for different datasets even for similar tasks.

Despite additional hyperparameter tuning, we have shown that the SFCN method together with the data augmentation and model regularisation methods are generalisable outside the UK Biobank dataset. However, this “generalisability” requires retraining or finetuning in the targeting dataset, and may not be feasible for smaller datasets (e.g., a dataset with 100-subject). Also, although PAC 2019 provides a true measurement for generalisability of models to unseen data (because the test set labels are hidden from the participants), this does not guarantee the generalisability to unseen scanning site (because the test set follows the same site and age distribution as the training set). For applications requiring site generalisability, see recent work aiming to address this specific issue (33).

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

To conclude, we have applied the lightweight convnet - SFCN model, data augmentation, regularisation, and bias correction techniques proposed in (24) to PAC 2019 challenge and achieved leading results. Besides initialising models randomly, we have shown that initialising weights pretrained in UK Biobank achieve better single-model results for the PAC 2019 dataset (after retraining/finetuning). For ensembles with multiple models, we have shown that the best ensemble comes from a large number of models taking the input of different pseudo-modalities.

REFERENCES

- Liu M, Zhang J, Adeli E, Shen D. Landmark-based deep multi-instance learning for brain disease diagnosis. *Med Image Anal.* (2018) 43:157–68. doi: 10.1016/j.media.2017.10.005
- Thomas RM, Gallo S, Cerliani L, Zhutovsky P, El-Gazzar A, van Wingen G. Classifying autism spectrum disorder using the temporal statistics of resting-state functional MRI data with 3D convolutional neural networks. *Front Psychiatry.* (2020) 11:1. doi: 10.3389/fpsy.2020.00440
- Zou L, Zheng J, Miao C, McKeown MJ, Wang ZJ. 3D CNN based automatic diagnosis of attention deficit hyperactivity disorder using functional structural MRI. *IEEE Access.* (2017) 5:23626–36. doi: 10.1109/ACCESS.2017.2762703
- Zeng LL, Wang H, Hu P, Yang B, Pu W, Shen H, et al. Multi-site diagnostic classification of schizophrenia using discriminant deep learning with functional connectivity MRI. *EBioMedicine.* (2018) 30:74–85. doi: 10.1016/j.ebiom.2018.03.017
- Shmueli G. To explain or to predict? *Stat Sci.* (2010) 25:289–310. doi: 10.1214/10-STS330
- Rosenberg MD, Casey BJ, Holmes AJ. Prediction complements explanation in understanding the developing brain. *Nat Commun.* (2018) 9:1–13. doi: 10.1038/s41467-018-02887-9
- Bzdok D, Varoquaux G, Steyerberg EW. Prediction, not association, paves the road to precision medicine. *JAMA Psychiatry.* (2021) 78:127–8. doi: 10.1001/jamapsychiatry.2020.2549
- Scheinost D, Noble S, Horien C, Greene AS, Lake EM, Salehi M, et al. Ten simple rules for predictive modeling of individual differences in neuroimaging. *Neuroimage.* (2019) 193:35–45. doi: 10.1016/j.neuroimage.2019.02.057
- LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature.* (2015) 521:436–44. doi: 10.1038/nature14539
- Lanka P, Rangaprakash D, Dretsch MN, Katz JS, Denney Jr TS, Deshpande G. Supervised machine learning for diagnostic classification from large-scale neuroimaging datasets. *Brain Imag Behav.* (2019) 14:2378–416. doi: 10.1007/s11682-019-00191-8
- Raghu M, Zhang C, Kleinberg J, Bengio S. Transfusion: understanding transfer learning for medical imaging. *arXiv [Preprint].* (2019) arXiv:1902.07208.

DATA AVAILABILITY STATEMENT

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author/s.

AUTHOR CONTRIBUTIONS

WG: conceptualisation, methodology, software, writing - review, and editing. CB, AV, and SS: conceptualisation, writing - review and editing, methodology, funding acquisition, and supervision. HP: conceptualisation, methodology, software, writing- original draft preparation, writing - review, and editing, and (co-)supervision. All authors contributed to the article and approved the submitted version.

FUNDING

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

ACKNOWLEDGMENTS

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

12. Brown TT, Kuperman JM, Chung Y, Erhart M, McCabe C, Hagler DJ, et al. Neuroanatomical assessment of biological maturity. *Curr Biol.* (2012) 22:1693–8. doi: 10.1016/j.cub.2012.07.002
13. Cole JH, Ritchie SJ, Bastin ME, Valdés Hernández MC, Muñoz Maniega S, Royle N, et al. Brain age predicts mortality. *Mol Psychiatry.* (2018) 23:1385–92. doi: 10.1038/mp.2017.62
14. Cole JH, Poudel RPK, Tsagkrasoulis D, Caan MWA, Steves C, Spector TD, et al. NeuroImage predicting brain age with deep learning from raw imaging data results in a reliable and heritable biomarker. *Neuroimage.* (2017) 163:115–24. doi: 10.1016/j.neuroimage.2017.07.059
15. Cole JH, Franke K. Predicting age using neuroimaging: innovative brain ageing biomarkers. *Trends Neurosci.* (2017) 40:681–90. doi: 10.1016/j.tins.2017.10.001
16. Dosenbach NUF, Nardos B, Cohen AL, Fair DA, Power JD, Church J, et al. Prediction of individual brain maturity using fMRI. *Science (80-.).* (2010) 329:1358–61. doi: 10.1126/science.1194144
17. Franke K, Ziegler G, Klöppel S, Gaser C. Estimating the age of healthy subjects from T1-weighted MRI scans using kernel methods: exploring the influence of various parameters. *Neuroimage.* (2010) 50:883–92. doi: 10.1016/j.neuroimage.2010.01.005
18. Levakov G, Rosenthal G, Shelef I, Raviv TR, Avidan G. From a deep learning model back to the brain—Identifying regional predictors and their relation to aging. *Hum Brain Mapp.* (2020) 41:3235–52. doi: 10.1002/hbm.25011
19. Neeb H, Zilles K, Shah NJ. Fully-automated detection of cerebral water content changes: study of age- and gender-related H2O patterns with quantitative MRI. *Neuroimage.* (2006) 29:910–22. doi: 10.1016/j.neuroimage.2005.08.062
20. Smith SM, Elliott LT, Alfaro-Almagro F, McCarthy P, Nichols TE, Douaud G, et al. Brain aging comprises many modes of structural and functional change with distinct genetic and biophysical associations. *Elife.* (2020) 9:802686. doi: 10.7554/eLife.52677
21. Cole J, Raffel J, Friede T, Eshaghi A, Brownlee W, Chard D, et al. Accelerated brain ageing and disability in multiple sclerosis. *bioRxiv.* (2019) 584888. doi: 10.1101/584888
22. Kaufmann T, van der Meer D, Doan NT, Schwarz E, Lund MJ, Agartz I, et al. Common brain disorders are associated with heritable patterns of apparent aging of the brain. *Nat Neurosci.* (2019) 22:1617–23. doi: 10.1038/s41593-019-0471-7
23. Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, et al. ImageNet large scale visual recognition challenge. *Int J Comput Vis.* (2015) 115:211–52. doi: 10.1007/s11263-015-0816-y
24. Peng H, Gong W, Beckmann CF, Vedaldi A, Smith SM. Accurate brain age prediction with lightweight deep neural networks. *Med Image Anal.* (2021) 68:101871. doi: 10.1016/j.media.2020.101871
25. Smith SM, Jenkinson M, Woolrich MW, Beckmann CF, Behrens TE, Johansen-Berg H, et al. Advances in functional and structural MR image analysis and implementation as FSL. *Neuroimage.* (2004) 23:S208–19. doi: 10.1016/j.neuroimage.2004.07.051
26. Miller KL, Alfaro-Almagro F, Bangerter NK, Thomas DL, Yacoub E, Xu J, et al. Multimodal population brain imaging in the UK Biobank prospective epidemiological study. *Nat Neurosci.* (2016) 19:1523–36. doi: 10.1038/nn.4393
27. Smith SM, Douaud G, Chen W, Hanayik T, Alfaro-Almagro F, Sharp K, et al. Enhanced Brain Imaging Genetics in UK Biobank. *bioRxiv* [Preprint]. (2020). doi: 10.1101/2020.07.27.223545
28. Alfaro-Almagro F, Jenkinson M, Bangerter NK, Andersson JLR, Griffanti L, Douaud G, et al. Image processing and Quality Control for the first 10,000 brain imaging datasets from UK Biobank. *Neuroimage.* (2018) 166:400–24. doi: 10.1016/j.neuroimage.2017.10.034
29. Kingma DP, Ba JL. *Adam: A Method for Stochastic Optimization.* *arXiv* [Preprint]. (2014) arXiv:1412.6980.
30. Krizhevsky A, Sutskever I, Hinton GE. Imagenet classification with deep convolutional neural networks. *Adv Neural Inf Process Syst.* (2012) 25:1097–105.
31. Dinsdale NK, Bluemke E, Smith SM, Arya Z, Vidaurre D, Jenkinson M, et al. Learning patterns of the ageing brain in MRI using deep convolutional networks. *Neuroimage.* (2021) 224:117401. doi: 10.1016/j.neuroimage.2020.117401
32. Smith SM, Vidaurre D, Alfaro-Almagro F, Nichols TE, Miller KL. Estimation of brain age delta from brain imaging. *Neuroimage.* (2019) 200:528–39. doi: 10.1016/j.neuroimage.2019.06.017
33. Dinsdale NK, Jenkinson M, Namburete AIL. Deep learning-based unlearning of dataset bias for MRI harmonisation and confound removal. *NeuroImage.* (2021) 228:117689. doi: 10.1016/j.neuroimage.2020.117689

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

Copyright © 2021 Gong, Beckmann, Vedaldi, Smith and Peng. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.