# Brain Structural Saliency Over The Ages

**Daniel Taylor**[1]                                                     TYLDAN008@MYUCT.AC.ZA

**Deshendran Moodley**[1]                                    DESHEN.MOODLEY@UCT.AC.ZA

**Matthias S. Treder**[2]                                            TREDERM@CARDIFF.AC.UK

**Jonathan Ipser**[1]                                             JONATHAN.IPSER@UCT.AC.ZA

**Jonathan Shock**[1,3,4]                                     JONATHAN.SHOCK@UCT.AC.ZA

[1] *University of Cape Town*

[2] *Cardiff University*

[3] *The National Institute for Theoretical and Computational Sciences, South Africa*

[4] *Institut national de la recherche scientifique, Canada*

## Abstract

Brain Age (BA) estimation via Deep Learning has become a strong and reliable bio-marker for brain health, but the black-box nature of Neural Networks does not easily allow insight into the features of brain ageing. We trained a ResNet model as a BA regressor on T1 structural MRI volumes from a small cross-sectional cohort of 524 individuals. Using Layer-wise Relevance Propagation (LRP) , we analysed the trained model to determine the most revealing structures over the course of brain ageing for the network. We show the change in attribution of relevance to different brain regions through the course of ageing. A tripartite pattern of relevance attribution to brain regions emerges. Some regions increase in relevance with age (e.g. the right Transverse Temporal Gyrus); some decrease in relevance with age (e.g. the right Fourth Ventricle); and others remained consistently relevant across ages. We also examine the effect of Brain Age Gap (BAG) on the distribution of LRP-assigned relevance within the brain volume. It is hoped that these findings will provide clinically relevant region-wise trajectories for normal brain ageing, and a baseline against which to compare brain ageing trajectories.

**Keywords:** Brain Age, Saliency Mapping, Layer-wise Relevance Propagation, MRI

## 1. Introduction

Brain Age (BA) prediction by Deep Learning (DL) methods using structural Magnetic Resonance Imaging (MRI) data has shown to be an accurate and reliable bio-marker for brain health (Cole et al., 2017). The Brain Age Gap (BAG) of an individual – the difference between the predicted BA and the chronological age – is an increasingly popular and predictive bio-marker as well. BAG predicts accelerated or slowed brain ageing (Smith et al., 2019; Dinsdale et al., 2021; Peng et al., 2021).

DL frameworks have shown great efficacy in the BA regression task (Cole et al., 2018; Jonsson et al., 2019; Kossaifi et al., 2020; Peng et al., 2021). An advantage of DL models is that they are able to analyse whole-brain structural images in their raw state. Other summary statistics such as measures of cortical thickness, volume and surface area (Zhao et al., 2019) treat the brain as being modular, and may therefore lack details contained in raw

MRI volumes (which represent the entire structure with all its integrated substructures). Convolutional Neural Networks (CNNs) have been shown to give extremely small Mean Absolute Error (MAE) between true age and BA. The state-of-the-art in BA regression at the time of writing is around 2.5y MAE (Kossaifi et al., 2020; Dinsdale et al., 2021; Peng et al., 2021). Such accuracies have thus far only been achieved on large datasets, and some of these datasets are limited in their age ranges, such as the UK Biobank (Sudlow et al., 2015) (45y-80y).

The black-box nature of Neural Networks makes it difficult to attribute BAG to specific brain regions. Many methods of Explainable Artificial Intelligence (XAI) have aimed at alleviating the black-box issue. These attempt to explain why the models make their predictions. A popular post-hoc method of reasoning for DL models is saliency mapping. This is a group of methods by which areas of input are assigned relevance scores proportional to their saliency to the model decision. We explore Layer-wise Relevance Propagation (LRP) (Bach et al., 2015) as a saliency mapping method for the BA task. In this work, we use the terms 'salience' and 'relevance' interchangeably.

Many factors can contribute to accelerated BA (positive BAG), such as the presence of neurodegenerative disease, type 2 diabetes or HIV, and even past physical activity (Cole et al., 2017). If we wish to create a model that accurately predicts BA, and therefore accurately assesses BAG, we must ensure that the model accurately captures a path of 'normal' brain ageing. It must thus be trained on data in which the subjects do not have any underlying pathologies that can affect BA independent of actual age (to the extent that this is possible). Even in a 'healthy' cohort, however, there will be non-negligible deviation from an average ageing trajectory. Smith et al. (2019) note that some meaningful BAG should exist between the predicted and chronological age of an individual, as long as the model does not badly over-fit.

In this work we aim to shed light on the contributions of specific brain regions to BA and BAG through the course of ageing. We explore the trajectories of relevance assignment across ages using LRP. It is hoped that this work will provide a clinically relevant baseline of salience trajectories for brain regions with respect to BA.

For our findings in the current literature, we expected that regions that are deemed highly relevant in general towards BA would increase in their relevance with age. Areas that are generally less salient would then necessarily decrease further in assigned relevance.

## 2. Related Work

Levakov et al. (2020) used SmoothGRAD (Smilkov et al., 2017) to create population-level saliency maps for a BA regression model. The authors examined which brain regions corresponded with the highest attribution of saliency by thresholding the top-1% of voxels in saliency maps and examining clusters with $> 100$ voxels. Most salience was attributed to the cisterns and ventricles.

Hofmann et al. (2021) used LRP to create saliency maps for two multi-ensemble BA regression models. They note that the SmoothGRAD method employed by Levakov et al. (2020) is not directional, while LRP highlights areas of input that both agree with and contradict the output. They utilised the $\text{LRP}_{\text{CMP}}$ saliency method with $\alpha = 1$ (see Appendix (C)). The authors verified the method on a simulated brain ageing model and found that

in the regression task, regions of higher relevance argue for greater age, while regions of lower relevance argue for lower age. After verifying the method, the authors created the LRP saliency maps for the real BA regression task. Like Levakov et al. (2020) they found that relevance was greatest in and around the ventricles. They also found a significant attribution to the grey-matter-dense regions at the cortical surface. The authors also contrasted the saliency maps of individuals in a younger and an older cohort to determine the difference in attribution of relevance.

While Hofmann et al. (2021) compared a younger and older cohort in their study, no previous work has focused on the change in relevance attributions to brain regions as a function of subject age. The authors also only examined statistical significance toward BAG in an older cohort, and not age-related contributors to BAG. Furthermore, the concern raised by Geirhos et al. (2019) and Sixt et al. (2019) about modified backpropagation algorithms like LRP has not been addressed. Their concern is that such methods attempt to recreate the input to the model, as opposed to focusing solely on areas of relevance.

## 3. Methods

### 3.1. BA Regression Model

The model used in the BA regression task was a ResNet (He et al., 2016) with filter number sequence $[32, 32, 64, 64, 128]$ in the main branches of the 5 residual blocks. This configuration of residual blocks is borrowed from Jonsson et al. (2019). We used a Softplus activation function for all the nonlinearities in the model, as recommended by Dombrowski et al. (2019), for the robustness of the saliency maps.

We used the cross-sectional Cam-CAN T1-weighted MRI dataset (Shafto et al., 2014; Taylor et al., 2017) for our experiments. The volumes are from 656 healthy individuals ranged 18-89y, with 49.4% or 324 male participants.

We trained our ResNet model on a random 80% training split of the dataset (524 subjects). 80 of these were randomly held out for validation (15% validation split from the beginning of training). The remaining 20% of the dataset (132 subjects) was used in testing. We used hyper-parameter tuning to find the best optimiser for the model, as well as the best loss function and the best starting learning rate. Taking the average of three trials for each hyperparameter configuration, it was found that the best configuration was using the Mean Squared Error (MSE) loss function with the Adam optimiser (Kingma and Ba, 2014) and an initial learning rate of $5 \times 10^{-3}$. For a set of $n$ predictions $[f(x_i)]$ and the corresponding set of $n$ true ages $[y_i]$, the MSE is defined as $\text{MSE}\left([f(x_i)], [y_i]\right) = \frac{\sum_i (f(x_i) - y_i)^2}{n}$. The learning rate was halved every 20 epochs using a learning rate schedule. We also utilised early stopping in training such that if there was no validation improvement for 20 epochs, the training would cease. Due to the large size of the data ($233 \times 189 \times 197$ voxels per scan), the batch size was limited to 4. The metrics used to evaluate the model were the MAE, MSE and Mean Average Percentage Error (MAPE). Appendix (A) details the pre-processing of the data.

3

### 3.2. Saliency Mapping

We performed saliency mapping to analyse the trained model using three LRP and two DeepLIFT (Shrikumar et al., 2017) methods. Details of these can be found in Appendix (C). We used the *Epsilon Rule* for dense layers and the $z^{\mathcal{B}}$-*Rule* for the input layer (Montavon et al., 2017).

After comparing the region-wise relevance attributions of the saliency mapping methods with previous findings (Levakov et al., 2020; Hofmann et al., 2021), we chose one method, LRP with $\alpha = 1$, for subsequent analyses.

We used permutation-based t-tests to determine statistically significant differences of the saliency maps from the input volumes. This is in order to address the concern of Geirhos et al. (2019) that the LRP methods simply recreate the input.

#### 3.2.1. REGIONAL RELEVANCE ATTRIBUTION

To analyse the distribution of relevance in the brain, we examined the attribution of top-1% relevance (T1R) to different brain regions across the test set. T1R refers to the the voxels with the first percentile of activations in the brain volume for a saliency map. This is similar to the methodology of Levakov et al. (2020). We count the number of such voxels in each region, and normalise by the region volume to get the proportion of the region assigned T1R. From the findings of Hofmann et al. (2021), we know that regions of high saliency argue for higher age and those of lower saliency argue for lower age in the context of BA regression, and so the directionality of relevance is important. We are most interested in regions which contribute to accelerated brain ageing, and so we consider only positive relevance in T1R (see the discussion of directionality by Hofmann et al. (2021)).

#### 3.2.2. REGIONAL RELEVANCE ATTRIBUTION BASED ON BAG

To examine the effect of large BAG on the distribution of relevance in the brain, we thresholded BAG in the test set and compared the T1R distributions of those lying above and below the threshold. This has not been done before in the literature, and so we made the simple choice of threshold for BAG of $\delta^* = \sigma$, where $\sigma$ is the standard deviation of the BAG values for the test set. We compared the distribution of T1R within the brains of those lying above the threshold ($\delta \geq \delta^*$) to those lying within the threshold ($|\delta| < \delta^*$). To calculate BAG, we used the age-orthogonal value $\delta_2$ of Smith et al. (2019). We compare the T1R distributions of three groups for LRP: $\delta_2 \geq \delta^*$ for younger individuals ($< 50y$), $|\delta| < \delta^*$, and $\delta_2 \geq \delta^*$ for older individuals ($> 50y$).

#### 3.2.3. REGIONAL RELEVANCE ATTRIBUTION ACROSS AGE BRACKETS

To assess the change in the distribution of relevance as a function of age, we examine how the T1R assignment by LRP changes between equally spaced age brackets. We do this for each region individually, by calculating the average proportion of T1R assigned to the region across the test set members who lie within each age bracket.

Of greatest interest are highly-relevant regions whose relevance changes most or least across ages. To quantify the greatest change, we examine those with the highest standard

deviation (SD) in proportion assigned T1R. To quantify the least change, we examine those with lowest coefficient of variation (CoV – standard deviation normalised by mean).

We do not focus on those regions with lowest SD since those tend to be the regions with lowest mean relevance attribution. Similarly we do not focus on those with highest CoV, since these also tend to be regions with low mean attributions, where a small change in relevance attribution can lead to a very large CoV. Our aim is to focus on the changes in relevance for highly-relevant structures. More is given on this choice of metrics in Appendix (D).

## 4. Results

### 4.1. BA Regression Model

The training was stopped by the early stopping callback at the end of the $143^{\text{rd}}$ epoch. The performance on the evaluation metrics is given in Table (2), and the regression plot is shown in Figure (3); both in Appendix (B).

### 4.2. Saliency Mapping

We report here the findings of the saliency mapping tasks. Figure (4a) in Appendix (C) shows sections of a composite MRI volume overlaid with aggregate saliency maps for LRP. The composite volume is created as the aggregate of all the brain volumes from the test set, and the overlaid map is similarly the aggregate of all the saliency maps for the test set. The aggregate map is thresholded to the top-1% of voxel values to show the average distribution of T1R over the test set. We also show in Figure (4b) an example of overlaid LRP salience on sections of the volume of a 60-year-old male subject.

#### 4.2.1. Regional Relevance Attribution

Here we show the results of the assignment of relevance to brain regions by LRP. We compare the results to previous findings as well as to some of the medical BA literature.

Figure (6) in Appendix (C) shows the proportion of each region that is assigned T1R. The relevance assignment in this case is averaged between the two hemispheres, since there is a large degree of symmetry in the distribution between hemispheres (see Appendix (C)). The ordering of the regions along the x-axis is in descending order of T1R proportion.

It is clear that LRP assigns high T1R to the ventricles (particularly the lateral and fourth), which agrees not only with the medical literature (Anderton, 2002; Raz and Rodrigue, 2006), but also with the findings of Levakov et al. (2020) and Hofmann et al. (2021). The transverse temporal gyrus was assigned T1R in large proportion, as were many limbic structures, such as the caudate nucleus, hippocampus, thalamus, parahippocampal gyrus and diencephalon. Relevance of the limbic system to ageing has been shown by Gunbey et al. (2014). A high proportion of T1R tended to be assigned to grey-matter-dense areas such as the vermal lobules of the cerebellum. This agrees with findings that grey matter density decreases with age, starting from late adolescence (Peters, 2006; Sowell et al., 2003; Raz and Rodrigue, 2006). On the other hand, the cerebellum white matter tends not to be assigned much relevance. Although white matter lesions can be indicative of age (Peters, 2006; Anderton, 2002; Sowell et al., 2003; Raz and Rodrigue, 2006), these are not shown

5

distinctively in T1-weighted images, and so we would not expect white matter regions to be assigned high relevance by this model. The optic chiasm tends to be assigned significant relevance. This may be due to its small size in conjunction with its proximity to other highly relevant structures such as the interpeduncular cistern (Levakov et al., 2020). The large amount of relevance assigned to the temporal gyrus agrees with findings by Lutz et al. (2007) of involvement with brain ageing, and this relevance may be due to the degradation with age of the auditory cortex.

The permutation-based t-tests showed that there were statistically significant ($p <$ 0.001) differences between the saliency maps and the input volumes for each method in almost the entire brain volume.

### 4.2.2. Regional Relevance Attribution Based on BAG

Here we compare the distribution of relevance in individuals with small BAG to those with large BAG, both older and younger. We examine how regional relevance associates with BAG for older and younger individuals. Figure (7) in Appendix (D) shows the distributions of the quantities $\delta_1$ and $\delta_2$, as defined by Smith et al. (2019), for the test set. The correction from orthogonalisation to age shifts the distribution to become approximately Gaussian. This allows us to threshold BAG reliably to compare BAG-salient regions in younger and older groups. On the test set, our threshold value is $\delta^* = \sigma = 11.58$.
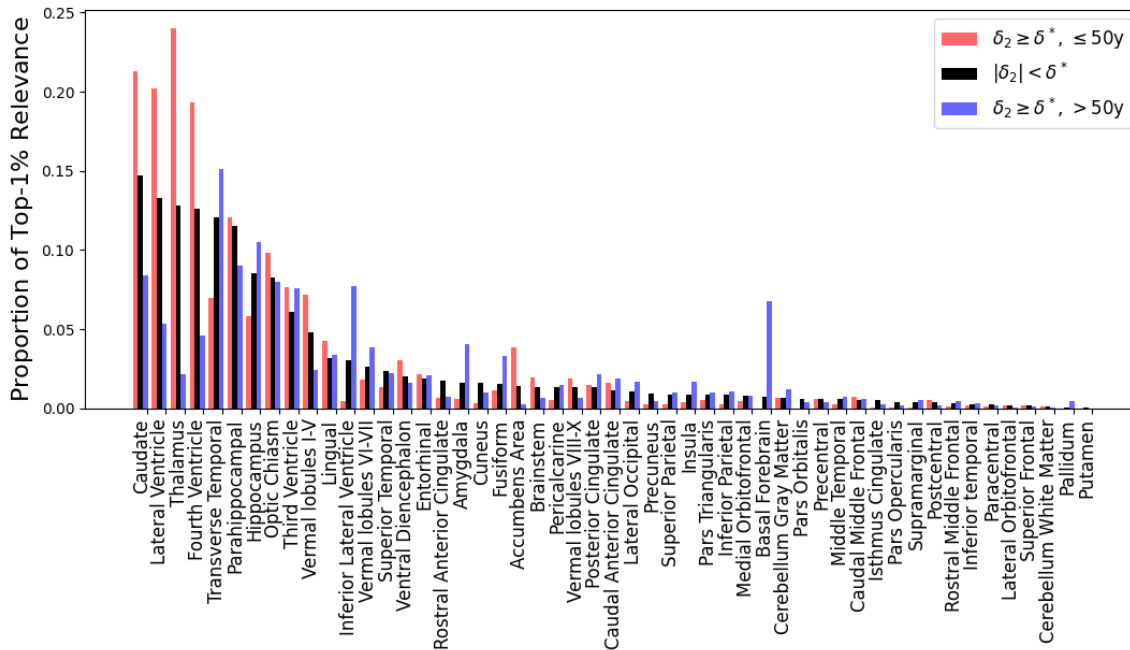


Figure 1: Effect of large BAG in young and old individuals, for LRP

Figure (1) compares the distributions of T1R for these groups. We see that several regions display significant changes in T1R assignment according to grouping by BAG and

age. By way of example from the figure, for large BAG values, LRP assigns to the transverse temporal gyrus less relevance in younger individuals than in older individuals. The opposite is true for the thalamus. Since increased relevance is predictive of higher BAG, this analysis suggests that markers of BAG change through the course of ageing. It appears that regions can be more indicative of BAG at some ages than others.

### 4.2.3. REGIONAL RELEVANCE ATTRIBUTION ACROSS AGE BRACKETS

Here we examine the change in relevance across age brackets for brain regions according to LRP. In Figure (8) of Appendix (E) we show the standard deviations and coefficients of variation in the proportion of regions assigned T1R over all age brackets.

One region which has a large standard deviation across ages is the transverse temporal gyrus. In Figure (2) we show the T1R across age brackets, normalised to have a maximum of 1, of the right transverse temporal gyrus. The figure shows a dramatic increase in relevance of the region with age.
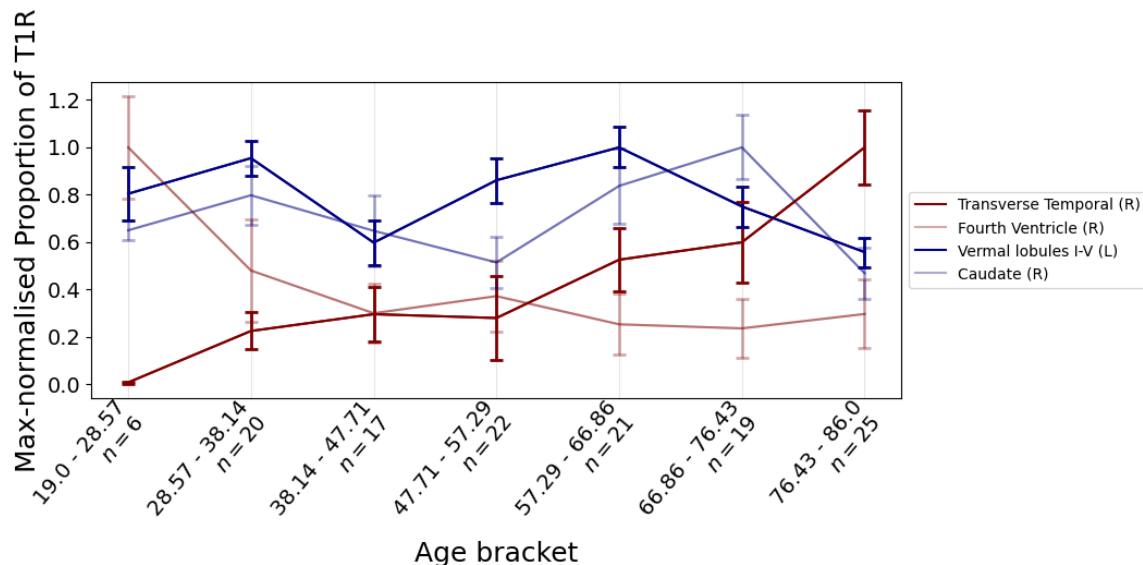


Figure 2: Proportion of T1R over age brackets for LRP, normalised to have max value 1. error bars show the standard deviations of relevance assignment within each age bracket.

Greater attribution of relevance in older subjects than younger subjects suggests that the structure bears more information about BA and BAG in older individuals. This agrees with our findings about relevance distribution as a function of BAG. The transverse temporal gyrus for example is a more salient indicator of accelerated brain ageing in older subjects than in younger subjects. On the other hand, we see in the figure the T1R across age brackets for the right fourth ventricle. This region too is shown to have significant standard deviation across age brackets. There is a marked decline in assignment of T1R from younger

to older ages. This suggests that while the fourth ventricle is universally relevant to BA and BAG, it is most telling of ageing in younger individuals.

We also see in the figure the T1R across age brackets for the left Vermal Lobules I-V. This region was shown to have uniformly low CoV across the methods. We see that there is a relatively uniform assignment of T1R across the age brackets in this case. Finally in the figure we see the T1R over age brackets for the right caudate nucleus. The region has a low CoV for LRP. Although there is some upward fluctuation in T1R in some of the older age brackets, there is no overarching upward or downward trend.

## 5. Discussion and Conclusions

We found that several regions were expressed as highly salient by LRP. As a group the ventricles were consistently assigned high proportions of T1R. Many limbic structures such as the caudate nucleus and the thalamus also were assigned high proportions of T1R.

We found that large BAG values are associated with distributions in relevance among brain regions that can vary to a large degree depending on age. This suggests that the degree of saliency of a region towards BAG is age-dependent. Upon inspection of age-bracket distribution of relevance, we showed that a tripartite pattern emerges of relevance attribution over age brackets to brain regions. Some regions increase in relevance with age, some decrease with age and others retain relatively uniform relevance with age. From our findings with BAG, decreases with age would suggest that the region is more informative of BAG in younger individuals than in older individuals (for example, the fourth ventricle). Similarly, increases with age would suggest that the region is more informative of BAG in older individuals than in younger ones (for example, the transverse temporal gyrus). Uniformity across ages in relevance would suggest that the region is consistently informative of BAG (for example, the caudate nucleus).

This tripartite pattern of trajectories was unexpected, since previous studies have only reported on increased relevance with age (Hofmann et al., 2021). Indeed, many regions decrease in the proportion assigned T1R with age. This can be explained for a structure like the right fourth ventricle, by the fact that although the structure is informative of BA at all ages, it is expected to dilate with age. A younger individual with dilated ventricles is clearly experiencing pathological brain ageing.

These findings may be used as baseline regional trajectories for BA salience in a clinical setting. As part of a DL pipeline for BA regression, clinicians can create saliency maps for patients' BA predictions and compare the regional distributions of T1R to these baselines. Large individual BAG can be assessed regionally by these methods. Saliency mapping can be performed post-hoc on the regression model in close to real-time, so as to have BA predictions accompanied by an explanation in a clinical setting.

Without access to this model and its accompanying results, clinicians will still be able to use the underlying methods – perhaps with more powerful models trained on larger datasets – to create and compare relevance trajectories for clinical use.

This work would benefit from development on a larger dataset and with a more powerful model, such as the ensemble model of Levakov et al. (2020). It would also be of interest in future to compare the saliency maps of several other, possibly more computationally expensive techniques, such as Integrated Gradients (Sundararajan et al., 2016, 2017).

## Acknowledgments

## References

Brian H. Anderton. Ageing of the brain. *Mechanisms of Ageing and Development*, 2002. ISSN 00476374. doi: 10.1016/S0047-6374(01)00426-2.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE*, 10(7):1–46, 2015. ISSN 19326203. doi: 10.1371/journal.pone.0130140.

J. H. Cole, S. J. Ritchie, M. E. Bastin, M. C. Valdés Hernández, S. Muñoz Maniega, N. Royle, J. Corley, A. Pattie, S. E. Harris, Q. Zhang, N. R. Wray, P. Redmond, R. E. Marioni, J. M. Starr, S. R. Cox, J. M. Wardlaw, D. J. Sharp, and I. J. Deary. Brain age predicts mortality. *Molecular Psychiatry*, 2018. ISSN 14765578. doi: 10.1038/mp.2017.62.

James H. Cole, Rudra P.K. Poudel, Dimosthenis Tsagkrasoulis, Matthan W.A. Caan, Claire Steves, Tim D. Spector, and Giovanni Montana. Predicting brain age with deep learning from raw imaging data results in a reliable and heritable biomarker. *NeuroImage*, 163: 115–124, 2017. ISSN 10959572. doi: 10.1016/j.neuroimage.2017.07.059. URL http://dx.doi.org/10.1016/j.neuroimage.2017.07.059.

DL Collins, AP Zijdenbos, WFC Baare, and AC Evans. ANIMAL+INSECT: Improved Cortical Structure Segmentation. *IPMI Lecture Notes in Computer Science*, 1613/1999: 210–223, 1999. doi: http://dx.doi.org/10.1007/3-540-48714-X_16.

Nicola K. Dinsdale, Emma Bluemke, Stephen M. Smith, Zobair Arya, Diego Vidaurre, Mark Jenkinson, and Ana I.L. Namburete. Learning patterns of the ageing brain in MRI using deep convolutional networks. *NeuroImage*, 224:117401, jan 2021. ISSN 1053-8119. doi: 10.1016/J.NEUROIMAGE.2020.117401.

Ann Kathrin Dombrowski, Maximilian Alber, Christopher J. Anders, Marcel Ackermann, Klaus Robert Müller, and Pan Kessel. Explanations can be manipulated and geometry is to blame. In *Advances in Neural Information Processing Systems*, 2019.

Fabian Eitel, Emily Soehler, Judith Bellmann-Strobl, Alexander U. Brandt, Klemens Ruprecht, René M. Giess, Joseph Kuchling, Susanna Asseyer, Martin Weygandt,

John Dylan Haynes, Michael Scheel, Friedemann Paul, and Kerstin Ritter. Uncovering convolutional neural network decisions for diagnosing multiple sclerosis on conventional MRI using layer-wise relevance propagation. *NeuroImage: Clinical*, 24, 2019. ISSN 22131582. doi: 10.1016/j.nicl.2019.102003.

VS Fonov, AC Evans, CR Almli, RC McKinstry, and DL Collins. Unbiased nonlinear average age-appropriate brain templates from birth to adulthood. *NeuroImage*, 47:s102, jul 2009. doi: http://dx.doi.org/10.1016/S1053-8119(09)70884-5.

Robert Geirhos, Claudio Michaelis, Felix A. Wichmann, Patricia Rubisch, Matthias Bethge, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. *7th International Conference on Learning Representations, ICLR 2019*, (c):1–22, 2019.

Irina Grigorescu, Lucilio Cordero-Grande, A David Edwards, Jo Hajnal, Marc Modat, and Maria Deprez. Interpretable Convolutional Neural Networks for Preterm Birth Classification. pages 1–4, 2019. URL http://arxiv.org/abs/1910.00071.

Hediye Pjnar Gunbey, Karabekir Ercan, Ayge Serap Fjndjkoglu, H Taner Bulut, Mustafa Karaoglanoglu, and Halil Arslan. The Limbic Degradation of Aging Brain: A Quantitative Analysis with Diffusion Tensor Imaging. 2014. doi: 10.1155/2014/196513. URL http://dx.doi.org/10.1155/2014/196513.

Sukrit Gupta, Yi Hao Chan, and Jagath C. Rajapakse. Decoding Brain Functional Connectivity Implicated in AD and MCI. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2019. ISBN 9783030322472. doi: 10.1007/978-3-030-32248-9_87.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016. ISBN 9781467388504. doi: 10.1109/CVPR.2016.90.

Simon M Hofmann, Frauke Beyer, Sebastian Lapuschkin, Markus Loeffler, Klaus-Robert Müller, Arno Villringer, Wojciech Samek, and A Veronica Witte. Towards the Interpretability of Deep Learning Models for Human Neuroimaging. *bioRxiv*, page 2021.06.25.449906, jan 2021. doi: 10.1101/2021.06.25.449906. URL http://biorxiv.org/content/early/2021/08/26/2021.06.25.449906.abstract.

M Jenkinson and SM Smith. A global optimisation method for robust affine registration of brain images. *Medical Image Analysis*, 5(2):143–156, 2001.

M Jenkinson, PR Bannister, JM Brady, and SM Smith. Improved optimisation for the robust and accurate linear registration and motion correction of brain images. *NeuroImage*, 17(2):825–841, 2002.

B A Jonsson, G Bjornsdottir, T E Thorgeirsson, L M Ellingsen, G Bragi Walters, D F Gudbjartsson, H. Stefansson, K Stefansson, and M O Ulfarsson. Brain age prediction using deep learning uncovers associated sequence variants. *Nature Communications*, 10

(1), 2019. ISSN 20411723. doi: 10.1038/s41467-019-13163-9. URL https://doi.org/10.1038/s41467-019-13163-9.

Diederik P. Kingma and Jimmy Lei Ba. Adam: A Method for Stochastic Optimization. *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*, dec 2014. URL https://arxiv.org/abs/1412.6980v9.

Maximilian Kohlbrenner, Alexander Bauer, Shinichi Nakajima, Alexander Binder, Wojciech Samek, and Sebastian Lapuschkin. Towards best practice in explaining neural network decisions with LRP. pages 1–5, 2019. URL http://arxiv.org/abs/1910.09840.

Jean Kossaifi, Arinbjörn Kolbeinsson, Aran Khanna, Tommaso Furlanello, and Anima Anandkumar. Tensor Regression Networks. *Journal of Machine Learning Research*, 21: 1–21, 2020. URL http://jmlr.org/papers/v21/18-503.html.

Gidon Levakov, Gideon Rosenthal, Ilan Shelef, Tammy Riklin Raviv, and Galia Avidan. From a deep learning model back to the brain—Identifying regional predictors and their relation to aging. *Human Brain Mapping*, 41(12):3235–3252, aug 2020. ISSN 10970193. doi: 10.1002/HBM.25011/FORMAT/PDF.

Juergen Lutz, Felix Hemminger, Robert Stahl, Olaf Dietrich, Martin Hempel, Maximilian Reiser, and Lorenz Jäger. Evidence of Subcortical and Cortical Aging of the Acoustic Pathway: A Diffusion Tensor Imaging (DTI) Study. *Academic Radiology*, 14(6):692–700, jun 2007. ISSN 1076-6332. doi: 10.1016/J.ACRA.2007.02.014.

AL Manera, M Dadar, VS Fonov, and DL Collins. CerebrA, registration and manual label correction of Mindboggle-101 atlas for MNI-ICBM152 template. *Scientific Data*, (7(1)): 1–9, 2020.

Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus Robert Müller. Explaining nonlinear classification decisions with deep Taylor decomposition. *Pattern Recognition*, 65(May 2016):211–222, 2017. ISSN 00313203. doi: 10.1016/j.patcog.2016.11.008. URL http://dx.doi.org/10.1016/j.patcog.2016.11.008.

Han Peng, Weikang Gong, Christian F. Beckmann, Andrea Vedaldi, and Stephen M. Smith. Accurate brain age prediction with lightweight deep neural networks. *Medical Image Analysis*, 68:101871, feb 2021. ISSN 1361-8415. doi: 10.1016/J.MEDIA.2020.101871.

Ruth Peters. Ageing and the brain, 2006. ISSN 00325473.

Theerasarn Pianpanit, Sermkiat Lolak, Phattarapong Sawangjai, Apiwat Ditthapron, Sanparith Marukatat, Ekapol Chuangsuwanich, and Theerawit Wilaiprasitporn. Interpreting deep learning prediction of the Parkinson's disease diagnosis from SPECT imaging. aug 2019. URL https://arxiv.org/abs/1908.11199http://arxiv.org/abs/1908.11199.

Naftali Raz and Karen M. Rodrigue. Differential aging of the brain: Patterns, cognitive correlates and modifiers, 2006. ISSN 01497634.

Meredith A. Shafto, Lorraine K. Tyler, Marie Dixon, Jason R. Taylor, James B. Rowe, Rhodri Cusack, Andrew J. Calder, William D. Marslen-Wilson, John Duncan, Tim Dalgleish, Richard N. Henson, Carol Brayne, Ed Bullmore, Karen Campbell, Teresa Cheung, Simon Davis, Linda Geerligs, Rogier Kievit, Anna McCarrey, Darren Price, David Samu, Matthias Treder, Kamen Tsvetanov, Nitin Williams, Lauren Bates, Tina Emery, Sharon Erzinçlioglu, Andrew Gadie, Sofia Gerbase, Stanimira Georgieva, Claire Hanley, Beth Parkin, David Troy, Jodie Allen, Gillian Amery, Liana Amunts, Anne Barcroft, Amanda Castle, Cheryl Dias, Jonathan Dowrick, Melissa Fair, Hayley Fisher, Anna Goulding, Adarsh Grewal, Geoff Hale, Andrew Hilton, Frances Johnson, Patricia Johnston, Thea Kavanagh-Williamson, Magdalena Kwasniewska, Alison McMinn, Kim Norman, Jessica Penrose, Fiona Roby, Diane Rowland, John Sargeant, Maggie Squire, Beth Stevens, Aldabra Stoddart, Cheryl Stone, Tracy Thompson, Ozlem Yazlik, Dan Barnes, Jaya Hillman, Joanne Mitchell, Laura Villis, and Fiona E. Matthews. The Cambridge Centre for Ageing and Neuroscience (Cam-CAN) study protocol: A cross-sectional, lifespan, multidisciplinary examination of healthy cognitive ageing. *BMC Neurology*, 2014. ISSN 14712377. doi: 10.1186/s12883-014-0204-1.

Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. *34th International Conference on Machine Learning, ICML 2017*, 7:4844–4866, 2017.

Leon Sixt, Maximilian Granz, and Tim Landgraf. When Explanations Lie: Why Many Modified BP Attributions Fail. (December 2019), 2019. URL http://arxiv.org/abs/1912.09818.

Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, Martin Wattenberg, Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. SmoothGrad: removing noise by adding noise. *arXiv*, page arXiv:1706.03825, jun 2017. ISSN 2331-8422. URL https://ui.adsabs.harvard.edu/abs/2017arXiv170603825S/abstract.

SM Smith. Fast robust automated brain extraction. *Human Brain Mapping*, 17(3):143–155, nov 2002.

SM Smith, M Jenkinson, MW Woolrich, CF Beckmann, TEJ Behrens, H Johansen-Berg, PR Bannister, M De Luca, I Drobnjak, DE Flitney, R Niazy, J Saunders, J Vickers, Y Zhang, N De Stefano, JM Brady, and PM Matthews. Advances in functional and structural MR image analysis and implementation as FSL. *NeuroImage*, 23(S1):208–19, 2004.

Stephen M. Smith, D. Vidaurre, F. Alfaro-Almagro, Thomas E. Nichols, and Karla L. Miller. Estimation of brain age delta from brain imaging. *NeuroImage*, 200:528–539, oct 2019. ISSN 1053-8119. doi: 10.1016/J.NEUROIMAGE.2019.06.017.

Elizabeth R. Sowell, Bradley S. Peterson, Paul M. Thompson, Suzanne E. Welcome, Amy L. Henkenius, and Arthur W. Toga. Mapping cortical change across the human life span. *Nature Neuroscience*, 2003. ISSN 10976256. doi: 10.1038/nn1008.

Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott, Jane Green, Martin Landray, Bette Liu, Paul Matthews, Giok Ong, Jill Pell, Alan Silman, Alan Young, Tim Sprosen, Tim Peakman, and Rory Collins. UK Biobank: An Open Access Resource for Identifying the Causes of a Wide Range of Complex Diseases of Middle and Old Age. *PLOS Medicine*, 12(3):e1001779, mar 2015. ISSN 1549-1676. doi: 10.1371/JOURNAL.PMED.1001779. URL https://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.1001779.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Gradients of Counterfactuals. nov 2016. URL https://arxiv.org/abs/1611.02639http://arxiv.org/abs/1611.02639.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. *34th International Conference on Machine Learning, ICML 2017*, 7:5109–5118, 2017.

Jason R. Taylor, Nitin Williams, Rhodri Cusack, Tibor Auer, Meredith A. Shafto, Marie Dixon, Lorraine K. Tyler, Cam-CAN, and Richard N. Henson. The Cambridge Centre for Ageing and Neuroscience (Cam-CAN) data repository: Structural and functional MRI, MEG, and cognitive data from a cross-sectional adult lifespan sample. *NeuroImage*, 144: 262–269, jan 2017. ISSN 1053-8119. doi: 10.1016/J.NEUROIMAGE.2015.09.018.

AM Winkler, GR Ridgway, MA Webster, SM Smith, and TE Nichols. Permutation inference for the general linear model. *NeuroImage*, 92:381–397, 2014.

Lu Zhao, William Matloff, Kaida Ning, Hosung Kim, Ivo D Dinov, and Arthur W Toga. Age-Related Differences in Brain Morphology and the Modifiers in Middle-Aged and Older Adults. 2019. doi: 10.1093/cercor/bhy300. URL http://biobank.ctsu.ox.

## Appendix A. Pre-processing

The standard pre-processing of the Cam-CAN dataset is detailed by Taylor et al. (2017). We used FSL's Brain Extraction Tool (BET) (Smith, 2002) to perform skull-stripping on the T1-weighted volumes, using the recommended fractional intensity of 0.5. We then aligned all volumes to MNI space by registering all to a standard MNI volume (Manera et al., 2020), using the FLIRT tool (Jenkinson and Smith, 2001; Jenkinson et al., 2002). This was done primarily for spatial correspondence with the region atlas (Collins et al., 1999; Fonov et al., 2009). Voxel values were normalised using Global Contrast Normalisation (GCN) to have unit variance within each volume.

## Appendix B. Regression Model

Figure (3) shows the regression plot for our ResNet model on the test set, and Table (2) shows the performance metrics. We note that the MAE of 6.5y is far from state-of-the-art, but also that the size of the dataset is far less than those of previous works. Table (1) shows the performances of some key studies for BA regression, with the number of individuals in the datasets used.

| Author | Ours | Cole et al. (2017) | Jónsson et al. (2019) | Kossaifi et al. (2020) | Levakov et al. (2020) | Hofmann et al. (2021) | Peng et al. (2021) |
|---|---|---|---|---|---|---|---|
| **Dataset Size** | 656 | 2001 | 12378 | 19100 | 10176 | 2016 | 14503 |
| **Test MAE (y) on T1 Volumes** | 6.55 | 4.65 | 4.00 | 2.69 | 3.02 | 3.95 | 2.14 |

Table 1: Test MAE performances of several key BA regression studies, as compared to ours. Our dataset is at least three times smaller than any other used here.



| Metric | Result |
|---|---|
| MAE | 6.5457 |
| MSE | 72.5481 |
| MAPE | 13.5305 |

Table 2: Metrics

Figure 3: Test set regression plot

Our ResNet had filter sizes $[32, 32, 64, 64, 128]$ in the main branches of the five residual modules. This configuration was adapted from the lightweight model of Jonsson et al. (2019).

## Appendix C. Saliency Maps

LRP and DeepLIFT are relevance decomposition methods which assign scores to a model's input space according to saliency to the model's decision. Both methods take into account the learned parameters of the model and the activations at each layer of the model specific to the inference of a given input. Relevance decomposition refers to the backward propagation of relevance from the output layer of a model to its input. We chose LRP and DeepLIFT because these are both fast and computationally inexpensive saliency mapping methods. Both methods have also been used before in reasoning for brain imaging (Eitel et al., 2019; Grigorescu et al., 2019; Hofmann et al., 2021; Pianpanit et al., 2019; Gupta et al., 2019), and offer directionality in their explanations, unlike simpler method such as SmoothGRAD (Smilkov et al., 2017)).

The decomposition rules for $\text{LRP}_{\text{CMP}}$ (Kohlbrenner et al., 2019) from a layer $(L+1)$ to a layer $L$ are as follows:

For fully-connected layers, the $\varepsilon$-rule (Bach et al., 2015) is used:

$$R_i^{(L)} = \sum_j \frac{z_{ij}}{\sum_{i'} z_{i'j} + \varepsilon} R_j^{(L+1)} \tag{1}$$

where $R_i^{(L)}$ is the relevance assigned to neuron $i$ in layer $L$, $z_{ij} = x_i w_{ij}$ is the product of previous layer activations and the learned weights between the layers, and $\varepsilon$ is a small constant for numerical stability.

For convolutional, BatchNorm and pooling layers (apart from the input layer), the $\text{LRP}_{\alpha\beta}$ rule (Bach et al., 2015) is used:

$$R_i^{(L)} = \sum_j \left( \alpha \frac{(x_i w_{ij})^+}{\sum_{i'} (x_i w_{ij})^+ + b_j^+} - \beta \frac{(x_i w_{ij})^-}{\sum_{i'} (x_i w_{ij})^- + b_j^-} \right) R_j^{(L+1)} \tag{2}$$

with the constraint that $\alpha - \beta = 1$. $w_{ij}^+$ refers to the positive parts only of the weights and $b_j^+$ to the positive parts only of the biases.

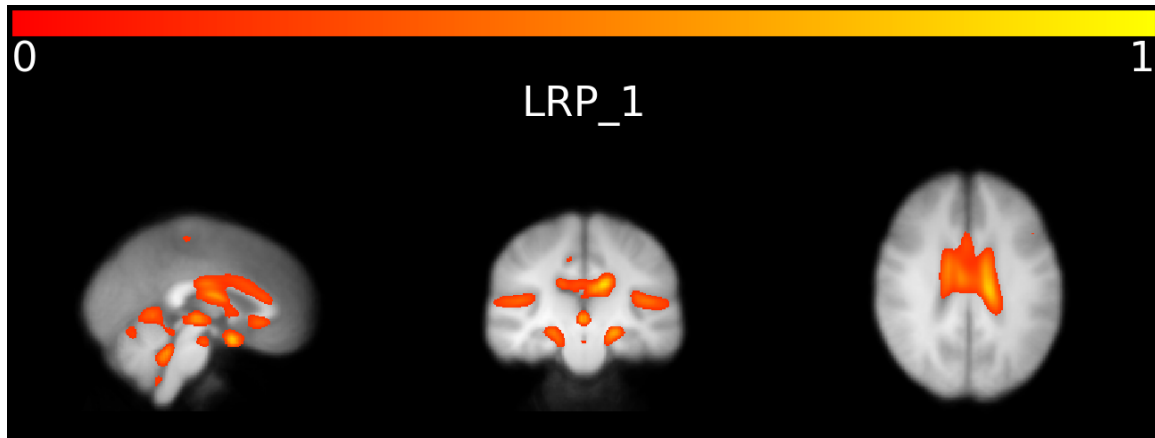For the input layer, the $z^{\mathcal{B}}$-rule (Montavon et al., 2017) is used:

$$R_i^{(L)} = \sum_j \frac{z_{ij} - l_i w_{ij}^+ - h_i w_{ij}^-}{\sum_{i'} z_{i'j} - l_{i'} w_{i'j}^+ - h_{i'} w_{i'j}^-} R_j^{(L+1)} \tag{3}$$

where $h_i$ and $l_i$ are the highest and lowest input values respectively.
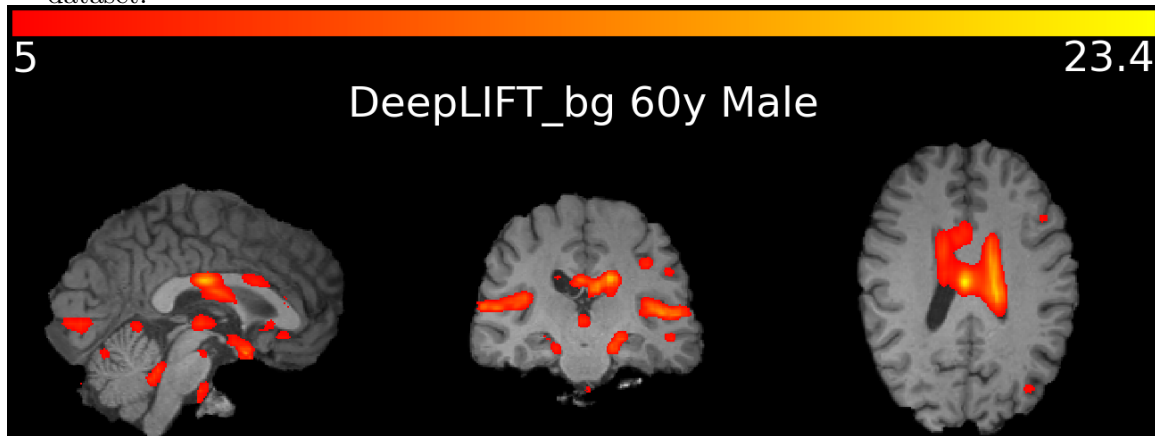
The DeepLIFT rules for relevance attribution can be found in (Shrikumar et al., 2017).

The parameters $\alpha$ and $\beta$ of LRP act as contrast parameters, with higher values of $\alpha$ giving higher contrast in the saliency maps (Bach et al., 2015).

Figure (5) shows the distributions of the correlation coefficients for relevance attribution to the same region in opposite hemispheres via LRP. We see that the majority of regions have

(*a*) Aggregate T1R map for LRP overlaid onto a composite MRI volume form the full test dataset.



(*b*) Individual saliency map for a 60-year-old male subject, thresholded to T1R and overlaid onto the subject's MRI volume

Figure 4: Sections of T1R from LRP, overlaid onto example MRI volumes.

high correlation coefficients between hemispheres. A few do not have high correlations, and these tend to be regions which are assigned little T1R. This correlation between hemispheres implies a symmetry across the hemispheres, and therefore allows us to average relevance between hemispheres in our analysis.



Figure 5: Left- and right-hemisphere correlations for the T1R attributions to brain regions.

## Appendix D. BAG

For the analyses of BAG and relevance trajectories, we only use the results of a single method. We chose to implement five methods of saliency mapping and determined which agreed best with previous findings from the literature. We chose to implement LRP for the further analyses for two reasons:

- The LRP methods all tend to assign relevance preferentially to the ventricles, as opposed to the surrounding structures (as is the case with DeepLIFT). This agrees better with the findings of previous work (Levakov et al., 2020; Hofmann et al., 2021).

- The choice of LRP$_{\text{CMP}}$ with $\alpha = 1$ is the most common parameterisation of LRP in the literature, and is considered best practice (Sixt et al., 2019; Kohlbrenner et al., 2019).
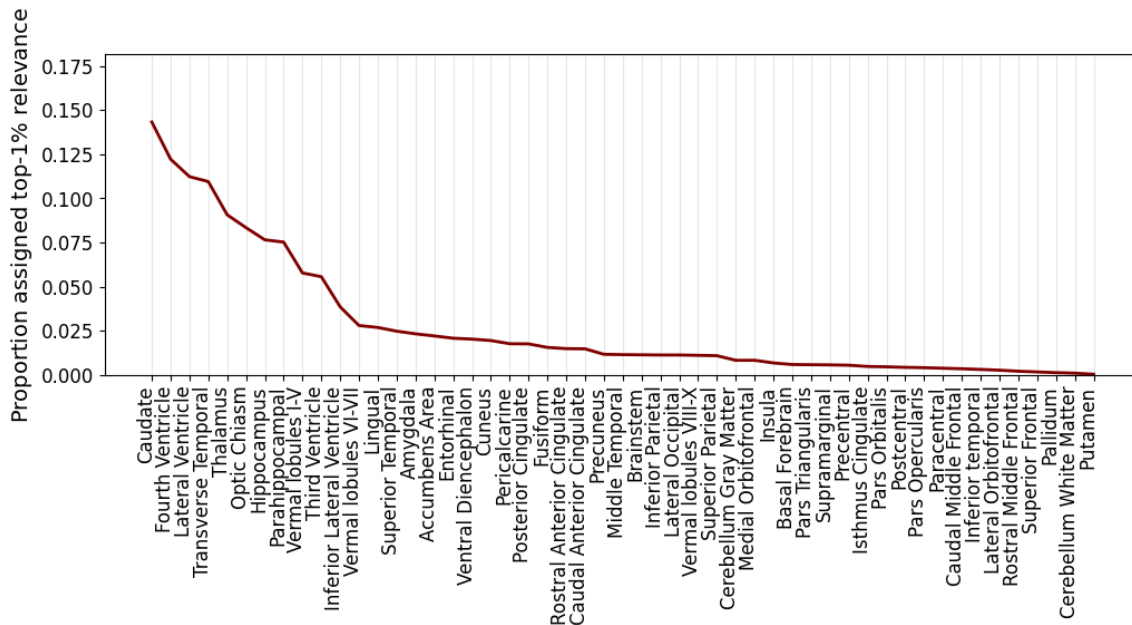
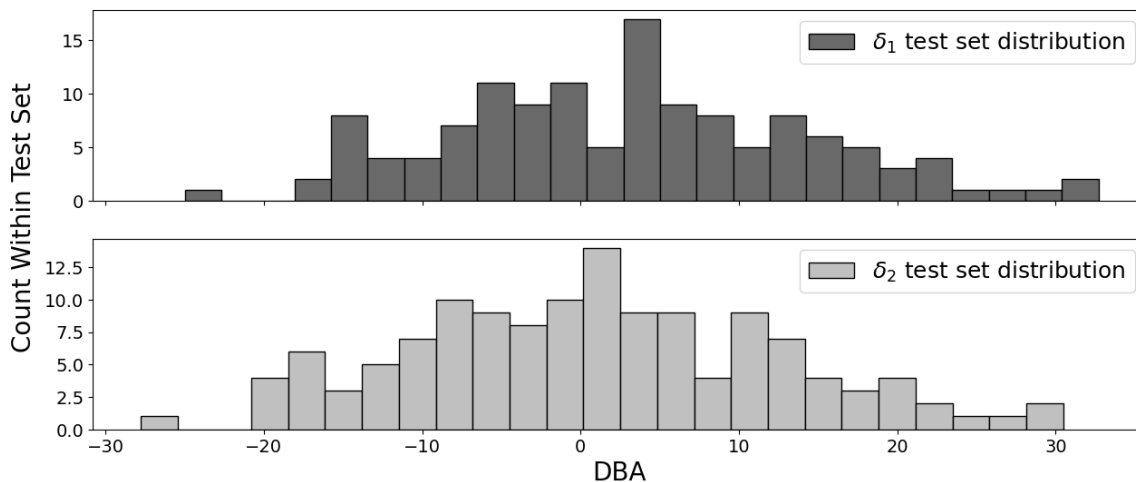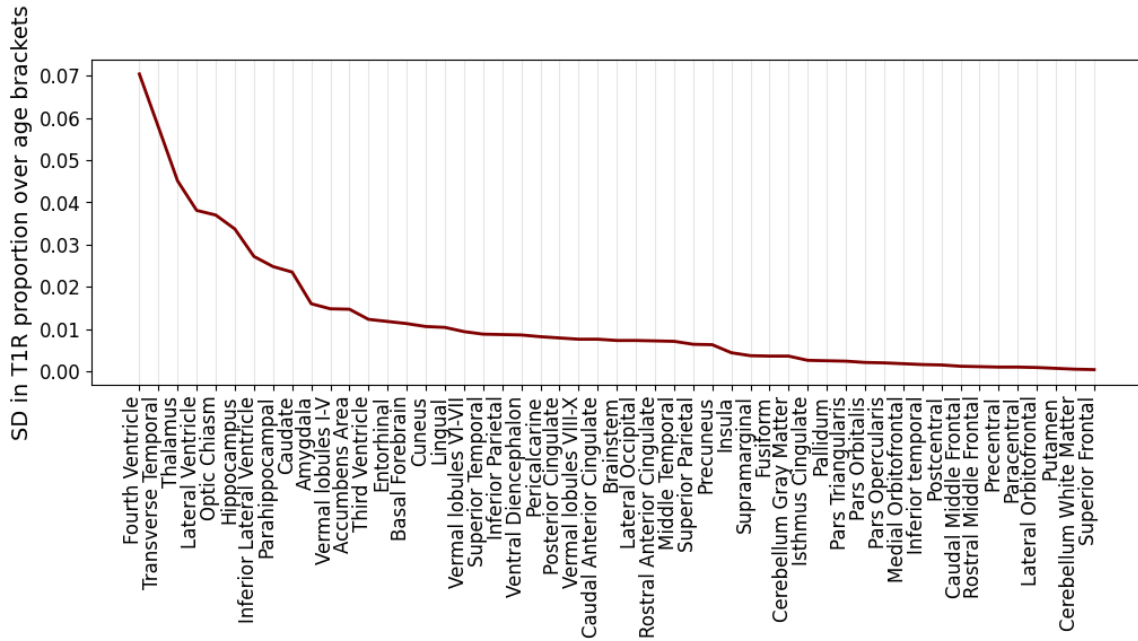Figure 6: Proportion of each region allocated T1R by LRP



Figure 7: Distributions of BAG by $\delta_1$ and $\delta_2$ (as defined by Smith et al. (2019)) on the test set. The age-orthogonality correction of $\delta_2$ shifts this distribution to be centered approximately at 0 (mean moves from 3.30y to 0.98y), and distributed more evenly to either side, becoming approximately Normal (range shifts from $(-24.98, 32.70)$ to $(-27.79, 30.47)$).
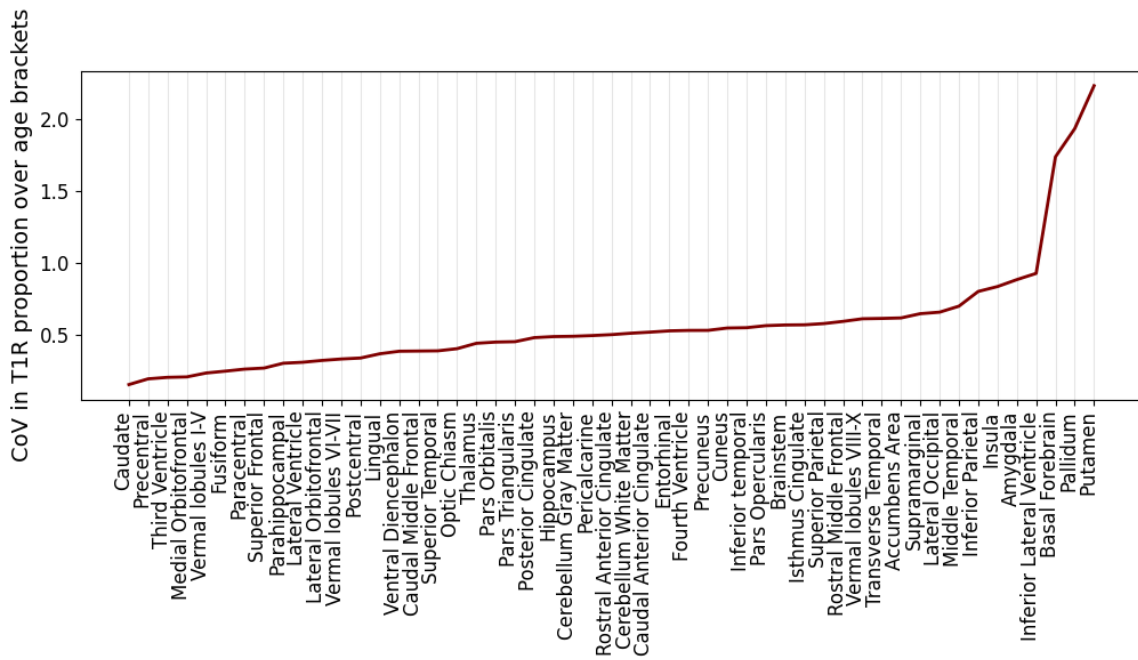
## Appendix E. Age Bracket Relevance

Figures (8a) and (8b) show the Standard Deviations and Coefficients of Variation respectively of the proportion of T1R in each brain region across age brackets for LRP.

Regions with largest SD are those salient areas which change most in their assignment of relevance over ages. Regions with low SD often tend to be regions that are simply assigned very low proportions of T1R, such as the Cerebellum White Matter.

We choose to use the metric of CoV to find those salient regions which change least in their relevance assignment over ages. The caudate nucleus, for example, has a large proportion of assigned T1R, and changes very little in this assignment across age brackets. On the other hand, regions with high CoV tend to be regions that are assigned low relevance overall. These regions can have small fluctuations in relevance assignment over age brackets, which leads to large CoVs. An example is the Putamen, which has very large CoV and very low proportions of assigned T1R (also, very small SD).

(*a*) Standard Deviation. Regions in descending order for LRP.



(*b*) Coefficient of Variation. Regions in ascending order for LRP.

Figure 8: Changes in proportion of each region allocated T1R for each method, over age brackets