

Identifying healthy individuals with Alzheimer neuroimaging phenotypes in the UK Biobank

Tiago Azevedo¹, Richard A.I. Bethlehem^{2,3}, David J. Whiteside⁴, Nol Swaddiwudhipong⁴, James B. Rowe⁴, Pietro Lió¹, and Timothy Rittman^{4,*}

¹Department of Computer Science and Technology, University of Cambridge, UK

²Brain Mapping Unit, Department of Psychiatry, University of Cambridge, UK

³Autism Research Centre, Department of Psychiatry, University of Cambridge, UK

⁴Department of Clinical Neurosciences and Cambridge University Hospitals NHS Trust, University of Cambridge, UK

*tr332@medschl.cam.ac.uk

ABSTRACT

Identifying prediagnostic neurodegenerative disease is a critical issue in neurodegenerative disease research, and Alzheimer's disease (AD) in particular, to identify populations suitable for preventive and early disease modifying trials. Evidence from genetic studies suggest the neurodegeneration of Alzheimer's disease measured by brain atrophy starts many years before diagnosis, but it is unclear whether these changes can be detected in sporadic disease. To address this challenge we train a Bayesian machine learning neural network model to generate a neuroimaging phenotype and AD-score representing the probability of AD using structural MRI data in the Alzheimer's Disease Neuroimaging Cohort (cut-off 0.5, AUC 0.92, PPV 0.90, NPV 0.93). We go on to validate the model in an independent real world dataset of the National Alzheimer's Coordinating Centre (AUC 0.74, PPV 0.65, NPV 0.80), and demonstrate correlation of the AD-score with cognitive scores in those with an AD-score above 0.5. We then apply the model to a healthy population in the UK Biobank study to identify a cohort at risk for Alzheimer's disease. This cohort have a cognitive profile in keeping with Alzheimer's disease, with strong evidence for poorer fluid intelligence, and with some evidence of poorer performance on tests of numeric memory, reaction time, working memory and prospective memory. We found some evidence in the AD-score positive cohort for modifiable risk factors of hypertension and smoking. This approach demonstrates the feasibility of using AI methods to identify a potentially prediagnostic population at high risk for developing sporadic Alzheimer's disease.

Introduction

A critical task in dementia research is to identify disease at the earliest possible timepoint, permitting early intervention with lifestyle change [1] or disease modifying therapies [2] at a time when the disease process could potentially be reversed or halted, and quality of life remains high. The difficulty in achieving early and accurate diagnosis has been highlighted as a major factor in the lack of success of clinical trials for neurodegenerative diseases, including Alzheimer's disease [2, 3]. Neuroimaging abnormalities in genetic dementia cohorts suggest that neurodegenerative pathologies begin decades before symptoms [4, 5]. Predicting disease with such certainty before symptom onset is not possible in sporadic forms of dementia, so an alternative strategy is needed to identify an at-risk population using disease biomarkers to find people with early stages of neuropathology who are at high risk of developing cognitive impairment in the future. This high risk group would be suitable for prevention studies or early disease modifying treatment trials [6].

The challenge of identifying disease at the earliest possible point has led to proposed criteria for at-risk or presymptomatic Alzheimer's disease that rely on biomarker evidence rather than a clinical syndrome [7]. One set of criteria propose an "ATN" classification of Alzheimer's disease, representing Amyloid (A), Tau (T) and Neuronal loss (N) as central pillars of Alzheimer pathology [8]. Many of the biomarkers to assess tau and amyloid pathology in life are expensive and not widely available. However, neuronal loss is readily measured *in vivo* using structural brain imaging.

Structural neuroimaging has been central in clinical diagnosis and in attempts to classify Alzheimer's disease using neuroimaging for many years [9, 10]. Loss of volume in the hippocampus is well described in Alzheimer's disease, and whole brain volume may also be relevant [11, 12]. Other specific brain regions are less well studied, yet may be relevant in identifying people with Alzheimer's disease - alone or in combination with hippocampal atrophy. More complex analytical approaches offer the opportunity to use all the available information from structural neuroimaging data for identifying a specific pattern of atrophy relevant to disease.

34 Artificial Intelligence (AI) and Machine learning (ML) describe computational algorithms that can make predictions
35 reflecting intuitive human thinking, and can ‘learn’ from new data. A class of AI models called deep learning methods use
36 multiple hierarchical levels of data abstraction to identify important features in order to make a prediction. These models
37 have been successfully applied to several different contexts in medicine [13, 14], leveraging neuroimaging datasets [15] to
38 address a multitude of neuroscientific questions [16]. AI approaches in structural MRI have facilitated the classification of
39 Alzheimer’s disease with a good degree of accuracy [17, 18, 19, 20, 21, 22, 23, 24, 25, 26], but few such studies have validated
40 their approach in an independent dataset [27, 28].

41 Despite these achievements in the neuroimaging field, there are challenges to the generalisability of deep learning
42 models [29, 30]. A relatively recent trend applies a probabilistic approach to deep learning by using measures to describe
43 Bayesian uncertainty [31, 32]. Such uncertainty measures allow for a better characterisation of the model’s output rather than
44 solely a deterministic value [33], ultimately strengthening the confidence in results derived from these stochastic models [34].

45 Probabilistic AI approaches are strong candidates to make the best use of all available information in structural neuroimaging.
46 The availability of large open access neuroimaging repositories permits us to use distinct datasets for training an AI model and
47 assessing its generalisability. In this work, we use a selective and well characterised dataset of Alzheimer’s disease to train the
48 model (Alzheimer’s Disease Neuroimaging Initiative dataset, ADNI), and a more ‘noisy’ real world clinical dataset with a
49 range of different diseases to assess generalisability (National Alzheimer’s Coordinating Center, NACC).

50 To identify a group of people at high risk of developing dementia, we use the trained model to find people with an
51 neuroimaging AI derived phenotype of Alzheimer’s disease in a healthy cohort without a diagnosis of dementia from the UK
52 Biobank study. We demonstrate poorer cognitive performance in people with an AD-like neuroimaging profile, suggesting that
53 this group have a high prevalence of early Alzheimer pathology and may be suitable for screening and selection into disease
54 modifying trials. We find that this group report poorer general health and identify hypertension and smoking as potential
55 modifiable risk factors in this cohort.

56 Methods

57 Datasets

58 The ADNI study recruits people with Alzheimer’s disease, Mild Cognitive Impairment (MCI), and control participants. It is
59 primarily a research cohort and has a well characterised population who have undergone high quality, standardised neuroimaging
60 with a standard battery of cognitive and clinical assessments. We used 736 baseline scan sessions from the ADNI dataset with a
61 diagnosis of Alzheimer’s disease (n=331) and Controls (n=405) whose demographic data is summarised in table 1.

Table 1. Summary of demographics of the ADNI dataset.

Diagnosis	n	mean age (sd)	Sex (male/female)
AD	331	75 (7.8)	181/150
Control	405	74.7 (5.7)	202/203

62 Because ADNI is a relatively select research cohort, it is vulnerable to selection bias [35]. We therefore used the NACC
63 dataset for validation. This is a “real world” memory clinic based cohort, including people with Alzheimer’s disease and a range
64 of other cognitive and non-cognitive disorders. Because of its pragmatic nature, the NACC dataset is more heterogeneous in the
65 quality of imaging and additional data collected. Therefore, it is an ideal dataset for validation of a tool developed in a more
66 ‘clean’ dataset such as ADNI. We used 5209 people from the NACC dataset whose demographics are summarised in table 2.

Table 2. Summary of demographics of the NACC dataset.

Diagnosis	n	mean age (sd)	Sex (male/female)
Control	2,824	68.6 (10.9)	938/1,886
AD	1,706	73.9 (9)	794/912
Other degenerative disorders	326	71.2 (9.9)	196/130
Other non-degenerative disorders	353	69.1 (10)	135/218

67 Finally, to apply the algorithm to a healthy cohort, we used the UK Biobank as a non-clinical cohort. This dataset is subject
68 to potential selection bias, tending to be a population with a low risk for disease [36]. Despite these limitations, the size of

69 the dataset, the age of participants and the high quality neuroimaging data makes it an ideal cohort in which to assess at-risk
70 features for neurodegenerative disease. A summary of the UK Biobank neuroimaging data is found in table 3.

Table 3. Demographics for those in the UK Biobank who underwent neuroimaging.

n	mean age (sd)	Sex (male/female)
37,104	55.3 (7.4)	19,493/17,611

71 ADNI Preprocessing

72 In each cohort, structural MRI Magnetization Prepared - Rapid Gradient Echo (MPRAGE) scans were acquired. Further details
73 of the individual imaging protocols are available for ADNI at <http://adni.loni.usc.edu/methods/documents/mri-protocols/>, for
74 NACC at <https://files.alz.washington.edu/documentation/rdd-imaging.pdf>, and for the UK Biobank at [37]. Scans underwent
75 estimation of regional cortical volume, regional cortical thickness, and estimated total intracranial volume using the FreeSurfer
76 tool box (version 6.0)¹. Given the size of the cohorts, the resulting segmentations were assessed for gross abnormalities, but
77 minor registration errors were not corrected. Results were obtained for cortical thickness and volume in the 68 surface-based
78 regions of the Desikan-Killiany atlas from both hemispheres. In addition, the brainstem volume was also extracted together with
79 9 volume features per hemisphere (cerebellum white matter, cerebellum cortex, thalamus proper, caudate, putamen, pallidum,
80 hippocampus, amygdala, and accumbens area). In total, 155 features were extracted per brain scan.

81 The ADNI dataset was divided into a training set with 662 samples, and a validation set with 74 samples, representing
82 approximately 90% and 10% of the original cohort respectively. This division approximately preserved the relative distributions
83 of diagnosis, estimated total intracranial volume, sex, and age.

84 To regress out confounds from each feature, independent linear regression models were fitted to the training set using
85 ordinary least squares (OLS) implemented in *statsmodels* [38]. For each one of the 68 cortical thickness features, the
86 independent variable to be regressed out was age. For the remaining volume features, the independent variables were age,
87 estimated total intracranial volume, and sex. These 155 regression models were saved to be later employed on the validation set
88 and other external datasets. For numerical stability when training a neural network, all features were independently scaled to
89 zero mean and unit variance using *Scikit-learn* [39]. Statistics were saved for each feature so they could be used to scale the
90 values in the validation set and other external datasets.

91 Bayesian Machine Learning

92 A supervised machine learning (ML) model learns a target function f_{θ} , parameterised by θ , such that it can predict $\mathbf{y} = f_{\theta}(\mathbf{x})$.
93 In the case of a classification task, the function is such that $f: \mathbb{R}^N \rightarrow \{1, \dots, k\}$, where k is the number of possible categories
94 (i.e. labels). For example, for a certain image with pixels represented in a feature vector \mathbf{x} , the function could try to predict
95 whether it contains a dog, a cat, or a bird ($k = 3$); in our context, the binary classification model predicts whether a patient has
96 Alzheimer’s disease or not ($k = 2$). Practically, this function f_{θ} learns how to predict labels \mathbf{y} from features \mathbf{x} by estimating the
97 probability distribution $p(\mathbf{y}|\mathbf{x})$ that generated those same labels.

98 The function f_{θ} can be modelled as a deep neural network. To train such a model with a particular dataset, one needs to
99 tune the learnable parameters of that model (i.e. θ) by minimising a loss function using stochastic gradient descent or another
100 optimisation algorithm. In contrast, under Bayesian ML the Bayes rule is used to infer model parameters θ from data \mathbf{x} :

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})}. \quad (1)$$

101 Here, the model parameters are represented by the posterior distribution $p(\theta|\mathbf{x})$, where the model parameters θ are
102 conditioned on the data \mathbf{x} . The goal of Bayesian ML is then to estimate this distribution given the likelihood $p(\mathbf{x}|\theta)$ and the
103 prior distribution $p(\theta)$ (i.e. belief of what the model parameters might be). The prior $p(\mathbf{x})$ cannot be generally computed but
104 as it is a normalising constant not dependent on θ and it stays the same for any model, it can be dropped from calculations
105 when estimating the posterior. The posterior distribution cannot usually be analytically calculated using big data in a practical
106 way, and therefore there are several methods to calculate these distributions and approximate the intractable posterior [40].

107 We use Monte Carlo dropout [41, 42] to approximate Bayesian inference by using dropout during the inference phase
108 of the model [43]. Dropout is a regularisation approach often employed in deep neural networks to avoid overfitting and it

¹<https://surfer.nmr.mgh.harvard.edu/fswiki/BrainstemSubstructures>

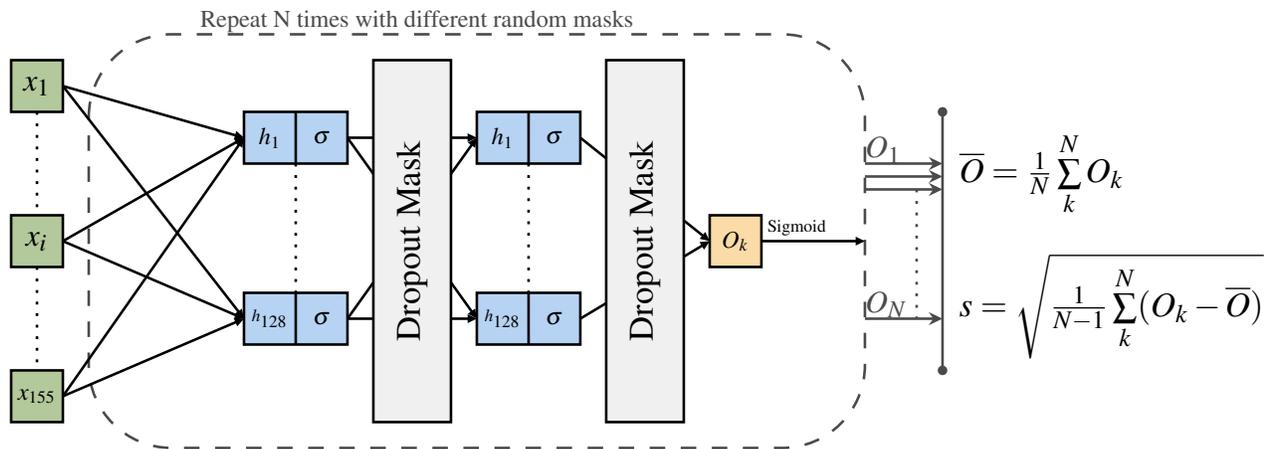


Figure 1. Architecture of the neural network model used in this paper. The neural network consists of two hidden layers of 128 dimensions and non-linear activation function $\sigma = \tanh()$. For each set of 155 inputs, $N = 50$ forward passes are run, each time with a different dropout mask sampled from a Bernoulli distribution. An AD likelihood score is generated as the mean, and model uncertainty as the standard deviation calculated from the 50 forward passes.

works by randomly dropping nodes during the training process. With Monte Carlo dropout, nodes are also randomly dropped during inference which means that for the same input, each forward pass will generate a different output; this is possible as for each pass a different Bernoulli mask is applied to the neural network's weights. Gal and Ghahramani [41] show that each forward pass on the neural network corresponding to a different dropout mask is a good approximation to sampling from the true posterior distribution $p(\theta|\mathbf{x})$.

With this simple yet powerful approximation, one can have the statistical power of a Bayesian ML model at very little added computational cost. Indeed, the Monte Carlo dropout method was chosen as it works well on a wide variety of previously trained neural networks, therefore could be used in other clinical contexts without the requirement for a full knowledge of Bayesian statistics. Furthermore, Monte Carlo dropout is known to bring advantages in modelling uncertainty [41], which is of paramount importance in a clinical context, as well as better overall performance for certain downstream tasks [44].

Deep Neural Network Implementation

As depicted in figure 1, we implemented a neural network with two hidden layers, each with 128 dimensions and using the hyperbolic tangent function ($\tanh()$) as the non-linear activation function to leverage both the positive and negative value ranges of the input. The *sigmoid* function was applied to the last output node to give a value between 0 and 1 to represent a probability that the individual has Alzheimer's disease. The dropout rate was set to 80% and Monte Carlo dropout was employed by sampling (i.e. making a forward pass) 50 times from the model, after which a mean and standard deviation was calculated. The mean corresponds to the final model prediction (i.e. probability of Alzheimer's disease) and standard deviation represents the uncertainty of the model [42]).

The model was implemented using Pytorch [45] and trained for 100 epochs using an Adam optimiser [46] with learning rate of 0.001, weight decay of 0.0001, and binary cross entropy loss. The training procedure took 9 seconds on a server with a TITAN X Pascal GPU and an Intel(R) Core(TM) i7-6900K CPU with 16 cores. The model with the smallest loss on the validation set during the training procedure was selected as the final model for evaluation. Inference time (i.e. 50 forward passes with output calculation) took an average of 12.7 ms (std: 1.78 ms) on GPU (average calculated over 1000 runs for the same batched input). The training log was saved using *Weights & Biases* [47]. In total, the model contained 36,609 trainable parameters.

Statistical Analysis

To assess group differences in the association between AD scores and clinical measures we used a Bayesian statistical approach given the different sizes of the cohorts used in this study, and the limitations of frequentist analysis in identifying statistically significant but clinically irrelevant group differences. We used Stan [48, 49] implemented in R (version 4.1.0) using linear regression and logistic regression implemented in the *brms* library [50, 51], and the *rstan* library for piecewise linear regression. To assess evidence for group differences we use the Region of Practical Equivalence (ROPE), which is an a priori effect size considered to be significant between groups. The 95% distribution of the Bayesian posterior is termed the Critical Interval

(CI); if the mean lies outside the ROPE there is some evidence to accept a hypothesis, between groups, and where the CI lies outside the ROPE there is strong evidence for accepting a hypothesis [52]. Where the CI lies completely within the ROPE, the null hypothesis can be accepted. The ROPE is either set by knowledge of the variable, or set to be 0.1 of the standard deviation of the control group. Model comparison used the *loo* package [53]. To assess the validity of our chosen breakpoint against variable or no breakpoint, we used the Expected Log Pointwise Predicted Density (ELPD) as the measure of model fit, assessing the difference in ELPD value between models and its standard error to consider whether there was evidence of a difference between models.

Results

Model Evaluation and Performance

We trained our deep learning model to detect Alzheimer’s disease from structural neuroimaging using the ADNI dataset. To evaluate model performance in the test set of the ADNI cohort and the NACC cohort, we report ROC curve analyses in table 4. For the NACC dataset we evaluated AD identification against two comparator groups: (1) controls alone, and (2) combined controls and non-AD diagnoses. As expected given the similarity to the training set and selective nature of the cohort, the highest accuracy was found in the ADNI test set. In NACC, a completely independent dataset, accuracy was lower, but still reasonable and in line with a previous similar study using SVM for out-of-distribution classification of AD [54]. Overall accuracy was above 0.7, although with a loss in positive predictive value (0.56). Of particular importance, the negative predictive value remained relatively high (0.83). This means our algorithm is balanced toward missing some people with Alzheimer’s disease, but is less likely to label healthy people as having an Alzheimer’s disease neuroimaging phenotype. This is a desirable property of the model given the application to UK Biobank data where the rate of Alzheimer’s disease will be substantially lower than either ADNI or NACC, so there is a greater risk of misclassifying healthy people as having Alzheimer’s disease.

Table 4. Performance metrics across datasets with a model trained on the ADNI training set, using a cut-off of and AD score of 0.5 and employing inference using MC Dropout with 50 samples. AUC=Area under the ROC curve. PPV=Positive predictive value. NPV=Negative predictive value.

Dataset	Accuracy	AUC	Sensitivity	Specificity	PPV/Precision	NPV
ADNI test set	0.92	0.97	0.90	0.93	0.90	0.93
NACC (only AD/Control)	0.74	0.79	0.68	0.78	0.65	0.80
NACC (AD/All)	0.72	0.76	0.68	0.73	0.56	0.83

We investigated the relationship between uncertainty measures generated by the model and the predicted value (AD score) in figure 2a. There was a wider range of uncertainty values when the average AD score was closer to 0.5 than closer to the extremes; in other words when the probability of classification was greater (towards 0 or 1) the AD score was more certain. Figure 2b demonstrates that when the model prediction was incorrect, its corresponding uncertainty value was higher on average compared to correct predictions.

We further compared the Bayesian ML model (including calculation of uncertainty with multiple passes) to a non-stochastic (single pass) one. In our analysis explained in detail in supplementary figures S1-S3, the Bayesian ML model consistently achieved better performance.

ADNI

Clinical scores

To assess the clinical validity of the AD score, we assessed the difference in clinical scores between those categorised as positive or negative by AD score using a cut-off of 0.5 and applying Bayesian regression models with age as a covariate; the posterior distributions are shown in figure 3. Skewed Gaussian families were used for MMSE and CDR Sum of Boxes, otherwise Gaussian distributions were assumed with cauchy distribution priors in all cases. All models converged well ($\hat{R} \approx 1.00$). We report four key cognitive measures from the ADNI dataset, finding very strong evidence for difference between AD score positive and negative groups in MMSE (Effect size -5.2, 95% Credible Interval -5.5 to -4.8), MoCA (-8.2, CI -9.1 to -7.4), CDR (3.7, CI 3.5 to 3.9), and Trails B (-13.0, CI -13.6 to -12.4).

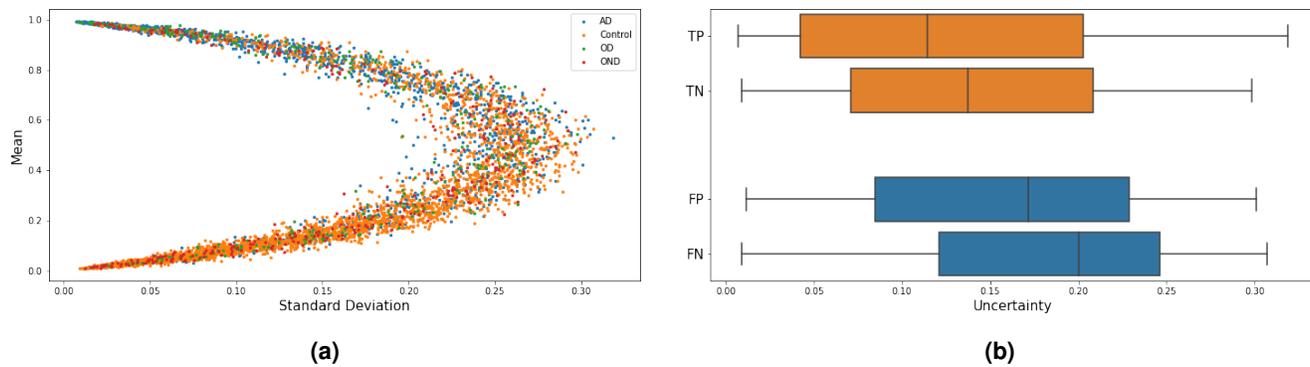


Figure 2. Model uncertainty for the NACC dataset, where uncertainty was measured as the standard deviation of the model's sampled outputs. (a) Relation of model's output and uncertainty. The model was more certain (i.e. smaller standard deviation) for more extreme mean outputs. For a mean output closer to 0.5, more variable and generally greater uncertainty was seen. (b) Uncertainty levels for different categories in the confusion matrix applying a cut-off of 0.5. On average, uncertainty levels are higher for incorrect predictions (i.e. FP, FN) when compared to correct predictions (i.e. TP, TN). There was a significant difference among these four groups (Kruskal-Wallis H-test, $p < 3.79 \times 10^{-58}$).

NACC

Clinical scores

We applied the trained model to the NACC datasets and assessed the relation of the model derived AD-score against clinical scores. Group differences were assessed with Bayesian analysis using the ROPE to assess the strength of evidence, shown in figure 4. There was strong evidence that people with a positive AD score had lower MMSE scores (-3.82, CI -4.62 to -3.02), MoCA scores (-7.00, CI -8.33 to -5.69), semantic fluency (-4.64, CI -5.49 to -3.80) and executive function (time taken to complete trails B 44.43, CI 33.36 to 55.63). For WAIS scores (-6.61, CI -9.12 to -4.10) and Boston naming test (-2.97, CI -3.95 to -1.98) there was moderate evidence of a difference in that the mean effect size of the AD score positive group fell outside the ROPE but the critical interval overlapped with the AD score negative, suggesting imprecision in the estimate of the AD score negative group; this may be explained by the relatively low Positive Predictive Value so that some people with Alzheimer's disease are included in the negative AD score group. Finally there was good evidence that the AD score does not predict forward (-0.19, CI -0.57 to 0.19) or backward (-0.46, CI to -0.86 to -0.06) digit span given the distribution of the AD score positive scores is completely contained within the critical interval of the AD score negative group.

To assess whether severity of disease was associated with the strength of expression of the AD neuroimaging phenotype we regressed the AD score against z-scored clinical measures. We used piecewise linear regression analysis given that we did not expect an association in the AD score negative group (below 0.5) compared with the AD score positive group. Firstly, we assessed whether the piecewise regression model was superior to a linear model, and whether our chosen breakpoint of 0.5 was reasonable by comparing piecewise linear regression models with a fixed breakpoint of 0.5, with variable breakpoint (permitted to vary between 0.25 and 0.75), and with no breakpoint (ie completely linear). The analysis presented in table 5 shows that models including a breakpoint were superior to the model without a breakpoint for all measures where we found evidence for difference between the AD score positive and AD score negative groups, specifically MMSE, MoCA, Semantic fluency. There was no substantial difference in whether the breakpoint was fixed at 0.5 or permitted to vary for almost all measures; for the Boston naming task, the variable breakpoint analysis was a better fit than the fixed breakpoint analysis, speculatively because executive cognitive function appears later than other cognitive impairments.

We therefore proceeded with our estimated breakpoint of 0.5 to differentiate AD score positive from AD score negative scores, shown in figure 4. There was evidence of a relationship between stronger expression in the AD score positive group than the AD score negative group of the AD score with more impaired cognitive function measured by MMSE, MoCA, forward digit span, trails B and semantic fluency and the Boston naming task, all with a credible interval lying outside the range -0.1 to 0.1 standard deviations of the control group mean.

UK Biobank

Using a cut-off for the AD score of 0.5, we divided the UK Biobank cohort into AD score positive or AD score negative groups. There were 1,304 (3.4%) with a positive AD score and 36,663 (96.6%) with a negative AD score.

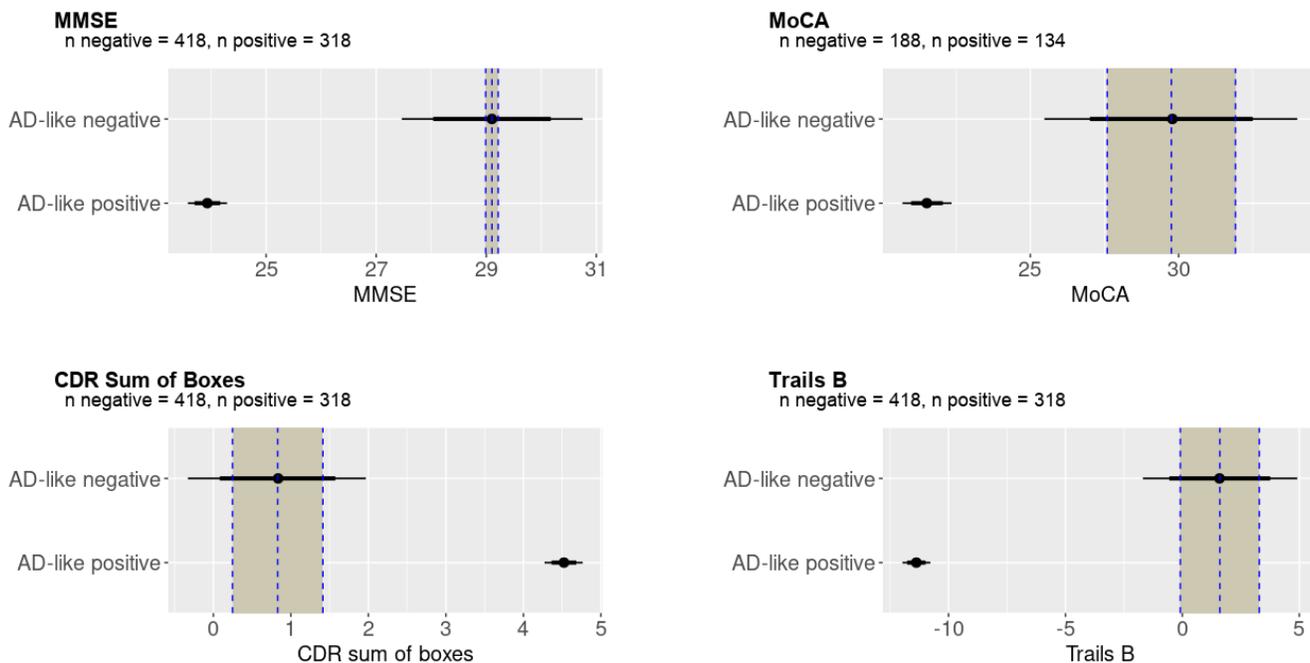


Figure 3. Bayesian analysis of cognitive tests in the ADNI dataset Bayesian posterior estimates of the mean of cognitive tests in AD score positive and negative groups with the Region Of Practical Equivalence (ROPE) as a shaded column. As expected in this well characterised dataset, for all measures there was very strong evidence of difference between groups classified as positive or negative by AD score derived from structural neuroimaging, indicated by mean AD score in the AD score positive group and the 95% credible intervals (indicated by the thin horizontal bars) falling outside the ROPE.

211 **AD scores predict cognitive differences in healthy individuals with an AD imaging phenotype**

212 To assess for differences in cognitive scores between the groups we used Bayesian linear or logistic regression models. All
 213 models achieved good convergence ($\hat{R} \approx 1$) and results are shown in figure 5.

214 There was strong evidence of worse fluid intelligence in the AD score positive group (-0.35, CI -0.46 to -0.21) with the
 215 95% CI lying completely outside the ROPE. There was moderate evidence to support poorer performance in matrix pattern
 216 completion (-0.35, CI -0.50 to -0.20), numeric memory (-0.17, CI -0.27 to -0.07), and reaction time for correct trials (13.11 ms,
 217 CI 7.12 to 19.33 ms), where the mean estimate was outside the ROPE, but the CI overlapped with the ROPE. On a working
 218 memory task (pairs matching) there was only weak evidence to suggest a poorer performance in the AD positive groups
 219 performance using logistic regression with an adjacent categories model; with AD positive participants slightly more likely to
 220 have 1 rather than 2 correct answers out of four (boundary effect size -0.14, CI -1.02 to 0.65), and slightly more likely to have
 221 2 rather than 3 correct answers (boundary effect size -1.74, I -5.19 to 0.79), and slightly more likely to have 3 rather than 4
 222 correct answers (boundary effect size 0.95, CI -0.36 to 3.46). There was also weak evidence to suggest poorer performance on
 223 a prospective memory task (increased probably of an incorrect answer in the AD positive group (0.09, CI 0.00 to 0.18).

224 On tests of executive function there was clear evidence of no difference in the number of errors on the Trails B test (0.12,
 225 CI -0.24 to 0.48) where the credible interval was completely within the ROPE, and weak evidence against an effect in tower
 226 rearranging (-0.31, CI -0.53 to -0.08) where the mean lies within the ROPE but the credible interval extends beyond the ROPE.

227 **AD score predicts worse reported overall health**

228 In non-cognitive measures, there was strong evidence that people in the AD group were more likely to report their overall
 229 health as 'poor' or 'fair' rather than 'good' or 'excellent'⁷ (probit 0.14, CI 0.09 to 0.19). There was weak evidence that hand
 230 grip was weaker in the AD positive group with a mean outside the ROPE but the CI overlapping with the ROPE (mean -1.10,
 231 CI -1.70 to -0.51). There was also weak evidence that the AD positive group were more likely to report one fall than no falls
 232 and more likely to report two or more falls than no falls (probit regression 0.07, CI 0.06 to 0.08).

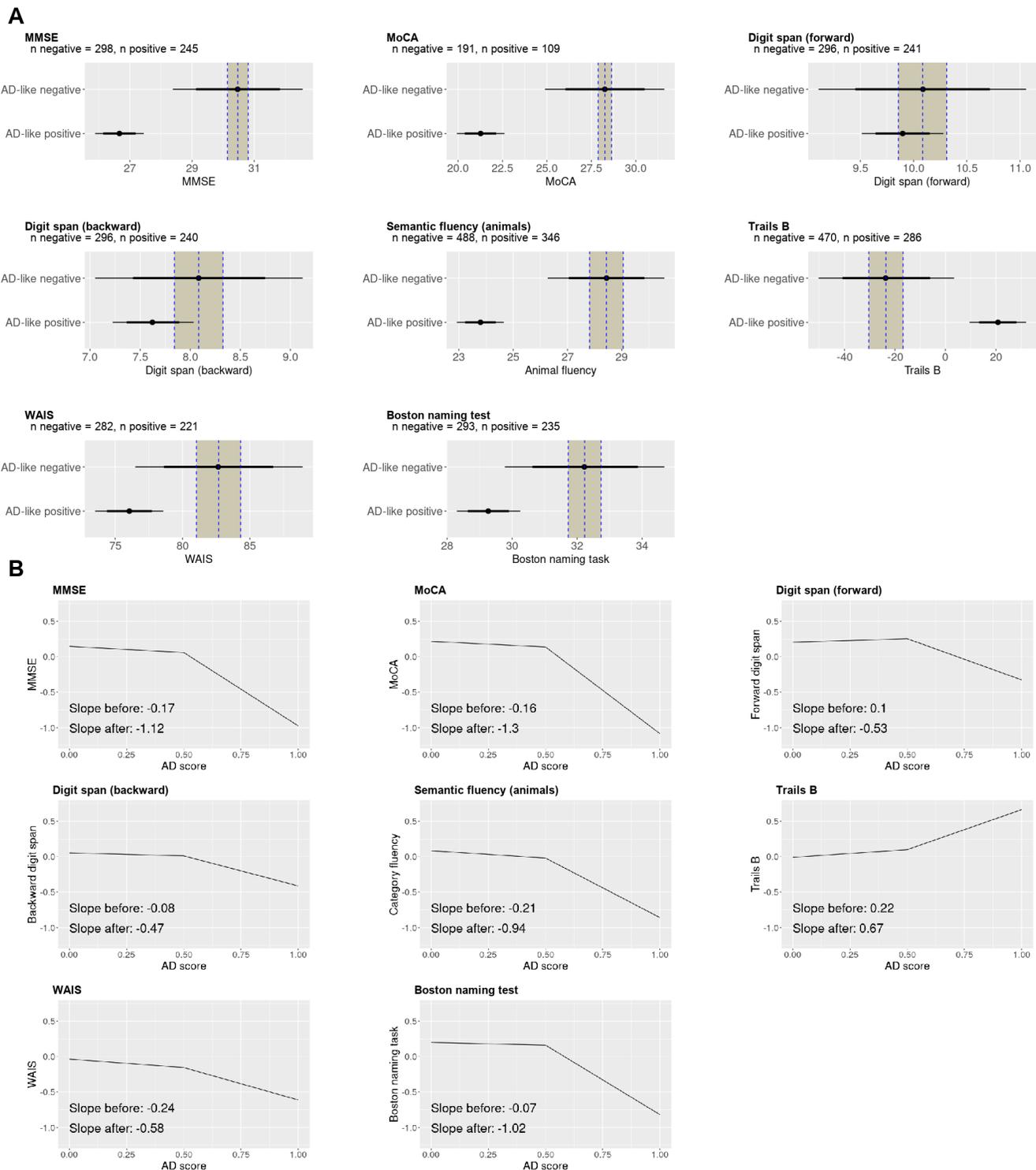


Figure 4. A: Bayesian analysis of the NACC clinical scores. There is strong evidence for impairment in the AD score positive group for MMSE, MoCA, and trails B since the posterior estimate of the effect size lies outside the 95% credible interval, and outside the Region Of Practical Equivalence (ROPE). There is good evidence of no difference for forward and backward digit span, since in both cases the distribution of the AD score positive group completely overlaps with the distribution of the AD score negative group. **B: Breakpoint analysis of the NACC clinical scores.** Disease severity correlated with the AD positive group (AD score >0.5) with evidence for difference in correlation from the AD negative (AD score <0.5) group in MMSE, MoCA, forward digits span, Trails B and the Boston naming task.

233 **AD scores are associated with modifiable risk factors**

234 Having identified a cohort potentially at risk of Alzheimer’s disease, the next step was to consider whether other health measures
 235 or modifiable risk factors are more common in this subgroup. We report the results of a number of risk factors in figure 6 and
 236 other health markers in figure 7.

237 There was some evidence of a difference in both diastolic blood pressure (1.12, CI 0.53 to 1.72) and systolic blood pressure
 238 (2.29, CI 1.26 to 3.30). Additionally, there was weak evidence that smoking (current or ex-smoker) was associated with AD
 239 positive score (0.06, CI -0.06 to 0.18) demonstrating a mean outside the ROPE, but a wide CI. Among those who smoked, there
 240 was moderate evidence that a greater smoking history (i.e. more pack years) was associated with an AD positive score (2.98, CI
 241 1.23 to 4.73).

242 There was moderately strong evidence for no difference in waist circumference (0.62, CI -0.11 to 1.35), consultation with
 243 GP for depression (logistic regression 0.03, CI -0.10 to 0.16), consultation with a psychiatrist for depression (logistic regression
 244 0.02, CI -0.19 to 0.23), hearing difficulties (logistic regression -0.01, CI -0.14 to 0.12). There was strong evidence of no
 245 difference in hip circumference (-0.12, CI -0.63 to 0.38), sleep duration (-0.01 hrs, CI 0.07 to 0.06), and neuroticism (0.11,
 246 -0.09 to 0.32) score.

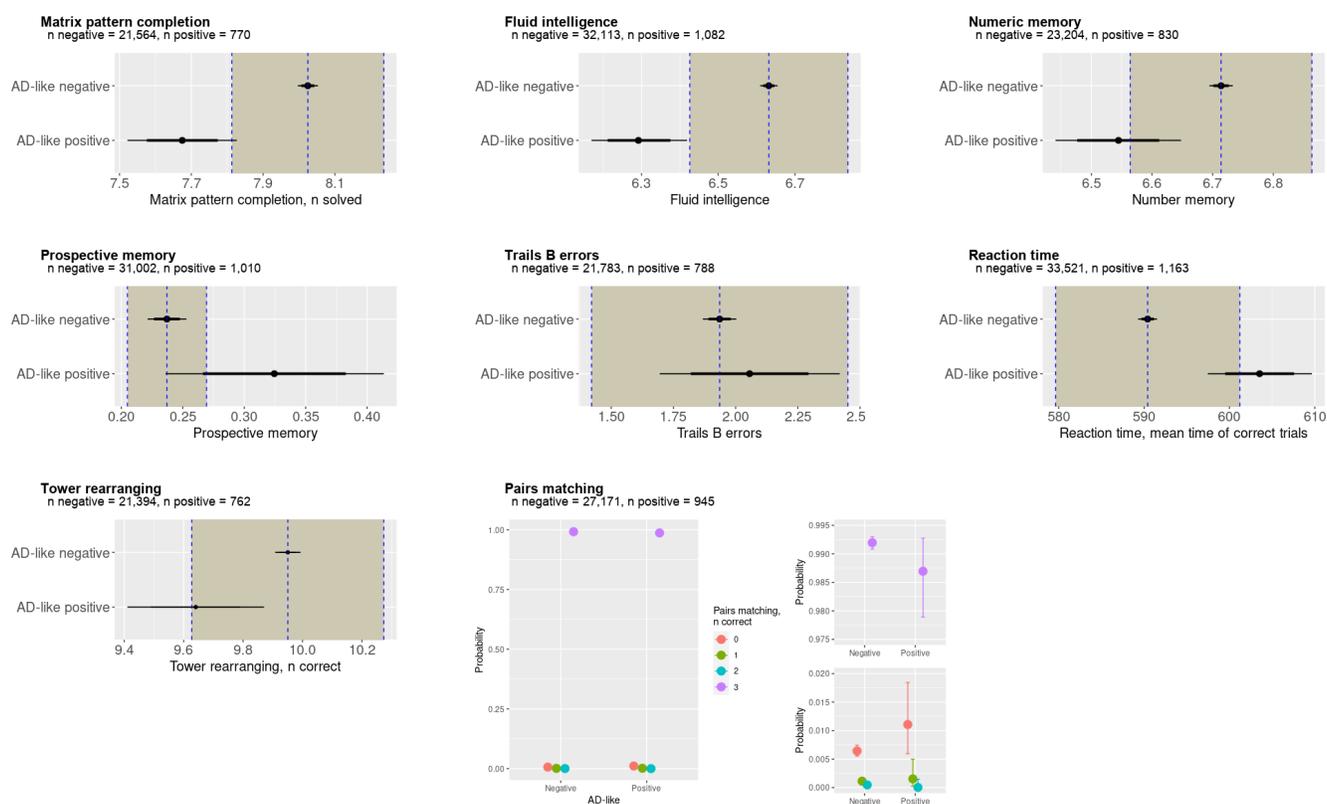


Figure 5. Bayesian analysis of cognitive tests in the UK Biobank group. It is possible to see a reduced cognitive function in participants with an AD score >0.5. In particular, there is strong evidence for impaired visual memory, and good evidence for impaired fluid intelligence, numeric memory and some evidence for impaired executive function. The shaded area represents the Region Of Practical Equivalence (ROPE) - if the distribution of the AD-positive group lies outside the ROPE there is strong evidence for a difference between the groups, and if the mean only lies outside the ROPE then there is some to good evidence for a group difference. There was strong evidence for no difference between groups in errors on the trials B task or reaction time, where the distributions lie completely inside the ROPE. For the pairs matching task we used Bayesian ordinal regression, plotting here the 95% credible intervals demonstrating only weak evidence of fewer correct answers in the AD positive group.

247 **Discussion**

248 We have identified a cohort of healthy individuals in the UK Biobank with an Alzheimer’s disease-like neuroimaging-based
 249 intermediate phenotype, by leveraging developments in Bayesian deep learning. Despite having no diagnosis or reported

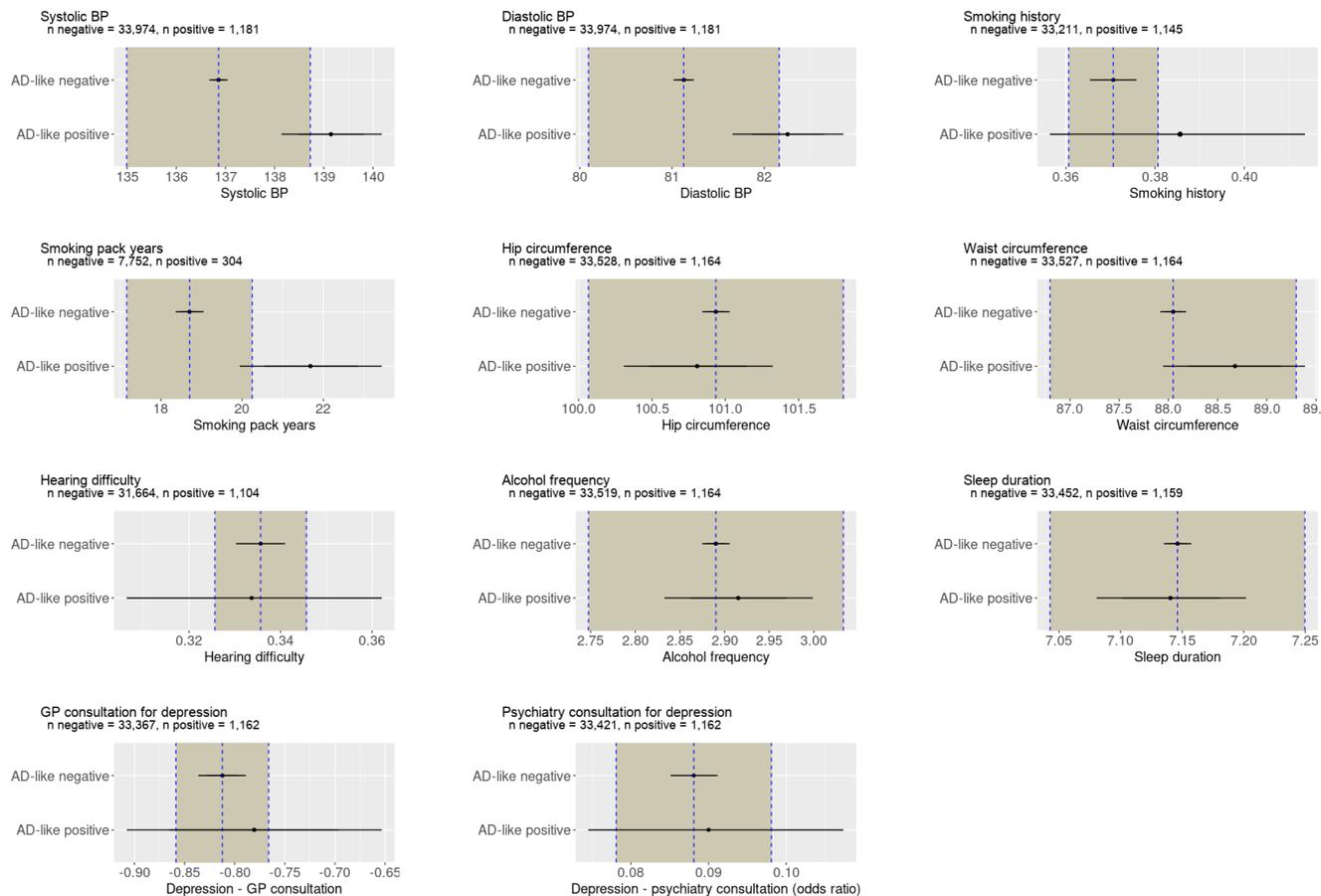


Figure 6. Results from Bayesian analysis of potentially modifiable risk factors in the UK Biobank population. There is partial evidence to support a higher diastolic and systolic blood pressure among participants with an AD score >0.5, indicated by a mean effect size lying outside the ROPE but with a distribution overlapping with the ROPE. No other risk factors were associated with a positive AD score.

250 symptoms of dementia, this AD-like cohort demonstrate a cognitive profile in keeping with early Alzheimer’s disease and
 251 report worse general health. In addition they have evidence of slightly higher blood pressure and longer smoking history as
 252 potentially modifiable risk factors.

253 Our approach offers the opportunity to identify and study presymptomatic idiopathic Alzheimer’s disease. The search for the
 254 earliest possible changes in Alzheimer’s disease has mainly focused on genetic forms of dementia [4, 10], with neuroimaging
 255 changes in presymptomatic genetic Alzheimer’s disease described since the 1990s using PET [55] or structural MRI [56]. We
 256 are aware of one promising study in idiopathic Alzheimer’s disease using a machine learning approach with multimodal imaging
 257 data to try to predict individualised presymptomatic disease in the ADNI cohort, currently in pre-print [57], an approach that
 258 will need independent validation. A study of cognitively normal adults over 70 years of age attempted to detect presymptomatic
 259 Alzheimer’s disease using FDG-PET, suggesting two-thirds of people in this age group had an abnormal FDG-PET scan which
 260 were associated with psychiatric symptoms [58]. This proportion of patients seems high for the age group under consideration,
 261 and abnormalities on PET have been associated with depression [59], so the relevance of these findings is unclear. In a small
 262 study using Pittsburgh Compound B (PiB) PET to detect presymptomatic Alzheimer’s disease in a healthy and MCI cohort there
 263 was a correlation between β -amyloid load and poorer episodic memory, though only one person converted to Mild Cognitive
 264 Impairment [60]. Another much larger study found a high rate of positive β -amyloid PET scans in otherwise cognitively
 265 normal older adults and no association with cognition, so the role and timing of β -amyloid PET abnormalities remain uncertain
 266 in the detection of presymptomatic Alzheimer’s disease [61].

267 In this context, our approach has improved on previous efforts by identifying individuals with possible early sporadic AD, a
 268 supposition that is supported by finding a cognitive profile in keeping with AD. Our findings are strengthened by identifying

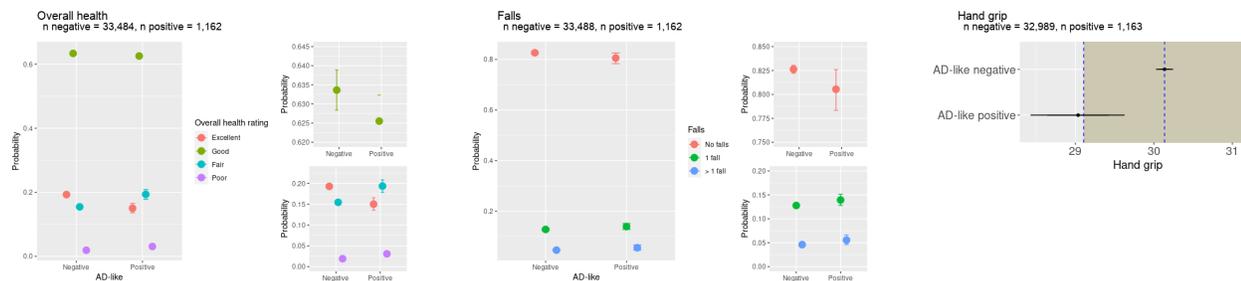


Figure 7. Other measures of health from the UK Biobank. People with positive AD scores were more likely to report their general health to be ‘fair’ or ‘poor’ and less likely to report their general health as ‘good’ or ‘excellent’. In addition, they had lower grip strength which has previously been associated with Alzheimer’s disease. There was weak evidence to suggest that people with a positive AD score were more likely to have had one or more falls in the previous year.

269 strong correlations between the AD scores and relevant cognitive tests in the independent NACC study. We found the AD score
 270 was associated with worse performance on global cognitive tests such as the MMSE and MoCA, and on more AD specific
 271 cognitive domains of memory and semantic fluency. In the UK Biobank cohort the AD score was associated with key cognitive
 272 domains of AD including memory and fluid intelligence.

273 In terms of the prospect for disease prevention, our results suggest that a smoking history, in particular a greater pack year
 274 history, and both systolic and diastolic hypertension as risk factors. Both smoking and hypertension are reported as risk factors
 275 in the 2020 Lancet Commission on Dementia [62]. Smoking is a particularly well established risk factor for dementia [63].
 276 In keeping with our findings, Rusanen et al [64] studied over 21,000 people finding that heavy smoking in middle age was
 277 associated with developing Alzheimer’s disease, and more specifically that greater cigarette use was associated with a higher
 278 risk of developing dementia. Our results suggest that the effect of smoking is mediated through structural volume loss in key
 279 brain regions.

280 The difference between blood pressure in the AD positive and AD negative groups was small, approximately 2.5mmHg for
 281 systolic BP and 1mmHg for diastolic BP. There has been much debate in the relationship between blood pressure and cognitive
 282 impairment, with studies finding both high and low diastolic blood pressure to be related to Alzheimer’s disease [65, 66]. More
 283 recent evidence from a meta-analysis has suggested that mid-life hypertension is a greater risk factor, with a systolic blood
 284 pressure above 140mmHg conferring a relative risk of 1.2 for developing dementia, and systolic blood pressure above 80mmHg
 285 conferring a relative risk of 1.54 [67]. However, the small increase in blood pressure we identified in the AD score positive
 286 group, and the overlap with the AD score negative group in both systolic and diastolic blood pressures suggests heterogeneity
 287 within the AD positive group.

288 We did not find differences in other potentially modifiable risk factors. Here the Bayesian approach is helpful, since we can
 289 confidently reject the possibility of some risk factors being associated with the AD neuroimaging phenotype in this group. For
 290 example, some risk-factors highlighted in the Lancet Commission 2020 report [62], were not identified as risks in the current
 291 study (i.e. alcohol frequency, hip circumference, sleep duration) since the distribution of the AD positive group lies wholly
 292 within the ROPE (see figure 6). For depression and hearing difficulty, there was a wide distribution of estimated risk beyond
 293 the ROPE suggesting an imprecise estimate of the risk. For these measures we cannot rule out an association with an AD
 294 neuroimaging phenotype.

295 Two factors may have limited our ability to identify potentially modifiable risk factors. Firstly, the UK Biobank has a sample
 296 bias towards people who are healthier with fewer disease risk factors than the general UK population [68]. For example, the
 297 proportion of people currently smoking in the UK Biobank population is 10.7% compared to 14.7% in the general population
 298 (data from the Office for National Statistics²).

299 Secondly, our model was biased towards a high negative predictive value, meaning that we may have ‘missed’ some people
 300 with early Alzheimer’s disease pathology. Whilst providing more confidence in the identification of an AD-like cohort, the
 301 potential classification of people with latent AD in the AD negative group may have reduced the power to detect a difference in
 302 risk factors between the AD positive and AD negative groups. We anticipate that combining neuroimaging with other risk
 303 biomarkers could improve the selection of a high risk group, for example blood biomarkers [69, 70] or polygenic risk scores
 304 [71, 72].

305 Despite these caveats, this approach has the potential to enrich dementia prevention trials. It is important to note that
 306 the impact of addressing risk factors on preventing dementia is not yet well established. The World Wide FINGERS study

²<https://www.ons.gov.uk>

307 has reported a trial of a multi-domain intervention with a small but significant effect size [1], although this was not targeted
308 at smoking cessation or lowering blood pressure specifically, and there was no difference in blood pressure between the
309 intervention and control groups at the end of the study. Our findings support the need for such trials, but raise some caution
310 about the prevalence and strength of the association between risk factors and AD pathology.

311 To identify the AD positive group we used a state-of-the-art Bayesian ML approximation method (i.e. Monte Carlo
312 dropout [41]) to identify the cohort of interest in the UK Biobank. The Bayesian approach allows a model to predict not only a
313 single AD-likelihood value as in typical deterministic neural networks, but also a measure of uncertainty (see figure 1). A key
314 advantage of this approach is the additional information about the generalisability of the model to challenging out-of-distribution
315 datasets, such as we have done in this paper; for example we were able to identify that greater uncertainty was associated with
316 incorrect predictions (see figure 2b).

317 Our approach is particularly well validated compared to other similar models. The model was trained only on the ADNI
318 dataset before validation on the completely independent and significantly more noisy NACC data, prior to application to the
319 UK Biobank. All the confound corrections on the input data were conducted in the training dataset (i.e. ADNI) alone, and
320 correction statistics are then applied to the external datasets; in this way we avoid biases that would have been introduced had
321 we corrected the model on all the available data.

322 There are limitations to our approach. Most importantly, we do not know at present whether the people identified as having
323 a positive AD score will go on to develop the syndrome of Alzheimer's disease. At the time of analysis only 17 people in the
324 neuroimaging cohort have developed dementia (6 of these self-reported at the baseline visit). The neuroimaging sub-study
325 began later than the main biobank study, so it may be some years before a sizeable population of people with dementia and
326 neuroimaging is available. Despite this caveat, we propose that the group we identified from their AD-like imaging phenotype
327 are at higher risk of future clinical Alzheimer's disease.

328 Whilst our model performed very well in the ADNI population in which it was trained, it performed, as expected, less
329 well in the independent NACC population. There are a number of reasons for this. Firstly, there is a recognised selection bias
330 when using the ADNI cohort which may lead to an overly optimistic classification [35]. Secondly, the NACC dataset relies on
331 clinical diagnosis rather than a defined set of diagnostic criteria without pathological information or biomarkers such as CSF,
332 therefore a lower diagnostic accuracy might be expected. Thirdly, the neuroimaging quality varies significantly in the NACC
333 dataset, for example both 1.5T and 3T MRI scans were included. For these reasons, it is not surprising that the classification
334 was poorer in the NACC dataset, though still with good metrics for the task at hand.

335 Using Bayesian statistics for group comparison and regression models provided several clear advantages for this study.
336 Firstly, given the unequal sizes of the positive and negative groups we were able to focus on the precision of parameter estimates
337 given the available data which differed between the two groups; this meant that we could distinguish a small effect size from an
338 imprecise parameter estimate. Secondly, we were able to use effect size to detect evidence of difference between groups; if we
339 had used a traditional frequentist approach we would have had difficult choices about correction for multiple comparisons and
340 concern about detecting small but clinically irrelevant differences. Finally, using Bayesian analysis enabled us to explicitly
341 accept the null hypothesis (i.e. no difference between groups) in a number of statistical comparisons.

342 In conclusion, we demonstrate an approach to identify a cohort of potentially presymptomatic sporadic Alzheimer's disease
343 using AI with structural neuroimaging to identify a neuroimaging phenotype.

344 Acknowledgements

345 T.A. was funded by the W. D. Armstrong Trust Fund, University of Cambridge, UK. TR is supported by the Cambridge Centre
346 for Parkinson-plus and NIHR Cambridge Biomedical Research Centre (BRC-1215-20014). The views expressed are those of
347 the authors and not necessarily those of the NIHR or the Department of Health and Social Care. JBR is supported by the Medical
348 Research Council (SUAG/051 R101400). Data collection and sharing for this project was funded by the Alzheimer's Disease
349 Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI (Department of Defense
350 award number W81XWH-12-2-0012). ADNI is funded by the National Institute on Aging, the National Institute of Biomedical
351 Imaging and Bioengineering, and through generous contributions from the following: AbbVie, Alzheimer's Association;
352 Alzheimer's Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir,
353 Inc.; Cogstate; Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its
354 affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research &
355 Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck & Co.,
356 Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation;
357 Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes
358 of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated
359 by the Foundation for the National Institutes of Health (www.fnih.org). The grantee organization is the Northern California
360 Institute for Research and Education, and the study is coordinated by the Alzheimer's Therapeutic Research Institute at the

361 University of Southern California. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of
362 Southern California. The NACC database is funded by NIA/NIH Grant U24 AG072122. NACC data are contributed by the
363 NIA-funded ADCs: P50 AG005131 (PI James Brewer, MD, PhD), P50 AG005133 (PI Oscar Lopez, MD), P50 AG005134
364 (PI Bradley Hyman, MD, PhD), P50 AG005136 (PI Thomas Grabowski, MD), P50 AG005138 (PI Mary Sano, PhD), P50
365 AG005142 (PI Helena Chui, MD), P50 AG005146 (PI Marilyn Albert, PhD), P50 AG005681 (PI John Morris, MD), P30
366 AG008017 (PI Jeffrey Kaye, MD), P30 AG008051 (PI Thomas Wisniewski, MD), P50 AG008702 (PI Scott Small, MD), P30
367 AG010124 (PI John Trojanowski, MD, PhD), P30 AG010129 (PI Charles DeCarli, MD), P30 AG010133 (PI Andrew Saykin,
368 PsyD), P30 AG010161 (PI David Bennett, MD), P30 AG012300 (PI Roger Rosenberg, MD), P30 AG013846 (PI Neil Kowall,
369 MD), P30 AG013854 (PI Robert Vassar, PhD), P50 AG016573 (PI Frank LaFerla, PhD), P50 AG016574 (PI Ronald Petersen,
370 MD, PhD), P30 AG019610 (PI Eric Reiman, MD), P50 AG023501 (PI Bruce Miller, MD), P50 AG025688 (PI Allan Levey,
371 MD, PhD), P30 AG028383 (PI Linda Van Eldik, PhD), P50 AG033514 (PI Sanjay Asthana, MD, FRCP), P30 AG035982 (PI
372 Russell Swerdlow, MD), P50 AG047266 (PI Todd Golde, MD, PhD), P50 AG047270 (PI Stephen Strittmatter, MD, PhD), P50
373 AG047366 (PI Victor Henderson, MD, MS), P30 AG049638 (PI Suzanne Craft, PhD), P30 AG053760 (PI Henry Paulson, MD,
374 PhD), P30 AG066546 (PI Sudha Seshadri, MD), P20 AG068024 (PI Erik Roberson, MD, PhD), P20 AG068053 (PI Marwan
375 Sabbagh, MD), P20 AG068077 (PI Gary Rosenberg, MD), P20 AG068082 (PI Angela Jefferson, PhD), P30 AG072958 (PI
376 Heather Whitson, MD), P30 AG072959 (PI James Leverenz, MD). This research has been conducted using data from UK
377 Biobank, a major biomedical database (<http://www.ukbiobank.ac.uk/>)

378 Author Contributions

379 T.A. conducted the experiments, T.A., J.R. and T.R. conceived the experiments and analysed the results, R.A.I.B. carried out
380 neuroimaging preprocessing and analysis, N.S. helped with statistical analysis, T.R. and P.L. supervised the study. All authors
381 reviewed the manuscript.

382 Competing Interests

383 The authors declare no competing interests.

384 Data Availability

385 All the data used in this study is available by application to the data managers of ADNI, NACC or the UK biobank. ADNI Data
386 used in the preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database
387 (adni.loni.usc.edu). The ADNI was launched in 2003 as a public-private partnership, led by Principal Investigator Michael W.
388 Weiner, MD. The primary goal of ADNI has been to test whether serial magnetic resonance imaging (MRI), positron emission
389 tomography (PET), other biological markers, and clinical and neuropsychological assessment can be combined to measure
390 the progression of mild cognitive impairment (MCI) and early Alzheimer's disease (AD). For up-to-date information, see
391 www.adni-info.org.

392 Code Availability

393 Code (for reproducibility) and results are publicly available on Github, with additional instructions for implementation:
394 https://github.com/tjiagoM/adni_phenotypes

395 Additional Information

396 The paper contains supplementary material.

References

1. Ngandu, T. *et al.* A 2 year multidomain intervention of diet, exercise, cognitive training, and vascular risk monitoring versus control to prevent cognitive decline in at-risk elderly people (FINGER): A randomised controlled trial. *The Lancet* **385**, 2255–2263 (2015).
2. Elmaleh, D. R. *et al.* Developing Effective Alzheimer's Disease Therapies: Clinical Experience and Future Directions. *Journal of Alzheimer's Disease* **71**, 715–732 (2019).
3. Aisen, P. S., Vellas, B. & Hampel, H. Moving towards early clinical trials for amyloid-targeted therapy in Alzheimer's disease. *Nature Reviews Drug Discovery* **12**, 324–324 (2013).

4. Rohrer, J. D. *et al.* Presymptomatic cognitive and neuroanatomical changes in genetic frontotemporal dementia in the Genetic Frontotemporal dementia Initiative (GENFI) study : A cross-sectional analysis. *Lancet neurology* **14**, 253–262 (2015).
5. Kinnunen, K. M. *et al.* Presymptomatic atrophy in autosomal dominant Alzheimer’s disease: A serial magnetic resonance imaging study. *Alzheimer’s and Dementia* **14**, 43–53 (2018).
6. Imtiaz, B., Tolppanen, A.-M., Kivipelto, M. & Soininen, H. Future directions in Alzheimer’s disease from risk factors to prevention. *Biochemical Pharmacology* **88**, 661–670 (2014).
7. Dubois, B. *et al.* Clinical diagnosis of Alzheimer’s disease: Recommendations of the International Working Group. *The Lancet Neurology* **20**, 484–496 (2021).
8. Jack Jr, C. R. *et al.* NIA-AA Research Framework: Toward a biological definition of Alzheimer’s disease. *Alzheimer’s & Dementia* **14**, 535–562 (2018).
9. Fox, N. C., Freeborough, P. A. & Rossor, M. N. Visualisation and quantification of rates of atrophy in Alzheimer’s disease. *The Lancet* **348**, 94–97 (1996).
10. Schott, J. M. *et al.* Reduced sample sizes for atrophy outcomes in Alzheimer’s disease trials: Baseline adjustment. *Neurobiology of aging* **31**, 1452–62, 1462.e1–2 (2010).
11. Sluimer, J. D. *et al.* Whole-brain atrophy rate and cognitive decline: Longitudinal MR study of memory clinic patients. *Radiology* **248**, 590–598 (2008).
12. Bethlehem, R. A. *et al.* Brain charts for the human lifespan. Preprint at <https://doi.org/10.1101/2021.06.08.447489> (2021).
13. Litjens, G. *et al.* A survey on deep learning in medical image analysis. *Medical Image Analysis* **42**, 60–88 (2017).
14. Piccialli, F., Somma, V. D., Giampaolo, F., Cuomo, S. & Fortino, G. A survey on deep learning in medicine: Why, how and when? *Information Fusion* **66**, 111–137 (2021).
15. Abrol, A. *et al.* Deep learning encodes robust discriminative neuroimaging representations to outperform standard machine learning. *Nature Communications* **12** (2021).
16. Yang, G. R. & Wang, X.-J. Artificial neural networks for neuroscientists: A primer. *Neuron* **107**, 1048–1070 (2020).
17. Amoroso, N. *et al.* Deep learning reveals Alzheimer’s disease onset in MCI subjects: Results from an international challenge. *Journal of Neuroscience Methods* **302**, 3–9 (2018).
18. Cui, R. & Liu, M. Hippocampus Analysis by Combination of 3-D DenseNet and Shapes for Alzheimer’s Disease Diagnosis. *IEEE journal of biomedical and health informatics* **23**, 2099–2107 (2019).
19. Cui, R. & Liu, M. RNN-based longitudinal analysis for diagnosis of Alzheimer’s disease. *Computerized Medical Imaging and Graphics* **73**, 1–10 (2019).
20. Gorji, H. T. & Kaabouch, N. A Deep Learning approach for Diagnosis of Mild Cognitive Impairment Based on MRI Images. *Brain Sciences* **9**, E217 (2019).
21. Huang, Y., Xu, J., Zhou, Y., Tong, T. & Zhuang, X. Diagnosis of Alzheimer’s Disease via Multi-Modality 3D Convolutional Neural Network. *Frontiers in Neuroscience* **13**, 509 (2019).
22. Jain, R., Jain, N., Aggarwal, A. & Hemanth, D. J. Convolutional neural network based Alzheimer’s disease classification from magnetic resonance brain images. *Cognitive Systems Research* **57**, 147–159 (2019).
23. Mendoza-Léon, R., Puentes, J., Uriza, L. F. & Hernández Hoyos, M. Single-slice Alzheimer’s disease classification and disease regional analysis with Supervised Switching Autoencoders. *Computers in Biology and Medicine* **116**, 103527 (2020).
24. Pan, D. *et al.* Early Detection of Alzheimer’s Disease Using Magnetic Resonance Imaging: A Novel Approach Combining Convolutional Neural Networks and Ensemble Learning. *Frontiers in Neuroscience* **0** (2020).
25. Suh, C. H. *et al.* Development and Validation of a Deep Learning–Based Automatic Brain Segmentation and Classification Algorithm for Alzheimer Disease Using 3D T1-Weighted Volumetric Images. *American Journal of Neuroradiology* (2020).
26. Spasov, S., Passamonti, L., Duggento, A., Liò, P. & Toschi, N. A parameter-efficient deep learning approach to predict conversion from mild cognitive impairment to alzheimer's disease **189**, 276–287 (2019).
27. Bzdok, D. & Yeo, B. T. Inference in the age of big data: Future perspectives on neuroscience **155**, 549–564 (2017).

28. Bzdok, D., Varoquaux, G. & Steyerberg, E. W. Prediction, not association, paves the road to precision medicine **78**, 127 (2021).
29. Hosseini, M. *et al.* I tried a bunch of things: The dangers of unexpected overfitting in classification of brain data. *Neuroscience & Biobehavioral Reviews* **119**, 456–467 (2020).
30. Mårtensson, G. *et al.* The reliability of a deep learning model in clinical out-of-distribution MRI data: A multicohort study. *Medical Image Analysis* **66**, 101714 (2020).
31. Ghahramani, Z. Probabilistic machine learning and artificial intelligence. *Nature* **521**, 452–459 (2015).
32. Wilson, A. G. & Izmailov, P. Bayesian deep learning and a probabilistic perspective of generalization. Preprint at <https://arxiv.org/abs/2002.08791> (2020).
33. Leibig, C., Allken, V., Ayhan, M. S., Berens, P. & Wahl, S. Leveraging uncertainty information from deep neural networks for disease detection. *Scientific Reports* **7** (2017).
34. Kompa, B., Snoek, J. & Beam, A. L. Second opinion needed: communicating uncertainty in medical machine learning. *npj Digital Medicine* **4** (2021).
35. Mendelson, A. F., Zuluaga, M. A., Lorenzi, M., Hutton, B. F. & Ourselin, S. Selection bias in the reported performances of AD classification pipelines. *NeuroImage: Clinical* **14**, 400–416 (2017).
36. Stamatakis, E. *et al.* Is Cohort Representativeness Passé? Poststratified Associations of Lifestyle Risk Factors with Mortality in the UK Biobank. *Epidemiology (Cambridge, Mass.)* **32**, 179–188 (2021).
37. Alfaro-Almagro, F. *et al.* Image processing and Quality Control for the first 10,000 brain imaging datasets from UK Biobank. *NeuroImage* **166**, 400–424 (2018).
38. Seabold, S. & Perktold, J. statsmodels: Econometric and statistical modeling with python. In *9th Python in Science Conference* (2010).
39. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011).
40. van de Schoot, R. *et al.* Bayesian statistics and modelling. *Nature Reviews Methods Primers* **1** (2021).
41. Gal, Y. & Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Balcan, M. F. & Weinberger, K. Q. (eds.) *Proceedings of The 33rd International Conference on Machine Learning*, vol. 48 of *Proceedings of Machine Learning Research*, 1050–1059 (PMLR, 2016).
42. Gal, Y. *Uncertainty in Deep Learning*. Ph.D. thesis, University of Cambridge (2016).
43. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* **15**, 1929–1958 (2014).
44. Ovadia, Y. *et al.* Can you trust your model’s uncertainty? Evaluating predictive uncertainty under dataset shift. *Advances in Neural Information Processing Systems* **32**, 13991–14002 (2019).
45. Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning library. In Wallach, H. *et al.* (eds.) *Advances in Neural Information Processing Systems* **32**, 8024–8035 (Curran Associates, Inc., 2019).
46. Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, Conference Track Proceedings* (2015).
47. Biewald, L. Experiment tracking with weights and biases (2020). URL <https://www.wandb.com/>. Software available from wandb.com.
48. Gelman, A., Lee, D. & Guo, J. Stan: A Probabilistic Programming Language for Bayesian Inference and Optimization. *Journal of Educational and Behavioral Statistics* **40**, 530–543 (2015).
49. Carpenter, B. *et al.* Stan: A Probabilistic Programming Language. *Journal of Statistical Software* **76**, 1–32 (2017).
50. Bürkner, P.-C. Advanced Bayesian Multilevel Modeling with the R Package brms. *The R Journal* **10**, 395–411 (2018).
51. Bürkner, P.-C. Brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal of Statistical Software* **80**, 1–28 (2017).
52. Kruschke, J. K. Rejecting or Accepting Parameter Values in Bayesian Estimation. *Advances in Methods and Practices in Psychological Science* **1**, 270–280 (2018). URL <https://doi.org/10.1177/2515245918771304>. Publisher: SAGE Publications Inc.

53. Vehtari, A., Gelman, A. & Gabry, J. Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing* **27**, 1413–1432 (2017).
54. Klöppel, S. *et al.* Applying Automated MR-Based Diagnostic Methods to the Memory Clinic: A Prospective Study. *Journal of Alzheimer's Disease* **47**, 939–954 (2015).
55. Kennedy, A. M. *et al.* Deficits in cerebral glucose metabolism demonstrated by positron emission tomography in individuals at risk of familial Alzheimer's disease. *Neuroscience Letters* **186**, 17–20 (1995).
56. Fox, N. C. *et al.* Presymptomatic hippocampal atrophy in Alzheimer's disease. A longitudinal MRI study. *Brain: A Journal of Neurology* **119** (Pt 6), 2001–2007 (1996).
57. Giorgio, J. *et al.* Predicting future regional tau accumulation in asymptomatic and early Alzheimer's disease. *Bioarxiv* 2020.08.15.252601 (2020).
58. Krell-Roesch, J. *et al.* FDG-PET and Neuropsychiatric Symptoms among Cognitively Normal Elderly Persons: The Mayo Clinic Study of Aging. *Journal of Alzheimer's Disease* **53**, 1609–1616 (2016).
59. Fu, C. *et al.* A combined study of 18F-FDG PET-CT and fMRI for assessing resting cerebral function in patients with major depressive disorder. *Experimental and Therapeutic Medicine* **16**, 1873–1881 (2018).
60. Villemagne, V. L. *et al.* Longitudinal assessment of A β and cognition in aging and Alzheimer disease. *Annals of Neurology* **69**, 181–192 (2011).
61. Rowe, C. C. *et al.* Amyloid imaging results from the Australian Imaging, Biomarkers and Lifestyle (AIBL) study of aging. *Neurobiology of aging* **31**, 1275–1283 (2010).
62. Livingston, G. *et al.* Dementia prevention, intervention, and care: 2020 report of the Lancet Commission. *The Lancet* **396**, 413–446 (2020).
63. Peters, R. *et al.* Smoking, dementia and cognitive decline in the elderly, a systematic review. *BMC Geriatrics* **8**, 36 (2008).
64. Rusanen, M., Kivipelto, M., Quesenberry, C. P., Jr, Zhou, J. & Whitmer, R. A. Heavy Smoking in Midlife and Long-term Risk of Alzheimer Disease and Vascular Dementia. *Archives of Internal Medicine* **171**, 333–339 (2011).
65. Razay, G., Williams, J., King, E., Smith, A. D. & Wilcock, G. Blood pressure, dementia and Alzheimer's disease: The OPTIMA longitudinal study. *Dementia and Geriatric Cognitive Disorders* **28**, 70–74 (2009).
66. Yuan, M., Chen, S.-J., Li, X.-L. & Xu, L.-J. Blood Pressure and the Risk of Alzheimer's Disease: Is There a Link? *American Journal of Alzheimer's Disease and Other Dementias* **31**, 97–98 (2016).
67. Ou, Y.-N. *et al.* Blood Pressure and Risks of Cognitive Impairment and Dementia: A Systematic Review and Meta-Analysis of 209 Prospective Studies. *Hypertension (Dallas, Tex.: 1979)* **76**, 217–225 (2020).
68. Fry, A. *et al.* Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank Participants With Those of the General Population. *American Journal of Epidemiology* **186**, 1026–1034 (2017).
69. Leonenko, G. *et al.* Polygenic risk and hazard scores for Alzheimer's disease prediction. *Annals of Clinical and Translational Neurology* **6**, 456–465 (2019).
70. Escott-Price, V., Myers, A., Huentelman, M., Shoai, M. & Hardy, J. Polygenic Risk Score Analysis of Alzheimer's Disease in Cases without APOE4 or APOE2 Alleles. *The Journal of Prevention of Alzheimer's Disease* **6**, 16–19 (2019).
71. Thijssen, E. H. *et al.* Diagnostic value of plasma phosphorylated tau181 in Alzheimer's disease and frontotemporal lobar degeneration. *Nature Medicine* **26**, 387–397 (2020).
72. Chhatwal, J. P. *et al.* Plasma N-terminal tau fragment levels predict future cognitive decline and neurodegeneration in healthy elderly individuals. *Nature Communications* **11**, 6024 (2020).

Table 5. Here we test the cut-off AD score value of 0.5 in the NACC dataset by applying linear analysis of the relationship between AD score and cognitive scores using Bayesian piecewise linear regression analysis. For the slope estimates we include the 95% Credible Interval (CI). Given the data are z-scored, we use 0.1 as the Region of Practical Equivalence (ROPE) which represent 0.1 of the standard deviation of the data. If the CI lies outside -0.1 to 0.1 we consider there is good evidence to of a relationship between the AD score and clinic score. For models with variable breakpoints, the breakpoint values obtained were similar to 0.5, and comparing models with the Expected Log Pointwise Predicted Density (ELPD) the difference between variable and fixed breakpoint models was negligible, except for the Boston naming task. There was good evidence for a breakpoint in MMSE, MoCA, forward digit span, semantic fluency and the Boston naming task, and further evidence supporting no difference between the two breakpoint models, but both being superior to the non-breakpoint model. These findings support our use of a breakpoint of 0.5. BP = Breakpoint, MMSE = Mini Mental State Examination, MoCA = Montreal Cognitive Assessment, WAIS = Wechsler Adult Intelligence Scale.

	Variable breakpoint				
	BP	Slope < BP (CI)	Slope > BP (CI)	Slope diff (CI)	ELPD diff (se)
MMSE	0.67	-0.19 (-1.8, -0.94)	-1.4 (-1.8, -0.94)	-1.2 (-1.7, -0.65)	0.0 (0.0)
MoCA	0.6	-0.18 (-0.35, 0.03)	-1.40 (-1.8, -1.0)	-1.2 (-1.8, -0.77)	0.0 (0.0)
Backward digit span	0.56	-0.08 (-0.25, 0.08)	-0.51 (-0.96, -0.13)	-0.43 (-1.0, 0.10)	0.0 (0.0)
Forward digit span	0.59	0.09 (-0.08, 0.26)	-0.62 (-1.1, -0.2)	-0.71 (-1.3, -0.16)	0.0 (0.0)
Semantic fluency	0.59	-0.22 (-0.34, -0.10)	-1.0 (-1.3, -0.74)	-0.8 (-0.18, -0.44)	0.0 (0.0)
Trails B	0.44	0.2 (0.079, 0.35)	0.66 (0.45, 0.89)	0.46 (0.13, 0.79)	0.0 (0.0)
WAIS	0.53	-0.24 (-0.42, -0.07)	-0.59 (-0.97, -0.25)	-0.35 (-0.82, 0.13)	-0.1 (0.1)
Boston naming task	0.66	-0.08 (-0.23, 0.056)	-1.3 (-1.7, -0.81)	-1.2 (-1.70, -0.64)	0.0 (0.0)
	Fixed breakpoint (0.5)				
	BP	Slope < BP (CI)	Slope > BP (CI)	Slope diff (CI)	ELPD diff (se)
MMSE	[0.5]	-0.17 (-0.32, -0.02)	-1.12 (-1.4, -0.82)	-0.95 (-1.4, -0.53)	-1.3 (1.0)
MoCA	[0.5]	-0.16 (-0.34, 0.03)	-1.3 (-1.6, -1.0)	-1.1 (-1.6, -0.72)	-0.5 (0.6)
Backward digit span	[0.5]	-0.08 (-0.26, 0.09)	-0.47 (-0.81, -0.13)	-0.38 (-0.87, 0.10)	-0.3 (0.2)
Forward digit span	[0.5]	0.10 (-0.07, 0.28)	-0.53 (-0.87, -0.2)	-0.63 (-1.1, -0.15)	-0.3 (0.3)
Semantic fluency	[0.5]	-0.21 (-0.34, -0.08)	-0.94 (-1.20, -0.72)	-0.74 (-1.1, -0.41)	-0.5 (0.4)
Trails B	[0.5]	0.22 (0.08, 0.36)	0.67 (0.46, 0.89)	0.46 (0.13, 0.77)	0.0 (0.2)
WAIS	[0.5]	-0.24 (-0.40, -0.07)	-0.58 (-0.87, -0.24)	-0.34 (-0.79, 0.13)	0.0 (0.0)
Boston naming task	[0.5]	-0.07 (-0.23, 0.09)	-1.0 (-1.3, -0.7)	-0.94 (-1.4, -0.49)	-1.7 (0.8)
	Linear model (no breakpoint)				ELPD diff (se)
	Slope (CI)				
MMSE	-0.47 (-0.55, -0.39)				-10.0 (5.3)
MoCA	-0.60 (-0.69, -0.51)				-12.4 (6.8)
Backward digit span	-0.2 (-0.29, -0.12)				-0.5 (1.6)
Forward digit span	-0.1 (-0.18, -0.02)				-2.6 (2.6)
Semantic fluency	-0.46 (-0.52, -0.40)				-9.6 (4.5)
Trails B	0.39 (0.36, 0.45)				-2.7 (2.8)
WAIS	-0.35 (-0.43, -0.26)				-0.3 (1.4)
Boston naming task	-0.37 (-0.45, -0.29)				-9.6 (4.7)

Table 6. Demographics of the UKBB AD score positive and negative groups the

AD score	n	Age (sd)	Male/female (%)	Handedness right/left/ambi
Negative	33529	63.9 (7.5)	52.8/47.2	88.9/9.3/1.9
Positive	1304	64.9 (8.1)	51.5/48.5	88.6/9.7/1.8