

1 **Title**

2 Constructing personalized characterizations of structural brain aberrations in patients with  
3 dementia and mild cognitive impairment using explainable artificial intelligence

4

5 **Authors**

6 Esten H. Leonardsen<sup>1,2,\*</sup>, Karin Persson<sup>3,4</sup>, Edvard Grødem<sup>1,5</sup>, Nicola Dinsdale<sup>6</sup>, Till  
7 Schellhorn<sup>7,8</sup>, James M. Roe<sup>1</sup>, Didac Vidal-Piñeiro<sup>1</sup>, Øystein Sørensen<sup>1</sup>, Tobias Kaufmann<sup>2,9,10</sup>,  
8 Eric Westman<sup>11</sup>, Andre Marquand<sup>12</sup>, Geir Selbæk<sup>3,4</sup>, Ole A. Andreassen<sup>2,13</sup>, Thomas  
9 Wolfers<sup>1,2,9,10,†</sup>, Lars T. Westlye<sup>1,2,13,†</sup>, Yunpeng Wang<sup>1,†</sup> for the Alzheimer's Disease  
10 Neuroimaging Initiative<sup>‡</sup> and the Australian Imaging Biomarkers and Lifestyle flagship study of  
11 ageing<sup>§</sup>

12

13 **Affiliations**

- 14 1. Department of Psychology, University of Oslo, Oslo, Norway  
15 2. Norwegian Centre for Mental Disorders Research (NORMENT), Oslo University Hospital &  
16 Institute of Clinical Medicine, University of Oslo, Oslo, Norway  
17 3. The Norwegian National Centre for Ageing and Health, Vestfold Hospital Trust, Norway  
18 4. Department of Geriatric Medicine, Oslo University Hospital, Oslo, Norway  
19 5. Computational Radiology & Artificial Intelligence (CRAI) Unit, Division of Radiology  
20 and Nuclear Medicine, Oslo University Hospital, Oslo, Norway  
21 6. Oxford Machine Learning in NeuroImaging (OMNI) Lab, University of Oxford, UK  
22 7. Institute of Clinical Medicine, University of Oslo, Oslo, Norway  
23 8. Division of Radiology and Nuclear Medicine, Oslo University Hospital, Oslo, Norway  
24 9. Department of Psychiatry and Psychotherapy, Tübingen Center for Mental Health, University  
25 of Tübingen, Germany  
26 10. German Center for Mental Health (DZPG)  
27 11. Division of Clinical Geriatrics, Department of Neurobiology, Care Sciences, and Society,  
28 Karolinska Institutet, Stockholm, Sweden  
29 12. Donders Institute for Brain, Cognition and Behaviour, Radboud University Medical Centre,  
30 Nijmegen, Netherlands  
31 13. KG Jebsen Center for Neurodevelopmental Disorders, University of Oslo, Oslo, Norway

32

33 † These authors contributed equally

34 ‡ Data used in preparation of this article were obtained from the Alzheimer’s Disease  
35 Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). As such, the investigators within  
36 the ADNI contributed to the design and implementation of ADNI and/or provided data but did  
37 not participate in analysis or writing of this report. A complete listing of ADNI investigators can  
38 be found at: [http://adni.loni.usc.edu/wp-](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)

39 [content/uploads/how to apply/ADNI Acknowledgement List.pdf](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)

40 § Data used in the preparation of this article was obtained from the Australian Imaging  
41 Biomarkers and Lifestyle flagship study of ageing (AIBL) funded by the Commonwealth  
42 Scientific and Industrial Research Organisation (CSIRO) which was made available at the ADNI  
43 database ([www.loni.usc.edu/ADNI](http://www.loni.usc.edu/ADNI)). The AIBL researchers contributed data but did not  
44 participate in analysis or writing of this report. AIBL researchers are listed at [www.aibl.csiro.au](http://www.aibl.csiro.au).

45

46 \* Corresponding author:

47 Esten H. Leonardsen,

48 Forskningsveien 3A, Harald Schjelderups hus, 0373 Oslo, Norway,

49 [estenhl@uio.no](mailto:estenhl@uio.no)

50

## 51 Abstract

52 Deep learning approaches for clinical predictions based on magnetic resonance imaging data  
53 have shown great promise as a translational technology for diagnosis and prognosis in  
54 neurological disorders, but its clinical impact has been limited. This is partially attributed to the  
55 opaqueness of deep learning models, causing insufficient understanding of what underlies their  
56 decisions. To overcome this, we trained convolutional neural networks on structural brain scans  
57 to differentiate dementia patients from healthy controls, and applied layerwise relevance  
58 propagation to procure individual-level explanations of the model predictions. Through extensive  
59 validations we demonstrate that deviations recognized by the model corroborate existing  
60 knowledge of structural brain aberrations in dementia. By employing the explainable dementia  
61 classifier in a longitudinal dataset of patients with mild cognitive impairment, we show that the  
62 spatially rich explanations complement the model prediction when forecasting transition to  
63 dementia and help characterize the biological manifestation of disease in the individual brain.  
64 Overall, our work exemplifies the clinical potential of explainable artificial intelligence in  
65 precision medicine.

## 66 Introduction

67 Since its invention in the 1970s, magnetic resonance imaging (MRI) has provided an opportunity  
68 to non-invasively examine the inside of the body. In neuroscience, images acquired with MRI  
69 scanners have been used to identify how the brains of patients with various neurological  
70 disorders differ from their healthy counterparts. Stereotypically, this has been done by collecting  
71 data from a group of patients with a given disorder and a comparable group of healthy controls,  
72 on which traditional statistical inference is applied to identify spatial locations of the brain where  
73 the groups differ <sup>1</sup>. Typically, these locations are not atomic locations identified by spatial  
74 coordinates, but rather morphological regions defined by an atlas, derived from empirical or  
75 theoretical insights of how the brain is structured. Differences between groups are described  
76 using morphometric properties like thickness or volume of these prespecified regions. A major  
77 benefit of this approach is the innate interpretability of the results: on average, patients with a  
78 given disorder deviate in a specific region of the brain in a comprehensible manner. Furthermore,  
79 the high degree of localization offered by modern brain scans allows for accurate  
80 characterization of where and how the brain of an individual deviates from an expected, typically

81 healthy, norm<sup>2</sup>. However, the effects which are found are typically small<sup>3</sup> with limited  
82 predictive power at the individual level<sup>4,5</sup>, which in turn has raised questions about whether  
83 these analytical methods are expressive enough to model complex mental or clinical phenomena  
84<sup>6</sup>. As an alternative, new conceptual approaches are proposed, advocating modelling frameworks  
85 with increased expressive power that allow for group differences through complex, non-linear  
86 interactions between multiple, potentially distant, parts of the brain<sup>7</sup>, with a focus on prediction  
87<sup>8</sup>. Such modelling flexibility is naturally achieved with artificial neural networks (ANNs), a class  
88 of statistical learning methods that combines aspects of data at multiple levels of abstraction, to  
89 accurately solve a predictive task<sup>9</sup>. However, while this often yields high predictive  
90 performance, e.g. by demonstrating clinically sufficient case-control classification accuracy for  
91 certain conditions, it comes at the cost of interpretation, as the models employ decision rules not  
92 trivially understandable by humans<sup>10</sup>. When the goal of the analysis is clinical, supporting the  
93 diagnosis and treatment of someone affected by a potential disorder, this opaqueness presents a  
94 substantial limitation. Thus, development and empirical validation of new methods within  
95 clinical neuroimaging that combine predictive efficacy with individual-level interpretability is  
96 imperative, to facilitate trust in how the system is working, and to accurately describe inter-  
97 individual heterogeneity.

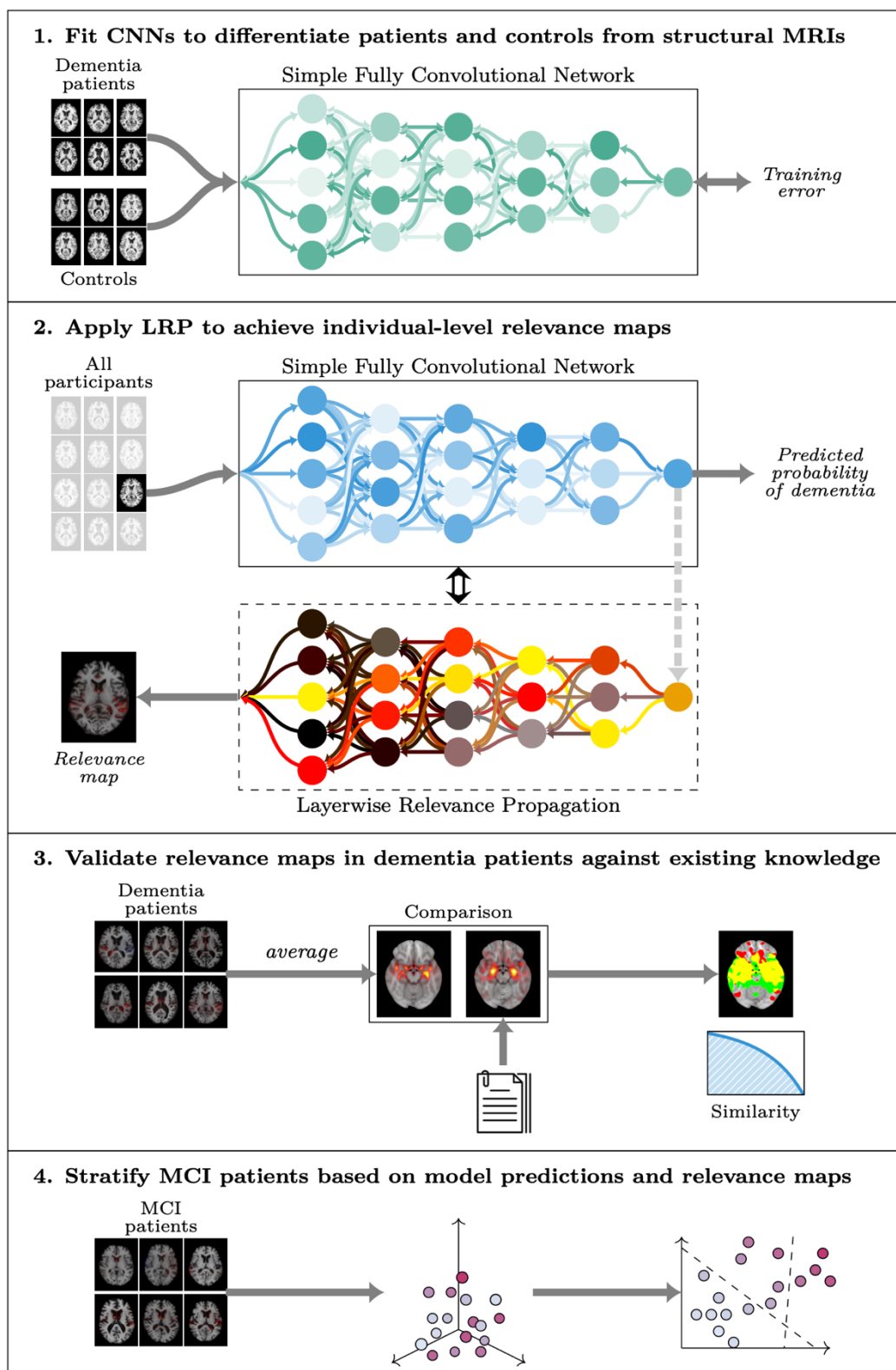
98  
99 With more than 55 million individuals afflicted worldwide<sup>11</sup>, over 25 million disability-adjusted  
100 life years lost<sup>12,13</sup> and a cost exceeding one trillion USD yearly<sup>14</sup>, dementia is a prime example  
101 of a neurological disorders that incur a monumental global burden. Due to the global aging  
102 population the prevalence is expected to nearly triple by 2050<sup>15</sup>, inciting a demand for  
103 technological solutions to facilitate handling the upcoming surge of patients. Dementia is a  
104 complex and progressive clinical condition<sup>16</sup> with multiple causal determinants and moderators.  
105 Alzheimer's disease (AD) is the most common form and accounts for 60%-80% of all cases<sup>11</sup>.  
106 However, the brain pathologies underlying different subtypes of dementia are not disjoint, but  
107 often co-occur<sup>17-19</sup>, and have neuropathological commonalities<sup>20</sup>. The most prominent is  
108 neurodegeneration, occurring in both specific regions like the hippocampus, and globally across  
109 the brain<sup>21</sup>, and inter-individual variations in the localization of atrophy has been associated with  
110 impairments in specific cognitive domains<sup>22,23</sup>. Thus, the biological manifestation of dementia in  
111 the brain is heterogeneous<sup>24</sup>, resulting in distinctive cognitive and functional deficits<sup>20</sup>,

112 highlighting the need for precise and personalized approaches to diagnosis. For patients with  
113 mild cognitive impairment (MCI), a potential clinical precursor to dementia, providing  
114 individualized characterizations of the underlying etiological disease at an early stage could  
115 widen the window for early interventions<sup>25</sup>, alleviate uncertainty about the condition, and help  
116 with planning for the future<sup>26</sup>.

117  
118 In dementia, ANNs, and particularly convolutional neural networks (CNNs), have been applied  
119 to brain MRIs to differentiate patients from controls<sup>27,28</sup>, prognosticate outcomes<sup>29</sup>, and  
120 differentially diagnose subtypes<sup>30</sup>. However, while research utilizing this technology has been  
121 influential, clinical translations are scarce<sup>31</sup>. Where techniques for segmenting brain tumours or  
122 detecting lesions typically produce segmentation masks that are innately interpretable, predicting  
123 a complex diagnosis would entail compressing all information contained in a high-dimensional  
124 brain scan into a single number. Using deep learning, the decisions underlying this immense  
125 reduction are obfuscated, both from the developer of the system, the clinical personnel using it,  
126 and the patient ultimately impacted by the decision. This black box nature is broadly credited for  
127 the low levels of adoption in safety-critical domains like medicine<sup>32</sup>. Responding to this  
128 limitation, explainable artificial intelligence (XAI) provides methodology to explain the  
129 behaviour of ANNs<sup>33</sup>. The nature of these explanations varies, e.g. by what type of model is to  
130 be explained, what conceptual level the explanation is at, and who it is tailored for<sup>34,35</sup>. In  
131 computer vision, XAI typically aims for post-hoc explanations of individual decisions,  
132 explaining why a model arrived at a given prediction for a given image. Explanations are often  
133 provided in a visual format, as a heatmap indicating how different regions of the image  
134 contribute to the prediction<sup>36</sup>. Layerwise Relevance Propagation (LRP) is a variant of such a  
135 method, based on propagating relevance from the prediction-space, backwards through all layers  
136 of the model to the image-space, to form a relevance map<sup>37</sup>. A major advantage of LRP is its  
137 intuitive interpretation: by construction, the total amount of relevance which denotes contribution  
138 to the prediction is kept fixed between layers. Thus, the relevance propagated back to an input  
139 voxel is directly indicative of the influence of that exact voxel to the prediction. Recently,  
140 several studies have applied both LRP and other explainable AI methods to dementia<sup>38</sup>, finding  
141 that the heatmaps generally highlight regions known to change in dementia<sup>39-42</sup>. However, the  
142 possibility of utilizing the fine-grained, individual, heatmaps produced by LRP to accurately

143 characterize individualized disease manifestations has not been explored, despite its potential for  
144 supporting clinical decisions towards precision medicine <sup>41,43</sup>.

145 In the present study, we applied techniques from deep learning and XAI on MRI scans of the  
146 brain to make explainable and clinically relevant predictions for dementia at the individual level  
147 (Figure 1). Using a state-of-the-art architecture for neuroimaging data, we trained CNNs to  
148 differentiate patients diagnosed with dementia from healthy controls based on T1-weighted  
149 structural MRIs. We implemented LRP on top of the trained models to form a computational  
150 pipeline producing individual-level explanations in the form of relevance maps alongside the  
151 model predictions. The relevance maps were validated in a subset of dementia patients, both in a  
152 qualitative comparison with existing knowledge of the anatomical distribution of structural  
153 aberrations, and in a quantitative, predictive context. Next, we applied the pipeline to a large,  
154 longitudinal dataset of MCI patients to create individual morphological records, a proposed data  
155 format for tracking and visualizing disease progression. Finally, we investigated the clinical  
156 utility of these records for stratifying patients, both in terms of their specific clinical profile, and  
157 progression of the disease. To facilitate reproducibility and improve the translational value of our  
158 work, the trained models and the complete explainable pipeline is made accessible on GitHub.



159

160 **Figure 1: Overview of the modelling process.** The modelling process consisted of four  
161 sequential steps. First, we fit multiple Simple Fully Convolutional Networks to classify dementia

162 *patients and healthy controls based on structural MRIs. Then we applied the best models to*  
163 *generate out-of-sample predictions and relevance maps for all participants. Next, we validated*  
164 *the relevance maps against existing knowledge using a meta-analysis to generate a statistical*  
165 *reference map. Finally, we employed the full pipeline in an exploratory analysis to stratify*  
166 *patients with mild cognitive impairment (MCI).*

## 167 Results

168 We compiled MRI data from multiple sources (Supplementary Table 1) into a dataset of  
169 heterogeneous dementia patients (n=854, age range=47-95, 47% females, Table 1) based on  
170 various diagnoses (Probable AD, vascular dementia, other/unspecified dementia) and diagnostic  
171 criteria for inclusion (Supplementary Table 2), and a set of controls strictly matched on site, age,  
172 and sex of equal size. We trained multiple CNNs to differentiate between the groups, employing  
173 a cross-validation approach utilizing all available timepoints for participants in three training  
174 folds and a single randomly selected timepoint for participants in separate validation and test  
175 folds. When stacking the out-of-sample predictions for all participants from all folds together  
176 (n=1708), for each fold using the model with the best validation performance, we observed  
177 satisfactory discrimination with a combined area under the receiver operating characteristics  
178 curve (AUC) of 0.908 (0.904-0.920 split across folds, Supplementary Figure 1), and an accuracy  
179 of 84.95% (83.04%-87.13%, Supplementary Table 3). This is slightly below with what is  
180 commonly achieved in similar studies classifying a specific subtype (typically AD) in a single  
181 dataset<sup>28</sup>.

182

CNN training and cross-validation			
Cohort	Participants	Mean age ( $\pm$ std)	Sex (F/M)
Healthy controls	854	75.13 $\pm$ 7.81	401/453
Dementia patients	854	74.82 $\pm$ 7.84	401/453
<b>Total</b>	<b>1708</b>	<b>74.98<math>\pm</math>7.82</b>	<b>802/906</b>
Downstream prognostic and correlational analyses			
Improved MCI	80	71.18 $\pm$ 8.14	37/43
Stable MCI	754	74.63 $\pm$ <b>7.66</b>	324/430
Progressive MCI	304	75.60 $\pm$ <b>7.46</b>	124/180
<b>Total</b>	<b>1138</b>	<b>74.67<math>\pm</math>7.73</b>	<b>485/653</b>



183 **Table 1: Cohorts.** Key characteristics of the cohorts used for training and testing the models,  
184 and further exploratory analyses.

185

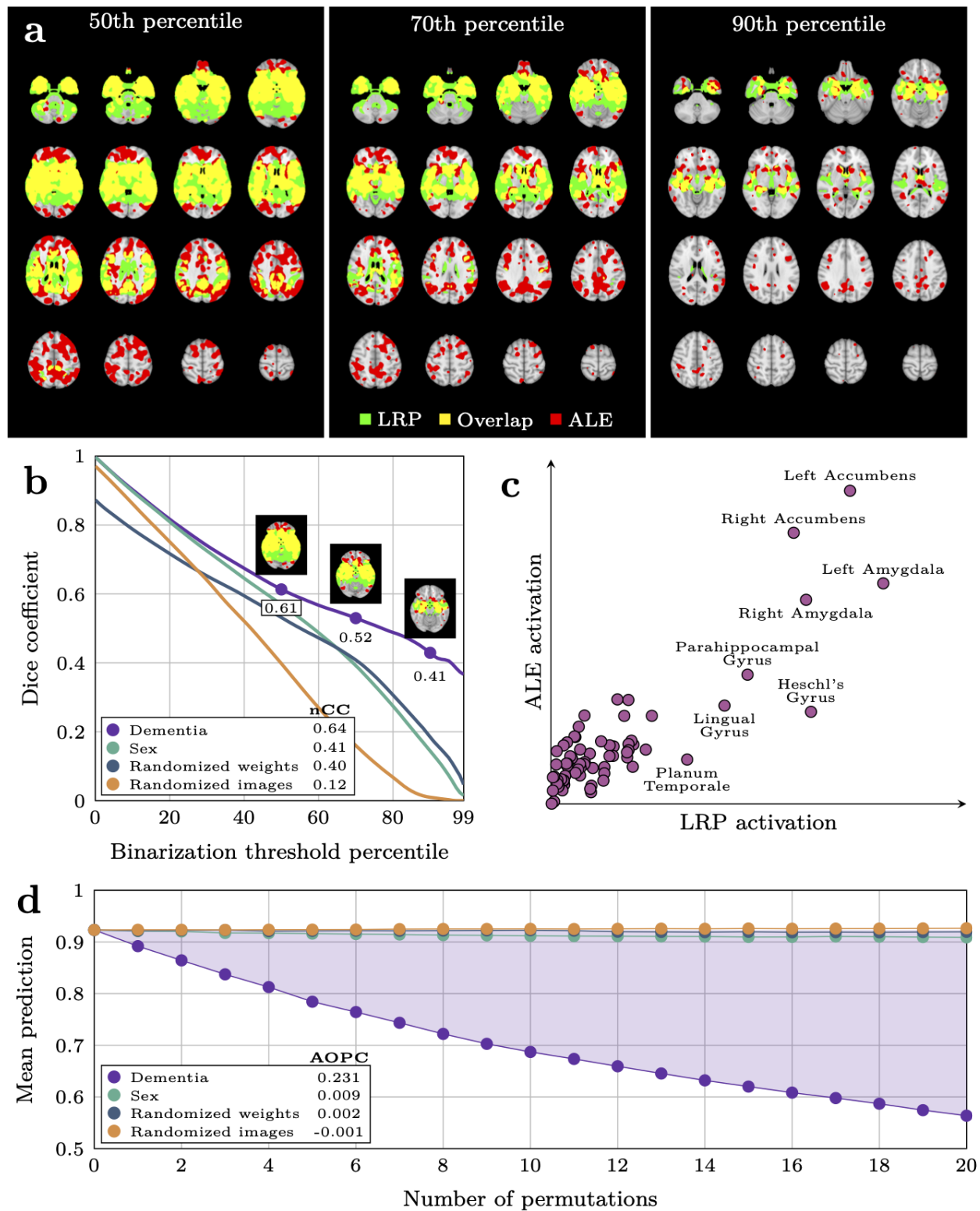
## 186 Relevance maps highlight predictive brain regions in individuals with dementia

187 Based on the classifiers with the highest AUCs in the validation sets, we built an explainable  
188 pipeline for dementia prediction,  $LRP_{dementia}$ , using composite LRP<sup>44</sup>, and a strategy to  
189 prioritize regions of the brain that contributed positively towards a prediction of dementia in the  
190 explanations. Using this pipeline, we computed out-of-sample relevance maps for all participants  
191 by applying the model for which the participant was unseen. Qualitatively, these maps  
192 corroborated known anatomical locations with structural aberrations in dementia, while still  
193 allowing for inter-individual variation (Supplementary Figure 2). We confirmed this apparent  
194 corroboration quantitatively by comparing a voxel-wise average map  $\bar{R}_{dementia}$  (Supplementary  
195 Figure 3), containing positive relevance from all correctly predicted dementia patients, with a  
196 statistical reference map  $G$  (Supplementary Figure 4) from an activation likelihood estimation  
197 (ALE) meta-analysis<sup>45</sup>, methodology established by an earlier study<sup>40</sup>. For sanity checks, we  
198 also computed average maps from three alternative pipelines,  $\bar{R}_{sex}$ ,  $\bar{R}_{randomized\ weights}$  and  
199  $\bar{R}_{randomized\ images}$ . The comparisons with the reference map were done by binarizing the maps  
200 on both sides of the comparison at various thresholds and measuring the Dice overlap (Figure  
201 2a). For the three alternative pipelines the amount of overlap decreased monotonically as the  
202 binarization threshold rose (Figure 2b), whereas for  $\bar{R}_{dementia}$  it stabilized as the maps grew  
203 sparser, indicating its higher similarity with  $G$ . This effect was reaffirmed by a normalized cross-  
204 correlation<sup>46</sup> of 0.64 for  $\bar{R}_{dementia}$ , compared to 0.41, 0.40 and 0.12 of  $\bar{R}_{sex}$ ,  
205  $\bar{R}_{randomized\ weights}$  and  $\bar{R}_{randomized\ images}$  respectively. In addition, we performed a region-  
206 wise, qualitative comparison of  $\bar{R}_{dementia}$  and  $G$ , also yielding general agreement (Figure 2c),  
207 with the most important regions in both maps being the nucleus accumbens, the amygdala, and  
208 the parahippocampal gyrus. Next, we tested the importance of the detected regions in a  
209 predictive context, by applying an iterative mask-and-predict procedure. For each participant, we  
210 produced a baseline dementia-prediction  $\hat{y}_0$  and relevance map  $R_{task}$  for each pipeline  $LRP_{task}$ .  
211 We then iteratively masked out the most important regions of the image according to the  
212 relevance map and recorded how the prediction changed as a function of the occlusion (Figure

213 2d). Using only true positives, the predictions should ideally start out at approximately 1.0  
214 (empirically found to be 0.89 on average) and trend towards 0.5 (random prediction) as a larger  
215 proportion of the image is occluded. The rate of decline is indicative of whether the masked  
216 regions contain information essential for the classifier to classify the image correctly. Over 20  
217 iterations we observed that the predictions based on maps from both  $LRP_{dementia}$ ,  $LRP_{sex}$  and  
218  $LRP_{randomized\ weights}$  decreased, but  $LRP_{dementia}$  at a distinctly steeper rate than the rest  
219 (Figure 2d). To quantify this observation we calculated an area over the perturbation curve  
220 (AOPC) of 0.231, 0.009, -0.001 and 0.002 for  $LRP_{dementia}$ ,  $LRP_{sex}$ ,  $LRP_{randomized\ images}$ ,  
221  $LRP_{randomized\ weights}$  respectively. Taken together, these results demonstrate that our pipeline  
222 generates maps with relevance in brain regions associated with changes in dementia.

223

224



225

226 *Figure 2: Validation of relevance maps from the dementia pipeline compared with three*  
 227 *alternative pipelines. a Visualization of the comparison between the binarized average relevance*

228 map  $\bar{\mathbf{R}}_{\text{dementia}}$  from the dementia-pipeline and the binarized statistical reference map  $\mathbf{G}$  from  
229 GingerALE, at different thresholds for binarization. **b** Overlap between the four average  
230 relevance maps  $\bar{\mathbf{R}}$  from our four pipelines and  $\mathbf{G}$  as a function of the binarization threshold. The  
231 numbers in the legend denote the normalized Cross Correlation (nCC) for each pipeline **c** Mean  
232 voxel-wise activation in  $\bar{\mathbf{R}}_{\text{dementia}}$  and  $\mathbf{G}$ , grouped by brain region. **d** Average participant-wise  
233 prediction from the dementia model after iteratively masking out regions of the image according  
234 to relevance maps from the four pipelines. Area over the permutation curve (AOPC) for the  
235 dementia map is indicated by the shaded area and denoted in the legend for all pipelines.

236

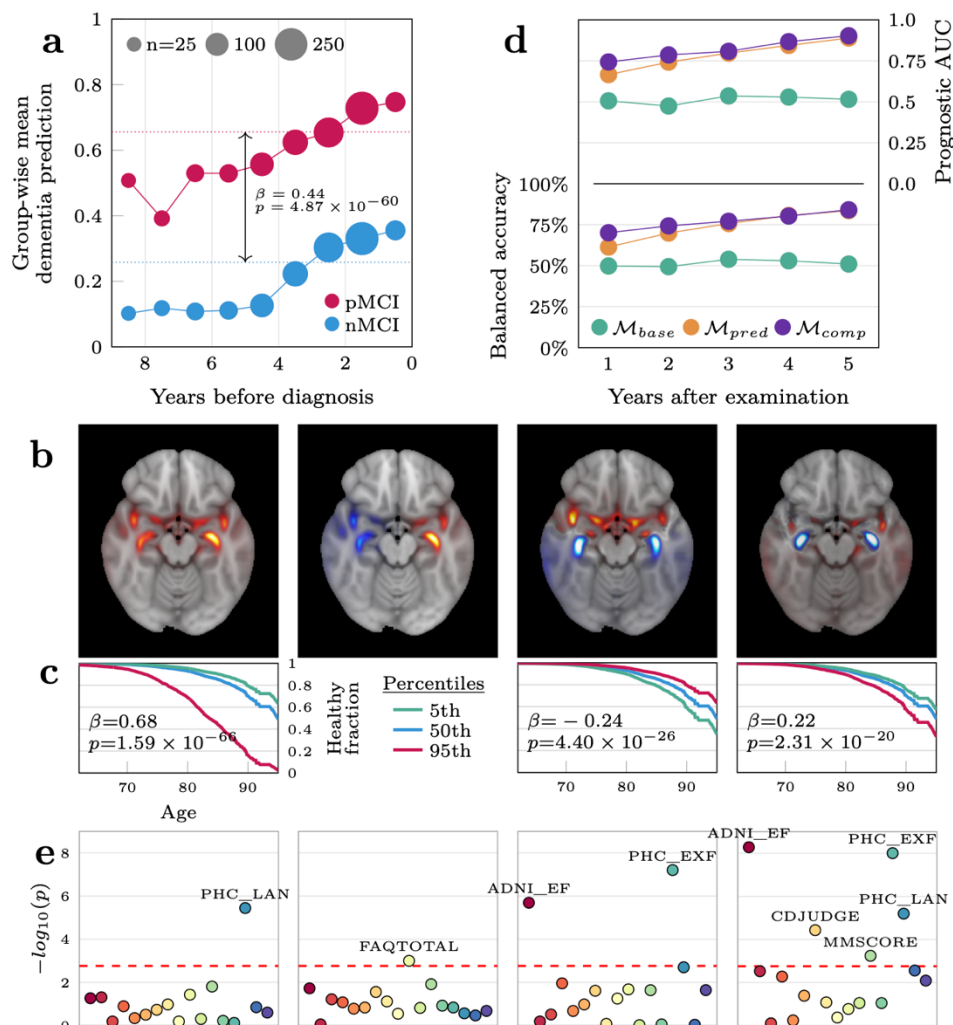
## 237 Output from the explainable dementia pipeline has prognostic value for MCI 238 patients

239 For the MCI patients (n=1256, timepoints=6448), previously unseen by all models, we built an  
240 averaging ensemble to procure a singular out-of-sample prediction and relevance map per patient  
241 per timepoint. Put together, we let this represent a morphological record (illustrated in Figure 4)  
242 visualizing the absolute quantity (indicated by the prediction) and location (indicated by the  
243 relevance map) of dementia-related pathology detected by the models over time. Qualitatively,  
244 both predictions and maps were relatively stable within a participant over time, while allowing  
245 enough variation to compose what resembled a trajectory. To investigate the prognostic value of  
246 our proposed morphological records we divided the MCI patients into three subgroups based on  
247 their trajectories in the follow-up period: those who saw improvement of their condition (n=80),  
248 those with a stable diagnosis throughout (sMCI, n=754), and those who progressed into dementia  
249 (pMCI, n=304). The remaining (n=118) had either a non-MCI diagnosis at the first timepoint, or  
250 a more complex diagnostic trajectory (e.g MCI -> AD -> CN) and were excluded from  
251 subsequent analyses. We observed that the predictions in the first group were generally very low  
252 (mean  $\hat{y} = 0.13$ , Supplementary Figure 5a), indicating that the models detected little, if any,  
253 evidence of dementia in these participants. For the stable patients the mean prediction was higher  
254 (mean  $\hat{y} = 0.33$ ), but still below the classification threshold of 0.5, whereas in the progressive  
255 group the model predicted the average patient to already have dementia (mean  $\hat{y} = 0.72$ ).  
256 Importantly, this was also true when considering only timepoints before these patients received  
257 the clinical diagnosis (mean  $\hat{y} = 0.65$ , Supplementary Figure 5b), suggesting that the model

258 found evidence of the disorder before the clinical symptoms surpassed the diagnostic threshold.  
259 To formally delineate the differences in predictions leading up to the potential diagnosis, we  
260 combined the improving and stable patients into a non-progressive group (nMCI, n=834), and  
261 sampled patients to match the progressive group based on their visiting histories, leading up to a  
262 terminal diagnosis timepoint (or a constructed non-diagnosis timepoint in the non-progressive  
263 group). In this matched dataset (n=550) we applied a linear mixed model controlling for age and  
264 sex and observed that the group difference was even greater than what we previously observed ( $\beta$   
265 = 0.47,  $p = 6.05 \times 10^{-71}$ , Figure 3a, Supplementary Table 4). Furthermore, we observed a  
266 significant difference in longitudinal slopes ( $\beta = 0.05$  increase in prediction per year,  $p =$   
267  $8.14 \times 10^{-17}$ ) indicating a greater rate of brain change detected by the model in those who  
268 would be diagnosed with dementia at a later point in time.

269  
270 The large group differences in the dementia predictions leading up to a potential diagnosis  
271 suggests this as a biomarker with innate prognostic value, yet the most salient part of our  
272 morphological records were the relevance maps. Thus, we performed exploratory analyses based  
273 on these to further differentiate the non-progressive and progressive groups and characterize both  
274 inter- and intra-group heterogeneity. However, given the high dimensionality of the maps and the  
275 relatively small number of patients, we first applied a principal component analysis (PCA) to  
276 relevance maps from all MCI patients, effectively compressing their information content into a  
277 smaller set of characteristic variables encoding facets of the maps, enabling the subsequent  
278 analyses. We retained the 64 components that explained the largest amount of variance and  
279 observed that they qualitatively clustered into three overarching categories. The first component  
280 was a generic component detecting general presence of relevance, resembling the average map  
281 from dementia patients, and thus made up a cluster by itself. The next cluster was comprised of  
282 the subsequent three components that captured high level, abstract patterns of relevance, namely  
283 differences in lateralization, along the sagittal axis and in subcortical regions (Figure 3b). The  
284 final cluster consisted of the remaining 60 components that captured specific, intricate patterns of  
285 presence/non-presence of relevance in regions revealed in the preceding analyses  
286 (Supplementary Figure 6). To investigate the potential of using the relevance maps for prognosis,  
287 we first performed a survival analysis using a Cox proportional hazards model where getting a  
288 diagnosis was considered the terminal event.

289



290

291 **Figure 3: Utility of the dementia pipeline for predicting progression and characterizing**  
 292 **individual-level deviations in the mild cognitive impairment cohort. a** Group-wise mean  
 293 predictions from the dementia-model in the progressive and non-progressive groups in the years  
 294 before a diagnosis was given. **b** The four first voxel-wise components of the principal component  
 295 analysis plotted in MNI152-space. **c** Survival curves for the average MCI patient (blue) and  
 296 fictitious patients at the extreme percentiles of the span for each component. The second  
 297 component was not significant and is not shown. **d** Predictive performance of the three models  
 298 predicting progression in the years following the MRI examination. The baseline model ( $M_{base}$ )  
 299 included only sex and age as covariates, the next model ( $M_{pred}$ ) included the prediction from the  
 300 dementia classifier as a predictor, while the final model ( $M_{comp}$ ) also added the component  
 301 vectors representing the relevance maps. **e** Significance levels of correlations between the each  
 302 of the four PCA components and various cognitive measures. The six annotated measures are

303 *composite language (PHC\_LAN) and executive function (PHC\_EXF) scores from the ADSP*  
304 *Phenotype Harmonization Consortium, total score from the Functional Activities Questionnaire*  
305 *(FAQTOTAL), composite executive function score from UW – Neuropsych Summary Scores*  
306 *(ADNI\_EF), clinical evaluation of impairment related to judgement and problem solving*  
307 *(CDJUDGE) from the Clinical Dementia Rating, and an overall measure of cognition from the*  
308 *Mini-Mental State Examination (MMSCORE, commonly referred to as MMSE).*

309 Specifically, we modelled the fraction of the population without a diagnosis as a function of age  
310 and used the subject-wise loadings of  $c_t$  as predictors. After Benjamini-Hochberg correction, 37  
311 of these components were significantly associated with staying undiagnosed (Figure 3c and  
312 Supplementary Table 5). However, we observed a correlation between the singular dementia  
313 prediction  $\hat{y}$  and the absolute magnitudes of these components (Supplementary Figure 7),  
314 indicating that the associations in the survival analysis could be induced by differences in the  
315 prediction rather than variability in the relevance maps. To mitigate this concern, we fit an  
316 equivalent model while stratifying on  $\hat{y}$ , observing that 29 associations remained significant, and  
317 that all coefficients had the same sign. Nonetheless, this analysis did not account for the  
318 predictions and relevance maps changing within a participant over time, so we reframed the  
319 question in a purely predictive setting, constructed to bear resemblance to a clinical scenario,  
320 using the same participants (nMCI=834, pMCI=304, total n=1138). For each MCI patient  $p$  at  
321 each timepoint  $t$  we asked whether we were able to predict, at yearly intervals  $\gamma$  up to five years  
322 into the future, whether  $p$  had progressed into dementia, using information from  $LRP_{dementia}$   
323 available at  $t$ . Importantly, all timepoints for all these participants were unseen by the dementia-  
324 model, yielding out of sample predictions and relevance maps from  $LRP_{dementia}$ , and we  
325 employed nested cross-validation to ensure the progression predictions were also out-of-sample.  
326 First, we fit a baseline model  $\mathcal{M}_{base}$  with age and sex as predictors, showing no predictive  
327 efficacy at any timepoint (all AUCs  $\approx 0.5$ , Supplementary Table 6), indicating that the dataset  
328 was not biased with respect to these variables. When adding the prediction from the dementia  
329 model  $\hat{y}_t$  as a predictor in model  $\mathcal{M}_{pred}$  we saw large improvements in prognostic efficacy at all  
330 yearly intervals, culminating with a fold-wise mean AUC of 0.889 after five years (Figure 3d). In  
331 the final model,  $\mathcal{M}_{comp}$ , also including the component vector  $c_t$  as predictors, we saw further  
332 improvements for all years, peaking at 0.903 after five years ( $p = 0.035$  when compared to

333  $\mathcal{M}_{pred}$  in a Wilcoxon signed-rank test across the outer folds). Overall, our best performing  
334 model predicted progression to dementia after five years with an AUC of 0.903, an accuracy of  
335 84.1%, a positive predicted value of 0.92, a sensitivity of 0.82 and a specificity of 0.86 (Table 2).  
336

Model	AUC	Balanced accuracy	PPV	Sensitivity	Specificity
$\mathcal{M}_{base}$	0.515	51.05%	0.14	0.09	0.93
$\mathcal{M}_{pred}$	0.889	83.61%	0.91	0.83	0.84
$\mathcal{M}_{comp}$	0.903	84.1%	0.92	0.82	0.86

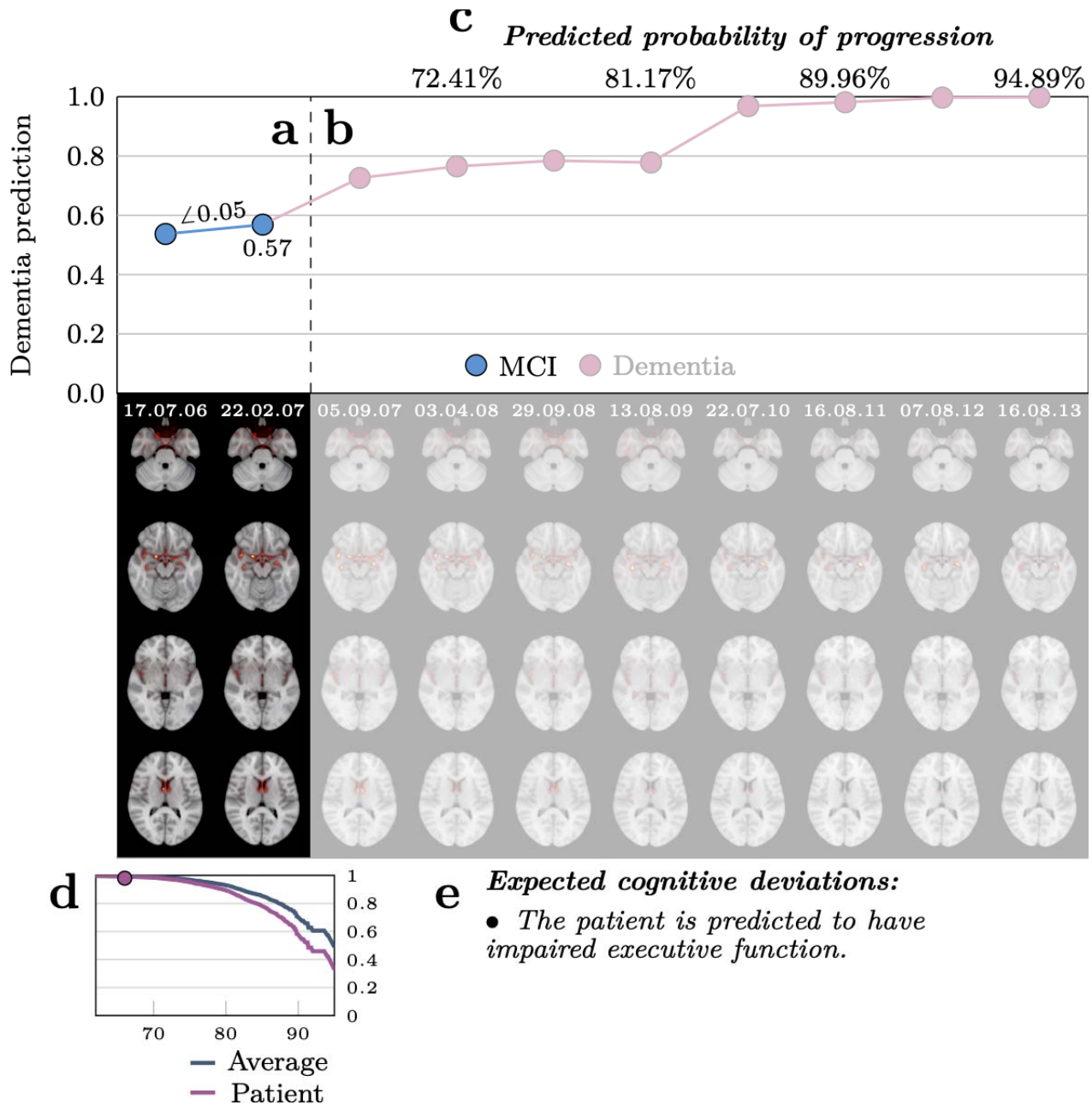
337 **Table 2: Predictive performance of the three models predicting progression five years into the**  
338 **future.** The baseline model  $\mathcal{M}_{base}$  used only age and sex as covariates.  $\mathcal{M}_{pred}$  also added the  
339 prediction from the dementia model at the current timepoint as a predictor, while  $\mathcal{M}_{comp}$   
340 additionally included the component vector  $c_t$  encoding information from the relevance maps.  
341

## 342 Facets of the relevance maps are associated with cognitive impairments in 343 distinct domains

344 Finally, we tested whether common features found in the relevance maps, represented by the  
345 PCA component, were correlated with impairments in distinct cognitive and functional domains.  
346 We extracted 17 summary measures from 7 neuropsychological tests (Supplementary Table 7  
347 and 8), performed approximately at the same time as an MRI examination, and tested for  
348 associations with the subject-wise loadings of  $c_t$  in 733 MCI patients using linear models. After  
349 FDR correction, while correcting for age, sex and  $\hat{y}$ , we found 48 significant correlations  
350 between 18 unique components and 14 of the cognitive measures (Figure 3e). Component 30 and  
351 the aggregate score from the Functional Activities Questionnaire (FAQTOTAL) had the highest  
352 number of significant hits among the components and measures respectively, both with six  
353 passing the threshold. Most importantly, the components showed distinct patterns of associations  
354 with the different cognitive measures. To ensure the significant associations were not driven by  
355 collinearity between components  $c_i$  and  $\hat{y}$ , we ran an equivalent analysis without including  $\hat{y}$  as a  
356 predictor, observing that only 5/48 of the previously significant hits had coefficients with the  
357 opposite sign. To summarize, the spatial features captured in our relevance maps, and  
358 subsequently in our component vectors, were associated with distinct patterns of performance on



359 neuropsychological tests relevant for characterizing phenotypic heterogeneity in dementia  
360 patients (Supplementary Figure 8).  
361



362

363 **Figure 4: A visualization of the proposed morphological record for a randomly selected**

364 **progressive MCI patient that was held out of all models and analyses. a** The top half shows the

365 prediction from the dementia model at each visit, while the bottom part displays the relevance

366 map underlying the prediction. The opaque sections (including c, d, and e) contain information

367 accessible at the imagined current timepoint (22.02.07) to support a clinician in a diagnostic

368 procedure. The angle ( ) represents the change in dementia prediction per year based on the

369 first two visits. b Translucent regions reveal the morphological record for the remaining follow

370 ups in the dataset, thus depicting the future. The ground truth diagnostic trajectory is encoded by

371 *the colour of the markers. c Predicted probabilities of progression at future follow-ups based on*  
372 *the prediction and relevance map at the current timepoint. d Survival curve of the patient*  
373 *compared to the average MCI patient calculated from the prediction and relevance map. The*  
374 *marker indicates the location of the patient at the current timepoint. e A list of cognitive domains*  
375 *where the patient is predicted to significantly differ from the average based on the prediction*  
376 *and relevance map.*  
377

## 378 Discussion

379 Given the huge burden of disease and expected increase in prevalence, innovative technological  
380 solutions for clinical decision making in dementia diagnostics and prognostics is urgently  
381 needed. Although commonly referred to as a homogenous condition or split into a few subtypes  
382 based on aetiology or pathophysiology<sup>17</sup>, dementia patients exhibit unique and complex  
383 deficiencies, disease trajectories, and cognitive deficits. To explore the potential of brain MRI  
384 and XAI to characterize heterogeneity in the brain underpinnings of dementia, we trained neural  
385 networks to differentiate dementia patients from healthy individuals, and derived relevance maps  
386 using Layerwise Relevance Propagation to explain the individual-level decisions of the classifier.  
387 The relevance maps were specific to the individual, spanned regions that were predictive of  
388 dementia and corroborated existing knowledge of the anatomical distribution of structural  
389 aberrations. In a cohort of MCI patients, it enabled characterization and differentiation of  
390 individual-level disease manifestations and trajectories linked to cognitive performance in  
391 multiple domains. While further validations in clinical contexts are needed, our XAI pipeline for  
392 dementia demonstrates how advanced predictive technology can be employed by clinicians to  
393 monitor and characterize disease development for individual patients.

394  
395 There is a multitude of XAI techniques available for explaining the decisions of an image  
396 classifier, many of which have yielded promising results for dementia classification<sup>38</sup>. We  
397 employed LRP due to its straightforward interpretation as well as earlier studies indicating  
398 robustness<sup>47</sup> and specificity<sup>42</sup>, properties we consider integral in a clinical decision support  
399 system. But while procuring explanations that are *ipso facto* meaningful is an important step  
400 towards adoption of AI in clinical neuroimaging, it is not in itself sufficient. There is a host of  
401 predictive models that are trivially explainable, but not understandable<sup>48</sup>, and there is genuine  
402 concern that XAI will lead to another level of systems that are formally well-defined, but opaque  
403 and obscure, and thus practically useless<sup>49</sup>. Thus, empirical explorations are imperative to  
404 investigate the nature of these explanations, examine how they may be useful and build essential  
405 trust<sup>50</sup>. In our validation, we observed that the explanatory maps produced by the dementia  
406 pipeline were more predictive and showed distinctly more agreement with existing knowledge of  
407 pathology than those produced by the three alternative pipelines. Given limitations that have  
408 been exposed in such methods earlier<sup>51,52</sup> these validations are crucial, and observing that our

409 results both corroborate earlier evidence<sup>40</sup> and extend upon it, provides confidence that the  
410 explanations derived from the model are meaningful. However, we emphasize that the ultimate  
411 validation should happen in actual implementations of the technology in end-user systems, with  
412 clinical personnel applying it in clinical scenarios on realistic data.

413  
414 We continued beyond validating the relevance maps by proposing them as a potential epistemic  
415 and clinical tool to characterize individual facets of dementia. To this end, we explored if the  
416 maps contributed to predicting imminent progression from MCI to dementia, and correlated them  
417 with different cognitive measures, extending upon the current literature<sup>38</sup>. In both analyses we  
418 found evidence, although modest, that the maps are informative beyond the predictions of the  
419 model. To illustrate the potential of the pipeline for clinical decision making we compiled its  
420 output into a proposed morphological record (visualized for a single patient in Figure 4) that can  
421 help clinicians localize morphological abnormalities during a diagnostic process. Identifying  
422 subtle pathophysiology through deep phenotyping could have a huge potential for charting the  
423 heterogeneity of dementia, providing precise biological targets to guide future research.  
424 Furthermore, for the individual patient, it can support personalized diagnosis to identify  
425 appropriate disease-modifying treatments, and in the future, hopefully, accurate therapeutic  
426 interventions.

427  
428 The regions with the highest density of relevance in our maps were the nucleus accumbens,  
429 amygdala and the parahippocampal gyrus, all of which are strongly affected in dementia<sup>53-55</sup>.  
430 While the two latter corroborate the established involvement of the medial temporal lobe<sup>56</sup> it is  
431 surprising that the hippocampus does not appear in our analyses, as it has frequently in similar  
432 studies<sup>38,41,42</sup>. While this could be caused by actual localization of pathology<sup>57</sup> we consider it  
433 more likely to be related to the internal machinery of the model. Specifically, the CNN relies on  
434 spatial context to identify brain regions before assessing their integrity, utilizing filters that span  
435 areas of the image larger than those containing the region itself. In the backwards pass, LRP uses  
436 these filters, and thus the localization of relevance is not necessarily voxel precise. Furthermore,  
437 we believe the model broadly can be seen as an atrophy detector, which necessarily entails  
438 looking for gaps surrounding regions instead of directly at the regions themselves. Therefore,

439 while the relevance maps provide important information, they depend on contextual information  
440 and thus rely on interpretation from clinicians to maximize their utility in clinical practice.

441  
442 We focused our analyses mainly on the relevance maps, but the results with largest, immediate,  
443 potential for clinical utility were the predictions from the dementia classifier. Other studies have  
444 shown the efficacy of machine learning models in differentiating dementia patients and healthy  
445 controls<sup>28</sup>, but it is intriguing that we see a large discrepancy in the predictions of the  
446 progressive and non-progressive MCI patients many years before the dementia diagnosis is  
447 given. This corroborates findings from theory-driven studies<sup>58</sup> and a recent deep learning study  
448<sup>27</sup>, implying detectable structural brain changes many years before the clinical diagnosis is given.  
449 This gives hope for advanced technology to contribute to early detection and diagnosis through  
450 MRI based risk scores, in our case supported by a visual explanation. If curative treatments  
451 prove efficacious and become accessible, early identification of eligible patients could be  
452 imperative<sup>59</sup>. Furthermore, timely access to interventions have shown efficiency in slowing the  
453 progress of cognitive decline<sup>60</sup>, in addition to improving the quality of life for those afflicted and  
454 their caregivers<sup>26,61</sup>. Widely accessible technology that allows for early detection with high  
455 precision could play a key role in the collective response to the impending surge of patients and  
456 provide an early window of opportunity for more effective treatments.

457  
458 While our results show a great potential for explainable AI, and particularly LRP, as a  
459 translational technology to detect and characterize dementia, there are limitations to our study.  
460 First, there are technical caveats to be aware of. Most importantly, there is an absolute  
461 dependence between the predictions of our model and the relevance maps. In our case, when we  
462 qualitatively assessed the relevance maps of the false negatives, they were indistinguishable from  
463 the true negatives. This emphasizes the fact that when the model is wrong, this is not evident  
464 from the explanations. Next, while the maps contain information sufficient to explain the  
465 prediction, they are not necessarily complete. Thus, they don't contain all evidence in the MRI  
466 pointing towards a diagnosis, a property which could prove essential for personalization. We  
467 have addressed this problem through pragmatic solutions, namely ensembling and targeted  
468 augmentations, but theoretical development of the core methodology might be necessary to  
469 theoretically guarantee complete maps. Beyond the fundamental aspects of LRP, there are

470 weaknesses to the present study that should be acknowledged. First, the dataset with dementia  
471 patients portrayed as heterogeneous mostly consists of ADNI and OASIS data, and thus patients  
472 with a probable AD diagnosis (although clinically determined). Thus, while we consider it likely,  
473 it is not necessarily true that the dimension of variability spanning from healthy controls to  
474 dementia patients portrayed by our model has the expressive power to extrapolate to other  
475 aetiologies. To overcome this in actual clinical implementations, we encourage the use of  
476 datasets that are organically collected from subsets of the population that are experiencing early  
477 cognitive impairments, for instance from memory clinics. Furthermore, it is not trivial to  
478 determine whether a clinical, broad, dementia-label is an ideal predictive target for models in  
479 clinical scenarios. Both ADNI and AIBL contain rich biomarker information with multiple  
480 variables known to be associated with dementia, such as amyloid positivity. It would be  
481 intriguing to see studies methodologically similar to ours with a biological predictive target, and  
482 we encourage investigations into whether this supports and complements the results we have  
483 observed here. Another limitation with the present study is out-of-sample generalization,  
484 especially related to scanners and acquisition protocols. Although we utilize data from many  
485 sites, which we have earlier shown to somewhat address this problem<sup>62</sup>, in combination with  
486 transfer learning, we did not explicitly test this by e.g., leaving sites out for validation. Again, we  
487 advise that clinical implementations should be based on realistic data, and thus at least be  
488 finetuned towards data coming from the relevant site, scanner, and protocol implemented in the  
489 clinic<sup>63</sup>. This also includes training models with class frequencies matching those observed in  
490 clinical settings, instead of naively balancing classes as we have done here. Next, we want to  
491 explicitly mention the cyclicity of our mask-and-predict validation. In a sense it trivially  
492 follows that regions that are considered important by a model are also the ones that are driving  
493 the predictions, and thus it is no surprise that the relevance maps coming from the dementia  
494 model are more important to the dementia model than the maps coming from e.g., the sex model.  
495 We addressed this by alternating the models for test and validation, but fully avoiding this  
496 circularity would require disjunct datasets, and more and larger cohorts. Finally, we highlight the  
497 potential drawbacks of including the improving MCI patients alongside the stable in the  
498 progression models. We believe this accurately depicts a realistic clinical scenario, where  
499 diagnostic and prognostic procedures happen based on currently available clinical information.  
500 However, that these patients improve could indicate that their condition is not caused by stable

501 biological aberrations. This could oversimplify the subsequent predictive task, inflating our  
502 performance measures. In summary, the predictive value we observed for the individual patient  
503 must be interpreted with caution. However, our extensive validation approach as well as our  
504 thorough explanation of the method and its limitations, and training on large datasets, provide a  
505 first step towards making explainable AI relevant for clinical decision support in neurological  
506 disorders. Nonetheless, it also reveals a complicated balance between validating against existing  
507 knowledge and allowing for new discoveries. In our case, confirming whether small details  
508 revealed in the relevance maps are important aspects of individualization or simply intra-  
509 individual noise requires datasets with a label-resolution beyond what currently exists. Thus, we  
510 reiterate our belief that the continuation of our work should happen at the intersection between  
511 clinical practice and research<sup>64</sup>, by continuously collecting and labelling data to develop and  
512 validate technology in realistic settings.

513

514 To conclude, while there are still challenges to overcome, our study provides an empirical  
515 foundation and a roadmap for implementations of brain MRI based explainable AI in  
516 personalized clinical decision support systems. Specifically, we show that deep neural networks  
517 trained on a heterogenous set of brain MRI scans can predict dementia, and that their predictions  
518 can be made human interpretable. Furthermore, our pipeline allows us to reason about structural  
519 brain aberrations in individuals showing early signs of cognitive impairment by providing  
520 personalized characterizations which can subsequently be used for precise phenotyping and  
521 prognosis, thus fulfilling a realistic clinical purpose.

## 522 **Materials and Methods**

### 523 **Data**

524 All data used in the present study have been obtained from previously published studies which  
525 have been approved by their respective institutional review board or relevant research ethics  
526 committee.

527

528 To train the dementia models we compiled a case-control dataset from seven different sources  
529 (Supplementary Table 1), consisting of patients with a dementia diagnosis and healthy controls  
530 from the same scanning sites. Because of the different diagnostic criteria used in the original



531 datasets we applied different rules to achieve a singular, heterogeneous dementia label  
532 (Supplementary Table 2). We extracted all participants with a dementia-diagnosis at all  
533 timepoints to comprise the patient group (n=854). Then, for each unique proxy site (In ADNI,  
534 due to the large number of scanners and acquisition protocols, and the work put into unifying  
535 them, we used field strength as a proxy for site), sex, and age-bin spanning 10 years, we sampled  
536 an equal number of healthy controls to form the matched control set (total n=1708, Table 1).  
537 Lastly, before modelling, we split the data into five equally sized folds stratified on diagnosis,  
538 site, sex, and age, such that all timepoints for a single participant resided in the same fold.

539  
540 For the MCI dataset we started with all participants from all ADNI waves with an MCI diagnosis  
541 (subjective memory complaint, MMSE between 24 and 30, CDR>0.5 with memory box>0.5,  
542 Weschler Memory Scale-Revised <9 for 16 years of education, <5 for 8-15 years of education  
543 and <3 for 0-7 years of education)<sup>65</sup>, on at least one timepoint. These were 12661 images from  
544 6448 visits for 1256 participants, none of which were used for model training. This selection  
545 criterion ensured all participants had an MCI diagnosis at one point in time, though it did not  
546 limit us to only those timepoints. Thus, in addition to those with a consistent, stable, MCI  
547 diagnosis (sMCI), we had a variety of diagnostic trajectories, including those transitioning from  
548 normal cognition to MCI, MCI to AD (pMCI) and various other combinations. Before the  
549 subsequent analyses we discarded all participants without an MCI diagnosis initially, and  
550 everyone with ambiguous trajectories (e.g. MCI->CN->AD), leaving 5607 visits from 1138  
551 participants.

552  
553 From these two datasets we extracted T1-weighted structural MRI data for each participant at  
554 each timepoint to use as inputs for the subsequent predictive models. Prior to modelling, the raw  
555 images were minimally processed using a previously developed pipeline<sup>2/21/2024 12:12:00 PM</sup>  
556 relying on FreeSurfer v5.3 and FSL v6.0<sup>66</sup> to perform skullstripping<sup>67</sup> and linear registration to  
557 MNI152-space<sup>68</sup> with six degrees of freedom. Consequently, the processed images consisted of  
558 normalized voxel values from the raw images, registered to a common spatial template and  
559 contained minimal non-brain tissue.

560

## 561 Modelling

562 All dementia models were variants of the PAC2019-winning simple fully convolutional network  
563 (SFCN) architecture<sup>69,70</sup>, modified to have a single output neuron with a sigmoid activation. The  
564 architecture is a simple, VGG-like convolutional neural network with 6 convolutional blocks and  
565 approximately 3 million parameters. We initialized the model with weights from a publicly  
566 accessible brain age model previously shown to have superior generalization capabilities when  
567 dealing with unseen scanning sites and protocols<sup>62</sup>. The models were trained on a single Nvidia  
568 A100 GPU with 40GB of memory, Tensorflow 2.6<sup>71</sup> through the Keras interface<sup>72</sup>. We used a  
569 vanilla stochastic gradient descent (SGD) optimizer with a learning rate defined by the  
570 hyperparameter settings (see next section), optimizing the binary cross-entropy loss. All models  
571 ran for 160 epochs with a batch size of 6, and for each run the epoch with the lowest validation  
572 loss was chosen. Varying slightly depending on the hyperparameters, a single model trained in  
573 approximately 4 hours.

574  
575 For each hold-out test fold we trained models on three of the remaining folds and validated on  
576 the fourth, akin to a cross-validation with an additional out-of-sample test set, to achieve out-of-  
577 sample predictions for all 1708 participants while allowing for hyperparameter tuning. The  
578 hyperparameters we optimized were dropout  $d \in \{0.25, 0.5\}$  and weight decay  $w \in$   
579  $\{10^{-2}, 10^{-3}\}$ . Additionally, we tested stepwise, one-cycle and multi-cycle learning rate  
580 schedules and a light and a heavy augmenter. Initial values for the learning rate were set  
581 manually based on a learning rate sweep<sup>73</sup>, though kept conservative to preserve the learned  
582 features from the pretraining. The hyperparameter search was implemented as a naive grid-  
583 search over the total 24 different configurations (Supplementary Figure 9). We selected the  
584 model procuring the best AUC in the validation set to produce out-of-sample predictions for the  
585 outer hold-out fold. In the final evaluation of the models, we compiled predictions for all  
586 participants, for each using the model where they belonged to the hold-out test set. Our main  
587 method for measuring performance was the AUC, but we also report accuracy, which, due to our  
588 matching procedure, is equivalent to balanced accuracy.

589

## 590 Relevance maps

591 We built a pipeline  $LRP_{dementia}$  for generating relevance maps by implementing LRP (Bach et  
592 al., 2015) on top of the trained classifier. LRP is a technique for explaining single decisions  
593 made by the model, and thus, when running the pipeline on input  $X$  a relevance map  $R$  is  
594 generated alongside the prediction  $\hat{y}$ .  $R$  is a three-dimensional volume, representing a visual  
595 explanation for  $\hat{y}$ , where each voxel  $r_{i,j,k} \in R$  has a spatial position  $i, j, k$  corresponding to the  
596 location of an input voxel  $x_{i,j,k} \in X$ . Furthermore, the intensity of  $r_{i,j,k}$  can be directly  
597 interpreted as how much voxel  $x_{i,j,k}$  contributes to  $\hat{y}$ , such that  $\sum_{r \in R} r = \hat{y}$ . In the original LRP-  
598 formulation, relevance  $r$  is propagated backwards between subsequent layers  $Z_l$  and  $Z_{l+1}$   
599 according to the relative contribution of one artificial neuron  $a_m \in Z_l$  in the first layer on  
600 relevance in all artificial neurons  $a_n \in Z_{l+1}$  in the following layer,

$$601$$
$$602 \quad r(a_m) = \sum_j \frac{a_m w_{mn}}{\sum_o a_o w_{on}} r(a_n),$$
$$603$$

604 where  $w_{mn}$  denotes the weight between  $a_m$  and  $a_n$ . We controlled the influence of different  
605 aspects of the explanations using a composite LRP strategy<sup>44</sup>, combining different formulations  
606 of the LRP-formula for the different layers in the model to enhance specific aspects of the  
607 relevance maps. Specifically, we employed a combination of alpha-beta and epsilon rules that  
608 have previously shown to produce meaningful results for dementia-classifiers<sup>41,42</sup>, described in  
609 detail in the Supplementary Methods. The resulting relevance maps produced by the pipeline  
610 were full brain volumes with the same dimensionality as the MRI data (167x212x160 voxels)  
611 containing mostly (see below) positive relevance.

612

613 Notation-wise we generally consider the relevance map  $R(X)$  for an image  $X$  to be a function of  
614 the model  $m_{task}$ , where  $task$  indicates which task the model was trained for, the LRP strategy  
615  $LRP_{composite}$  and the image  $X$ ,

$$616$$
$$617 \quad R(X) = f(m_{task}, LRP_{composite}, X).$$
$$618$$

619 Because the composite LRP strategy described above is kept fixed in our pipeline, we contract  
620 this to

$$621 \\ 622 R(X) = f(m_{task}, X). \\ 623$$

624 Furthermore, we let the model-specifier task annotate the map for a further simplification,

$$625 \\ 626 R_{task}(X) = f(X). \\ 627$$

628 Thus,  $LRP_{task}$  is used to annotate the full pipeline for a given task, while  $R_{task}(X)$  denotes a  
629 single relevance map generated by this pipeline for image  $X$ . When the task is given by the  
630 context, we sometimes simplify this further to  $R(X)$ , and when a general image is considered, we  
631 simply use  $R$  to denote its relevance map.

632  
633 While we generally discuss our pipeline as a singular one, there were in reality five  
634 approximately equivalent pipelines (corresponding to the models trained for the five test folds),  
635 and which one is used depends on what image was used as input. Specifically, for each  
636 participant diagnosed with dementia, the pipeline is chosen where the participant was part of the  
637 hold-out test set while training the model, and both the relevance maps and the predictions are  
638 thus always out-of-sample. For participants used in the MCI analysis, which are all out-of-  
639 sample for all models, we created an ensemble by averaging the predictions and the voxel-wise  
640 relevance across all models.

641  
642 Before implementing the LRP procedure we made two slight modifications to the models to  
643 facilitate the backwards relevance propagation, both leaving the functional interface of the model  
644 unchanged. First, we removed the sigmoid activation in the final layer, so that the output of the  
645 model changed from a bounded continuous variable  $\hat{y} \in [0, 1]$  to an unbounded prediction  
646  $\hat{y}_\sigma \in [-\infty, \infty]$ . In this space a raw prediction of  $\hat{y}_\sigma = 0$  is equivalent to a sigmoid-transformed  
647 prediction of  $\hat{y} = 0.5$ , and thus  $\hat{y}_\sigma < 0$  means that the model predicts control status for the  
648 given participant, and oppositely  $\hat{y}_\sigma > 0$  implies that the model predicts a dementia diagnosis.  
649 Furthermore, this means that all positive relevance  $r \in R$ ,  $r > 0$  can be interpreted as visual

650 evidence in favour of a dementia diagnosis. Secondly, we modified the model by fusing all batch  
651 normalization layers with their preceding convolutional layers, adjusting their weights and biases  
652 to match the shift and scaling previously performed by the normalization layer<sup>74,75</sup>.

653  
654 After generation, the relevance maps are in the same stereotaxic space as their corresponding,  
655 linearly registered, input MRIs. To ensure intra-individual comparisons were done in the same  
656 space we non-linearly registered the maps to MNI152- space before subsequent statistical  
657 analyses were run. First, we registered the preprocessed MRIs  $X$  used as inputs to the 1mm  
658 MNI152 template packaged with FSL using `fnirt` with `splineorder=2`. We then applied the  
659 transformation computed for  $X$  to  $R(X)$  using `applywarp`. We also restrained our relevance maps  
660 to contain strictly positive relevance, evidence in favour of a dementia prediction, by clipping  
661 them to a minimum value of 0. Furthermore, to remove edge-effects from our analyses, we  
662 enforce that there is no relevance in non-brain tissue by nullifying all relevance outside the brain:

663  
664 
$$\forall(i, j, k)[x_{i,j,k} = 0 \Rightarrow r_{i,j,k} = 0].$$

665  
666 All visualized relevance maps are plotted after non-linear registration, overlaid on the MNI152-  
667 template. As the maps are three-dimensional, we generally plot a collection of distributed axial  
668 slices. The relevance is coloured by the `nibabel v3.2.2`<sup>76</sup> `cold_hot` colourmap. Since the absolute  
669 relevance values vary between maps, all maps are normalized to the intensity range  $[0, 1]$  in the  
670 visualizations.

671  
672 **Validating the relevance maps**  
673 Earlier studies have shown that interpretability techniques in general are prone to generate visual  
674 explanations that do not capture salient parts of the input<sup>51,52</sup>. To investigate the extent of this for  
675 our pipeline  $LRP_{dementia}$  we employed two analyses to assess the sanity of the relevance maps.  
676 The first was an established task-specific technique comparing the relevance maps to existing  
677 knowledge of the pathology of dementia<sup>40</sup>. The second was a purely quantitative analysis  
678 examining how important the regions found by the pipeline are for the dementia prediction  $\hat{y}$ . In  
679 both cases we contrasted the relevance maps generated from the main pipeline with three

680 alternative pipelines representing variants of a null hypothesis, all expected to produce relevance  
681 maps with no significant association to dementia.

682  
683  $LRP_{random\ images}$  represents the simplest alternative pipeline, and is built around the dementia-  
684 model, but with an additional preprocessing step scrambling the input,

685  
686 
$$R_{random\ images}(X) = R_{dementia}(\mathcal{X}), \text{ where}$$
  
687 
$$\mathcal{X} = \mathcal{N}(\bar{X}, \sigma_X).$$

688  
689  $LRP_{random\ images}$  is expected to generate relevance maps where the relevance is evenly  
690 distributed across the entire image. In the next pipeline  $LRP_{random\ weights}$  we replaced the  
691 dementia-model with a model with random weights,

692  
693 
$$R_{random\ weights}(X) = R(m_\theta, X).$$

694  
695  $m_\theta$  has not been trained for any task, and thus has random weights initialized by the default  
696 Keras "Glorot Uniform" weight-initializer. This pipeline is expected to produce relevance maps  
697 which correlate with the raw voxel intensities, e.g. high intensity in the input should entail more  
698 (absolute) relevance, thereby reflecting aspects of morphology. The final and most realistic  
699 alternative pipeline was  $LRP_{sex}$ , where we replaced the dementia-model with a binary sex-  
700 classifier,

701  
702 
$$R_{sex}(X) = R(m_{sex}, X).$$

703  
704 The sex-classifier was trained to differentiate males from females in one of the splits from the  
705 dementia-dataset, achieving an out-of-sample AUC of 0.956 and a balanced accuracy of 89.40%.  
706 We did not do any hyperparameter optimization for this model but used the best configuration  
707 from the dementia cross-validation in the same fold. The heatmaps from this pipeline should  
708 reflect regions where there is intra-individual variation in morphology, which are predictive of  
709 sex but with minimal association with dementia.

710

711 As a proxy for existing knowledge in the literature we performed an ALE meta-analysis using  
712 Sleuth v3.0.4<sup>77</sup> and GingerALE v3.0.2<sup>45</sup>. We used Sleuth to search for relevant articles with the  
713 query

714  
715 Imaging Modality is MRI  
716 AND  
717 Context is disease  
718 AND  
719 Diagnosis is Dementia OR Alzheimer’s Disease OR Lewy Body Dementia OR Frontotemporal  
720 Dementia OR Non-Aphasic Frontotemporal Dementia

721  
722 in the Voxel-based morphometry database, yielding 394 experiments from 124 articles. These  
723 experiments contained 3972 foci, 280 of which were outside the MNI152 mask, leaving 3692 to  
724 be loaded into GingerALE. Then the reference map  $G$ , with voxels  $g_{i,j,k}$ , was generated by an  
725 ALE meta-analysis using the default parameters: Cluster-level FWE=0.01, Threshold  
726 Permutations=1000, P Value=0.001. The reference map is visualized in Supplementary Figure 4.

727  
728 We performed four pairwise comparisons to estimate the amount of overlap between each of the  
729 pipelines and  $G$ . For each pipeline the comparison was performed by computing an average map  
730  $\bar{R}$ , binarizing both it and  $G$ , and computing the Dice overlap between the two. The employed  
731 approach closely resembles the method of Wang et al.<sup>40</sup>, but with multiple thresholds of  
732 binarization also for  $G$ , and allowed us to plot similarity as a function of the threshold. The full  
733 details of the procedure is described in the Supplementary Methods. Additionally, to have a  
734 singular numerical basis for comparison, we computed the normalized cross-correlation<sup>46</sup>  
735 between the (non-binarized) average maps  $\bar{R}$  and the reference map  $G$ ,

736  
737

$$nCC(\bar{R}, G) = \frac{\sum_{i,j,k} (r_{i,j,k} - \bar{r})(g_{i,j,k} - \bar{g})}{\sqrt{\sum_{i,j,k} (r_{i,j,k} - \bar{r})^2 * \sum_{i,j,k} (g_{i,j,k} - \bar{g})^2}}, r \in \bar{R}, g \in G.$$

738

739 To facilitate an intuitive understanding of what parts of the brain the dementia-model is focusing  
740 on, we also performed a similar, region-wise comparison. This was done by extracting a subset  
741 of voxels from the average relevance map  $\bar{R}_{dementia}$ ,

742

$$743 \quad \bar{R}_\rho = \{r_{i,j,k} \mid (i, j, k) \in \rho\},$$

744

745 where  $\rho$  is one of 69 regions defined in the Harvard-Oxford cortical and subcortical atlases<sup>78</sup>.

746 We did the same for  $G$  and let the mean activation per region for both constitute a tuple

747

$$\left( \frac{\sum_{r \in R_\rho} r}{|R_\rho|}, \frac{\sum_{g \in G_\rho} g}{|G_\rho|} \right)$$

748

749 plotted Figure 2c. However, since it is non-trivial to determine which aggregation method  
750 corresponds to the most understandable and intuitive interpretation, we also created plots for  
751 tuples of sums,

752

$$\left( \sum_{r \in R_\rho} r, \sum_{g \in G_\rho} g \right)$$

753

754 and maximum values

755

$$\left( \max_{r \in R_\rho} r, \max_{g \in G_\rho} g \right)$$

756

757 per region in Supplementary Figure 10.

758 To quantify the importance of the spatial locations captured by the various LRP pipelines for  
759 predicting dementia, we implemented a procedure for iteratively occluding parts of the image  
760 based on the relevance maps and observing how the prediction from the dementia model changed  
761<sup>79</sup>. Still using the true positives, for each pipeline  $LRP_{task}$  for each MRI  $X_0$  we generated a  
762 baseline dementia-prediction  $\hat{y}_0$  and relevance map  $R_{task}$ . Then we located the voxel with the  
763 highest amount of relevance in  $R_{task}$  and replaced a 15x15x15 cube centred around the voxel



764 with random uniform noise  $\mathcal{U}(0, 1)$ , effectively concealing all information contained in this  
765 region. Next, we ran the modified image  $X_{task}^1$  through the dementia-model to see how the  
766 prediction  $\hat{y}_{task}^1$  changed as a function of the occlusion. Note that injecting a box of random  
767 noise into the image is not trivially equivalent to removing information, however we specifically  
768 applied the same modification in the random box-augmentation during training and are thus  
769 hopeful that the model is invariant to the injection beyond the information removal. We  
770 iteratively applied this modify-and-predict procedure, also masking out the regions from the  
771 relevant maps between each iteration to minimize overlap of occlusion windows, for 20  
772 iterations, producing a list of predictions  $[\hat{y}_0, \hat{y}_{task}^1, \hat{y}_{task}^2, \dots, \hat{y}_{task}^{19}]$  plotted for each pipeline in  
773 Figure 2d (averaged across all true positives). The rate of decline in these traces indicate the  
774 importance of the regions found in the respective relevance maps. We quantified the differences  
775 between the pipelines  $LRP_{task}$  by calculating the area over the area over their perturbation  
776 curves <sup>79</sup>,

777

$$778 \quad AOPC_{task} = \frac{1}{20} \left( \sum_{i=1}^{20} \hat{y}_0 - \hat{y}_{task}^i \right).$$

779

## 780 Exploratory analyses in the MCI cohort

781 In the exploratory MCI analyses we used  $LRP_{dementia}$  to generate predictions and relevance  
782 maps for participants from ADNI who were given an MCI diagnosis at inclusion. We first  
783 compiled the predictions and relevance maps (and the corresponding timestamps) for each  
784 participant at all timepoints into a single data structure we called a morphological record. We  
785 then tried to utilize this data structure to differentiate three groups: stable MCI patients (sMCI),  
786 progressive MCI patients (pMCI), and those who saw improvement in their cognition throughout  
787 the data collection phase. The remaining participants, e.g. those who either passed through all  
788 three diagnostic stages, or bounced between diagnoses, were excluded. Furthermore, we  
789 combined the stable and improving cohorts into a non-progressive group (nMCI) to facilitate  
790 binary group comparisons in the subsequent analyses.

791

792 For the first analysis comparing predictions in the two groups, due to variability in the total  
793 number and the frequency of visits between participants, we aimed to create a matched dataset  
794 based on visit history from the nMCI and pMCI cohorts to compare the predictions in the two

795 groups with reference to a specific timepoint. We first started with all the progressive patients  
796  $p_p \in pMCI$  who got a diagnosis at timepoint  $t_{n+1}$ , and, for each patient individually, compiled  
797 all previous visits  $t_m$ ,  $m \leq n$  into a vector  $h_p$  representing the time of the visits. The entries  
798  $d_{t_m}$  of the vector were the number of days until the diagnosis was given,  $t_{n+1} - t_m$ . For  
799 simplicity we also appended  $d_{t_{n+1}} = 0$  to  $h_p$ , such that for a single patient

800

$$801 \quad h_p = [d_{t_0}, d_{t_1}, \dots, d_{t_n}, 0].$$

802

803 Then, for each of the non-progressive patients  $p_n \in nMCI$  who didn't have a time of diagnosis  
804 (e.g.  $t_{n+1}$  is not given) we compiled a set  $H_p$  of all possible history vectors  $h_p$  by varying which  
805 visit was chosen as  $t_0$  and a terminal non-diagnosis timepoint  $t_{n+1}$ . Next, we defined a cost-  
806 criterion for matching two histories (with an equal number of visits) as the sum of absolute  
807 pairwise differences between the vectors,

808

$$809 \quad cost(h_1, h_2) = \sum_{m=0}^n |d_{t_m}^{h_1} - d_{t_m}^{h_2}|.$$

810

811 For each pair of progressive and non-progressive patients  $(p_p, p_n)$  this allowed us to calculate a  
812 best possible match, given that the stable patient had a total number of visits equal to or larger  
813 than the number of visits for the progressive patient:

814

$$815 \quad match(p_p, p_s) = \begin{cases} \min_{h \in H_{p_s}} cost(h_{p_p}, h) & \exists h \in H_{p_s} (|h| = |h_{p_p}|) \\ \infty & else \end{cases}$$

816

817

818 Finally, we compiled the cost of the optimal match from all pairs into a matrix and found the best  
819 complete matching by minimizing the total cost across this matrix using the Hungarian algorithm  
820 implemented in `scipy v1.6.3`<sup>80</sup>, such that each patient occurs in at most one pair.

821

822 We estimated differences in predictions  $\hat{y}$  between the two groups using a linear mixed model.

823 Specifically, we modelled  $\hat{y}$  at all timepoints before the terminal timepoint  $t_{n+1}$  as a function of

824 age, sex (as controlling variables), years to diagnosis, categorical group membership (nMCI,  
825 pMCI), and an interaction between years to diagnosis and group. In addition, we had an  
826 independent intercept and slope per participant. The model was fit the formula API of  
827 statsmodels v0.13.2<sup>81</sup> using the formula

$$y \sim \text{age} + \text{sex} + \text{years to diagnosis} + C(\text{group}) + \text{years to diagnosis}:C(\text{group}) \\ + (1 + \text{years to diagnosis} \mid \text{subject})$$

829  
830 on the matched dataset. A full overview of coefficients and p-values can be found in  
831 Supplementary Table 4.

832  
833 Due to the high dimensionality of the relevance maps, we decomposed them with a principal  
834 component analysis (PCA) before the final analyses. To fit the PCA we used the non-linearly  
835 registered relevance maps from a randomly selected timepoint for all MCI patients. Before fitting  
836 the model, all relevance maps were smoothed with a constant 3x3x3 blurring kernel using the  
837 convolution operation from Tensorflow 2.6 to strengthen the signal-to-noise ratio. The PCA was  
838 computed using scikit-learn v1.0.2<sup>82</sup>, retaining 64 components (out of 1137 maximally possible)  
839 in a component vector  $c = [c_0, c_1, \dots, c_{63}]$ . An axial slice from each of the 64 components  
840 visualized in MNI152-space is shown in Supplementary Figure 6.

841  
842 We fit Cox proportional hazard models using the component vectors as predictors to assess the  
843 association between the relevance maps and progression as a function of age. In addition to the  
844 components, representing the maps, we controlled for sex in the model. The p-values and  
845 coefficients can be found in Supplementary Table 5. To account for covariance between the  
846 components and the dementia-prediction  $\hat{y}$  we ran an additional model where we divided the  
847 patients into ten strata based on  $\hat{y}$ . Both models were fit using lifelines v0.27.1<sup>83</sup>.

848  
849 To further explore the prognostic efficacy of our pipeline we set up a predictive analysis for  
850 predicting progression at multiple, fixed timepoints a given number of months in the future. For  
851 each participant  $p$  with visits at timepoints  $t^p$ , we denoted the last timepoint with an MCI

852 diagnosis  $t_{neg}^p$  and the first timepoint with a dementia diagnosis (if present)  $t_{pos}^p$ . Using a fixed  
853 set of years into the future,  $\gamma \in \{1, 2, 3, 4, 5\}$ , we constructed a target variable  $z_\gamma(t^p)$  such that  
854

$$z_\gamma(t^p) = \begin{cases} 1 & t^p + \gamma \geq t_{pos}^p \\ 0 & t^p + \gamma \leq t_{neg}^p \\ NA & \text{else} \end{cases}$$

855  
856  
857 where the NAs allow for exclusion of all patients where the status at timepoint  $t^p + \gamma$  is  
858 unknown. For each  $\gamma$  we constructed the target vector  $z_\gamma$  across all timepoints for all participants  
859 with  $z_\gamma \neq NA$  and split the constituent patients  $p$  into five folds stratified on  $z_\gamma$ , sex and age,  
860 such that all timepoints from a participant resided in the same fold. Using these folds, we fit  
861 logistic regression models to predict  $z_\gamma$  with an  $l_1$ -penalty in a nested cross-validation loop,  
862 allowing us to both tune the regularization parameter  $\lambda$  and have out-of-sample predictions for  
863 all participants. For eligible participants we used all timepoints for training the models, but  
864 during testing we sampled a random timepoint per participant to ensure independence between  
865 datapoints in the final evaluation. For each  $\gamma$  we fit three models: a baseline model

866

$$\mathcal{M}_{base} := z_\gamma \sim age_{t^p} + sex + age_{t^p} \times sex$$

867  
868 to assess the bias in the dataset with respect to age at the given timepoint  $t^p$  and sex, a model  
869 using the prediction  $\hat{y}_{t^p}$  from the dementia classifier at  $t^p$  as a predictor

870

$$\mathcal{M}_{pred} := z_\gamma \sim age_{t^p} + sex + age_{t^p} \times sex + \hat{y}_{t^p} + age_{t^p} \times \hat{y}_{t^p}$$

871  
872 and a model including the relevance maps from  $t^p$ , represented by the component vector  $c_{t^p}$ ,

873

$$\mathcal{M}_{comp} := z_\gamma \sim age_{t^p} + sex + age_{t^p} \times sex + \hat{y}_{t^p} + age_{t^p} \times \hat{y}_{t^p} + c_{t^p}.$$

874  
875  
876 All models were fit and tuned using the LogisticRegressionCV interface of sklearn v1.0.2<sup>82</sup>. We  
877 compared models by measuring the mean AUC across the five folds (Supplementary Table 6).

878 To evaluate clinical applicability we also report accuracy, positive predictive value, sensitivity,  
879 and specificity (Table 2). To determine whether the more complex models represented  
880 significant improvements we employed a one-sided Wilcoxon signed-rank test from `scipy v1.9.3`  
881 <sup>80</sup> to do pairwise comparisons between  $\mathcal{M}_{base}$  and  $\mathcal{M}_{pred}$ , and  $\mathcal{M}_{pred}$ , and  $\mathcal{M}_{comp}$  across the  
882 five out-of-sample AUCs independently.

883  
884 To assess whether the relevance maps were associated with specific cognitive functions we  
885 associated aspects of them with performance on various cognitive tests. We first extracted test  
886 results from seven neuropsychological batteries which spanned all ADNI waves and contained  
887 high-level summary scores from the ADNI website (Supplementary Table 7). We then manually  
888 extracted 17 summary scores spanning different, but overlapping, cognitive domains  
889 (Supplementary Table 8). The component vectors  $c$  were used as proxies for the relevance maps,  
890 where each  $c_i$  represented a template for localization of pathology. We matched 2402 component  
891 vectors with test results from 733 MCI patients, forming a basis for the comparison. We then  
892 calculated the univariate association between cognitive performance according to each of the 17  
893 with each of the dimensions  $c_i \in c$ , while including age and sex as covariates for correction. To  
894 isolate the effect of the localization we also corrected for dementia-prediction,  $\hat{y}$ . When a patient  
895 had multiple potential matches, a random timepoint was selected, and the final number of  
896 datapoints used in the analyses varied from 518 to 675. Correction for multiple testing was done  
897 with the Benjamini-Hochberg procedure. To ensure the associations were not confounded by  
898 collinearities between  $c$  and  $\hat{y}$ , we also performed an equivalent analysis without correction to  
899 observe whether the sign of the coefficients changed.

900

## 901 References

- 902 1. Woo, C.-W., Chang, L. J., Lindquist, M. A. & Wager, T. D. Building better biomarkers: brain  
903 models in translational neuroimaging. *Nat Neurosci* **20**, 365–377 (2017).
- 904 2. Bethlehem, R. a. l. *et al.* Brain charts for the human lifespan. *Nature* **604**, 525–533 (2022).

- 905 3. Marek, S. *et al.* Reproducible brain-wide association studies require thousands of  
906 individuals. *Nature* **603**, 654–660 (2022).
- 907 4. Arbabshirani, M. R., Plis, S., Sui, J. & Calhoun, V. D. Single subject prediction of brain  
908 disorders in neuroimaging: Promises and pitfalls. *NeuroImage* **145**, 137–165 (2017).
- 909 5. Sui, J., Jiang, R., Bustillo, J. & Calhoun, V. Neuroimaging-based Individualized Prediction of  
910 Cognition and Behavior for Mental Disorders and Health: Methods and Promises. *Biological*  
911 *Psychiatry* **88**, 818–828 (2020).
- 912 6. Davatzikos, C. Why voxel-based morphometric analysis should be used with great caution  
913 when characterizing group differences. *NeuroImage* **23**, 17–20 (2004).
- 914 7. Westlin, C. *et al.* Improving the study of brain-behavior relationships by revisiting basic  
915 assumptions. *Trends in Cognitive Sciences* **27**, 246–257 (2023).
- 916 8. Bzdok, D. & Ioannidis, J. P. A. Exploration, Inference, and Prediction in Neuroscience and  
917 Biomedicine. *Trends in Neurosciences* **42**, 251–262 (2019).
- 918 9. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
- 919 10. Rudin, C. Stop explaining black box machine learning models for high stakes decisions and  
920 use interpretable models instead. *Nat Mach Intell* **1**, 206–215 (2019).
- 921 11. Gauthier S, Webster C, Servaes S, Morais JA, Rosa-Neto P. World Alzheimer Report 2022 –  
922 Life after diagnosis: Navigating treatment, care and support. (2022).
- 923 12. Nichols, E. *et al.* Global, regional, and national burden of Alzheimer’s disease and other  
924 dementias, 1990–2016: a systematic analysis for the Global Burden of Disease Study 2016.  
925 *The Lancet Neurology* **18**, 88–106 (2019).

- 926 13. Vos, T. *et al.* Global burden of 369 diseases and injuries in 204 countries and territories,  
927 1990–2019: a systematic analysis for the Global Burden of Disease Study 2019. *The Lancet*  
928 **396**, 1204–1222 (2020).
- 929 14. World Health Organization. Global status report on the public health response to dementia.  
930 (2021).
- 931 15. Nichols, E. *et al.* Estimation of the global prevalence of dementia in 2019 and forecasted  
932 prevalence in 2050: an analysis for the Global Burden of Disease Study 2019. *The Lancet*  
933 *Public Health* **7**, e105–e125 (2022).
- 934 16. Feldman, H. H. *et al.* Diagnosis and treatment of dementia: 2. Diagnosis. *CMAJ* **178**, 825–  
935 836 (2008).
- 936 17. Karantzoulis, S. & Galvin, J. E. Distinguishing Alzheimer’s disease from other major forms of  
937 dementia. *Expert Rev Neurother* **11**, 1579–1591 (2011).
- 938 18. Echávarri, C. *et al.* Co-occurrence of Different Pathologies in Dementia: Implications for  
939 Dementia Diagnosis. *Journal of Alzheimer’s Disease* **30**, 909–917 (2012).
- 940 19. Schneider, J. A. Neuropathology of Dementia Disorders. *CONTINUUM: Lifelong Learning in*  
941 *Neurology* **28**, 834 (2022).
- 942 20. Ryan, J., Fransquet, P., Wrigglesworth, J. & Lacaze, P. Phenotypic Heterogeneity in  
943 Dementia: A Challenge for Epidemiology and Biomarker Studies. *Front Public Health* **6**, 181  
944 (2018).
- 945 21. Ikram, M. A. *et al.* Brain tissue volumes in relation to cognitive function and risk of  
946 dementia. *Neurobiology of Aging* **31**, 378–386 (2010).

- 947 22. McDonald, C. R. *et al.* Relationship between regional atrophy rates and cognitive decline in  
948 mild cognitive impairment. *Neurobiology of Aging* **33**, 242–253 (2012).
- 949 23. Ferreira, D., Nordberg, A. & Westman, E. Biological subtypes of Alzheimer disease: A  
950 systematic review and meta-analysis. *Neurology* **94**, 436–448 (2020).
- 951 24. Verdi, S., Marquand, A. F., Schott, J. M. & Cole, J. H. Beyond the average patient: how  
952 neuroimaging models can address heterogeneity in dementia. *Brain* **144**, 2946–2953  
953 (2021).
- 954 25. Rasmussen, J. & Langerman, H. Alzheimer’s Disease – Why We Need Early Diagnosis.  
955 *Degener Neurol Neuromuscul Dis* **9**, 123–130 (2019).
- 956 26. Robinson, L., Tang, E. & Taylor, J.-P. Dementia: timely diagnosis and early intervention. *BMJ*  
957 **350**, h3029 (2015).
- 958 27. Lu, B. *et al.* A practical Alzheimer’s disease classifier via brain imaging-based deep learning  
959 on 85,721 samples. *Journal of Big Data* **9**, 101 (2022).
- 960 28. Mirzaei, G. & Adeli, H. Machine learning techniques for diagnosis of alzheimer disease, mild  
961 cognitive disorder, and other types of dementia. *Biomedical Signal Processing and Control*  
962 **72**, 103293 (2022).
- 963 29. Mirabnahrzam, G. *et al.* Predicting time-to-conversion for dementia of Alzheimer’s type  
964 using multi-modal deep survival analysis. *Neurobiology of Aging* **121**, 139–156 (2023).
- 965 30. Castellazzi, G. *et al.* A Machine Learning Approach for the Differential Diagnosis of  
966 Alzheimer and Vascular Dementia Fed by MRI Selected Features. *Frontiers in*  
967 *Neuroinformatics* **14**, (2020).



- 968 31. Yao, A. D., Cheng, D. L., Pan, I. & Kitamura, F. Deep Learning in Neuroradiology: A  
969 Systematic Review of Current Algorithms and Approaches for the New Wave of Imaging  
970 Technology. *Radiology: Artificial Intelligence* **2**, e190026 (2020).
- 971 32. Kundu, S. AI in medicine must be explainable. *Nat Med* **27**, 1328–1328 (2021).
- 972 33. Samek, W., Montavon, G., Vedaldi, A., Hansen, L. K. & Müller, K.-R. *Explainable AI:  
973 Interpreting, Explaining and Visualizing Deep Learning*. (Springer Nature, 2019).
- 974 34. Barredo Arrieta, A. *et al.* Explainable Artificial Intelligence (XAI): Concepts, taxonomies,  
975 opportunities and challenges toward responsible AI. *Information Fusion* **58**, 82–115 (2020).
- 976 35. Samek, W. & Müller, K.-R. Towards Explainable Artificial Intelligence. in *Explainable AI:  
977 Interpreting, Explaining and Visualizing Deep Learning* (eds. Samek, W., Montavon, G.,  
978 Vedaldi, A., Hansen, L. K. & Müller, K.-R.) 5–22 (Springer International Publishing, Cham,  
979 2019). doi:10.1007/978-3-030-28954-6\_1.
- 980 36. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep Inside Convolutional Networks: Visualising  
981 Image Classification Models and Saliency Maps. Preprint at  
982 <https://doi.org/10.48550/arXiv.1312.6034> (2014).
- 983 37. Bach, S. *et al.* On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise  
984 Relevance Propagation. *PLOS ONE* **10**, e0130140 (2015).
- 985 38. Martin, S. A., Townend, F. J., Barkhof, F. & Cole, J. H. Interpretable machine learning for  
986 dementia: A systematic review. *Alzheimer's & Dementia* (2023).
- 987 39. Böhle, M., Eitel, F., Weygandt, M. & Ritter, K. Layer-Wise Relevance Propagation for  
988 Explaining Deep Neural Network Decisions in MRI-Based Alzheimer's Disease Classification.  
989 *Frontiers in Aging Neuroscience* **11**, (2019).

- 990 40. Wang, D. *et al.* Deep neural network heatmaps capture Alzheimer’s disease patterns  
991 reported in a large meta-analysis of neuroimaging studies. *NeuroImage* **269**, 119929 (2023).
- 992 41. Dyrba, M. *et al.* Improving 3D convolutional neural network comprehensibility via  
993 interactive visualization of relevance maps: evaluation in Alzheimer’s disease. *Alzheimer’s*  
994 *Research & Therapy* **13**, 191 (2021).
- 995 42. Dyrba, M., Pallath, A. H. & Marzban, E. N. Comparison of CNN Visualization Methods to Aid  
996 Model Interpretability for Detecting Alzheimer’s Disease. in *Bildverarbeitung für die Medizin*  
997 *2020* (eds. Tolxdorff, T. et al.) 307–312 (Springer Fachmedien, Wiesbaden, 2020).  
998 doi:10.1007/978-3-658-29267-6\_68.
- 999 43. Martin, S. A., Townend, F. J., Barkhof, F. & Cole, J. H. Interpretable machine learning for  
1000 dementia: A systematic review. *Alzheimer’s & Dementia* **19**, 2135–2149 (2023).
- 1001 44. Kohlbrenner, M. *et al.* Towards Best Practice in Explaining Neural Network Decisions with  
1002 LRP. in *2020 International Joint Conference on Neural Networks (IJCNN)* 1–7 (2020).  
1003 doi:10.1109/IJCNN48605.2020.9206975.
- 1004 45. Eickhoff, S. B., Bzdok, D., Laird, A. R., Kurth, F. & Fox, P. T. Activation Likelihood Estimation  
1005 meta-analysis revisited. *Neuroimage* **59**, 2349–2361 (2012).
- 1006 46. Briechele, K. & Hanebeck, U. D. Template matching using fast normalized cross correlation. in  
1007 (eds. Casasent, D. P. & Chao, T.-H.) 95–102 (Orlando, FL, 2001). doi:10.1117/12.421129.
- 1008 47. Eitel, F. & Ritter, K. Testing the Robustness of Attribution Methods for Convolutional Neural  
1009 Networks in MRI-Based Alzheimer’s Disease Classification. in *Interpretability of Machine*  
1010 *Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision*

- 1011        *Support* (eds. Suzuki, K. et al.) 3–11 (Springer International Publishing, Cham, 2019).
- 1012        doi:10.1007/978-3-030-33850-3\_1.
- 1013    48. Erasmus, A., Brunet, T. D. P. & Fisher, E. What is Interpretability? *Philos. Technol.* **34**, 833–
- 1014        862 (2021).
- 1015    49. Ghassemi, M., Oakden-Rayner, L. & Beam, A. L. The false hope of current approaches to
- 1016        explainable artificial intelligence in health care. *The Lancet Digital Health* **3**, e745–e750
- 1017        (2021).
- 1018    50. Amann, J. *et al.* To explain or not to explain?—Artificial intelligence explainability in clinical
- 1019        decision support systems. *PLOS Digital Health* **1**, e0000016 (2022).
- 1020    51. Adebayo, J. *et al.* Sanity Checks for Saliency Maps. *arXiv:1810.03292 [cs, stat]* (2020).
- 1021    52. Kindermans, P.-J. *et al.* The (Un)reliability of Saliency Methods. in *Explainable AI:*
- 1022        *Interpreting, Explaining and Visualizing Deep Learning* (eds. Samek, W., Montavon, G.,
- 1023        Vedaldi, A., Hansen, L. K. & Müller, K.-R.) 267–280 (Springer International Publishing, Cham,
- 1024        2019). doi:10.1007/978-3-030-28954-6\_14.
- 1025    53. Nie, X. *et al.* Subregional Structural Alterations in Hippocampus and Nucleus Accumbens
- 1026        Correlate with the Clinical Impairment in Patients with Alzheimer’s Disease Clinical
- 1027        Spectrum: Parallel Combining Volume and Vertex-Based Approach. *Frontiers in Neurology* **8**,
- 1028        (2017).
- 1029    54. Poulin, S. P., Dautoff, R., Morris, J. C., Barrett, L. F. & Dickerson, B. C. Amygdala atrophy is
- 1030        prominent in early Alzheimer’s disease and relates to symptom severity. *Psychiatry Res* **194**,
- 1031        7–13 (2011).

- 1032 55. Van Hoesen, G. W., Augustinack, J. C., Dierking, J., Redman, S. J. & Thangavel, R. The  
1033 parahippocampal gyrus in Alzheimer's disease. Clinical and preclinical neuroanatomical  
1034 correlates. *Ann N Y Acad Sci* **911**, 254–274 (2000).
- 1035 56. Visser, P. J. *et al.* Medial temporal lobe atrophy and memory dysfunction as predictors for  
1036 dementia in subjects with mild cognitive impairment. *J Neurol* **246**, 477–485 (1999).
- 1037 57. Echávarri, C. *et al.* Atrophy in the parahippocampal gyrus as an early biomarker of  
1038 Alzheimer's disease. *Brain Struct Funct* **215**, 265–271 (2011).
- 1039 58. Dickerson, B. C. *et al.* Alzheimer-signature MRI biomarker predicts AD dementia in  
1040 cognitively normal adults. *Neurology* **76**, 1395–1402 (2011).
- 1041 59. Rafii, M. S. & Aisen, P. S. Detection and treatment of Alzheimer's disease in its preclinical  
1042 stage. *Nat Aging* **3**, 520–531 (2023).
- 1043 60. Frisoni, G. B. *et al.* Dementia prevention in memory clinics: recommendations from the  
1044 European task force for brain health services. *The Lancet Regional Health – Europe* **26**,  
1045 (2023).
- 1046 61. de Vugt, M. E. & Verhey, F. R. J. The impact of early dementia diagnosis and intervention on  
1047 informal caregivers. *Progress in Neurobiology* **110**, 54–62 (2013).
- 1048 62. Leonardsen, E. H. *et al.* Deep neural networks learn general and clinically relevant  
1049 representations of the ageing brain. *NeuroImage* **256**, 119210 (2022).
- 1050 63. Mårtensson, G. *et al.* The reliability of a deep learning model in clinical out-of-distribution  
1051 MRI data: A multicohort study. *Medical Image Analysis* **66**, 101714 (2020).
- 1052 64. Herzog, C. On the Ethical and Epistemological Utility of Explicable AI in Medicine. *Philos.*  
1053 *Technol.* **35**, 50 (2022).

- 1054 65. Petersen, R. C. *et al.* Alzheimer’s Disease Neuroimaging Initiative (ADNI): clinical  
1055 characterization. *Neurology* **74**, 201–209 (2010).
- 1056 66. Jenkinson, M., Beckmann, C. F., Behrens, T. E. J., Woolrich, M. W. & Smith, S. M. FSL.  
1057 *NeuroImage* **62**, 782–790 (2012).
- 1058 67. Ségonne, F. *et al.* A hybrid approach to the skull stripping problem in MRI. *Neuroimage* **22**,  
1059 1060–1075 (2004).
- 1060 68. Jenkinson, M. & Smith, S. A global optimisation method for robust affine registration of  
1061 brain images. *Med Image Anal* **5**, 143–156 (2001).
- 1062 69. Gong, W., Beckmann, C. F., Vedaldi, A., Smith, S. M. & Peng, H. Optimising a Simple Fully  
1063 Convolutional Network for Accurate Brain Age Prediction in the PAC 2019 Challenge.  
1064 *Frontiers in Psychiatry* **12**, (2021).
- 1065 70. Peng, H., Gong, W., Beckmann, C. F., Vedaldi, A. & Smith, S. M. Accurate brain age  
1066 prediction with lightweight deep neural networks. *Medical Image Analysis* **68**, 101871  
1067 (2021).
- 1068 71. Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed  
1069 Systems. **19** (2015).
- 1070 72. Chollet, F. & others. Keras. <https://github.com/fchollet/keras> (2015).
- 1071 73. Smith, L. N. Cyclical Learning Rates for Training Neural Networks. *arXiv:1506.01186 [cs]*  
1072 (2017).
- 1073 74. Guillemot, M., Heusele, C., Korichi, R., Schnebert, S. & Chen, L. Breaking Batch  
1074 Normalization for better explainability of Deep Neural Networks through Layer-wise  
1075 Relevance Propagation. *arXiv:2002.11018 [cs, stat]* (2020).

- 1076 75. Montavon, G., Binder, A., Lapuschkin, S., Samek, W. & Müller, K.-R. Layer-Wise Relevance  
1077 Propagation: An Overview. in *Explainable AI: Interpreting, Explaining and Visualizing Deep*  
1078 *Learning* (eds. Samek, W., Montavon, G., Vedaldi, A., Hansen, L. K. & Müller, K.-R.) 193–209  
1079 (Springer International Publishing, Cham, 2019). doi:10.1007/978-3-030-28954-6\_10.
- 1080 76. Brett, M. *et al.* nipy/nibabel: 3.2.2. Zenodo <https://doi.org/10.5281/zenodo.6617121>  
1081 (2022).
- 1082 77. Laird, A. R., Lancaster, J. L. & Fox, P. T. BrainMap: the social evolution of a human brain  
1083 mapping database. *Neuroinformatics* **3**, 65–78 (2005).
- 1084 78. Desikan, R. S. *et al.* An automated labeling system for subdividing the human cerebral  
1085 cortex on MRI scans into gyral based regions of interest. *Neuroimage* **31**, 968–980 (2006).
- 1086 79. Samek, W., Binder, A., Montavon, G., Lapuschkin, S. & Müller, K.-R. Evaluating the  
1087 Visualization of What a Deep Neural Network Has Learned. *IEEE Transactions on Neural*  
1088 *Networks and Learning Systems* **28**, 2660–2673 (2017).
- 1089 80. Virtanen, P. *et al.* SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python.  
1090 *Nature Methods* **17**, 261–272 (2020).
- 1091 81. Seabold, S. & Perktold, J. Statsmodels: Econometric and Statistical Modeling with Python. in  
1092 92–96 (Austin, Texas, 2010). doi:10.25080/Majora-92bf1922-011.
- 1093 82. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine Learning*  
1094 *Research* **12**, 2825–2830 (2011).
- 1095 83. Davidson-Pilon, C. lifelines: survival analysis in Python. *Journal of Open Source Software* **4**,  
1096 1317 (2019).
- 1097

## 1098 Author contributions:

1099 Conceptualization: EHL, TW, LTW, YW. Data curation: KP, EW, GS. Formal analysis: EHL.  
1100 Funding acquisition: OAA, YW. Investigation: EHL, JMR, DVP, TK, AM, OAA, TW, LTW,  
1101 YW. Methodology: EHL, EG, ND, TS, ØS, TW, LTW, YW. Project administration: GS, OAA,  
1102 LTW, YW. Software: EHL. Supervision: TW, LTW, YW. Validation: EHL. Visualization: EHL.  
1103 Writing – original draft: EHL, TW, LTW, YW. Writing – review & editing: KP, EG, ND, TS,  
1104 JMR, DVP, ØS, TK, EW, AM, GS, OAA.

## 1105 Acknowledgements

1106 This work was funded by the UiO:LifeScience Convergence Environment (project: 4MENT), the  
1107 Research Council of Norway (302854), and the European Research Council under the European  
1108 Union's Horizon 2020 research and Innovation program (802998). The Southern and Eastern  
1109 Norway Regional Health Authority supported the study through funding for KP but was not  
1110 involved in conducting the study or in preparation of the manuscript. TW acknowledges funding  
1111 from the German Research Foundation (DFG) Emmy Noether: 513851350. The work was  
1112 performed on the Service for Sensitive Data (TSD) platform, owned by the University of Oslo,  
1113 operated, and developed by the TSD service group at the University of Oslo IT-Department  
1114 (USIT). We also acknowledge the computational resources provided by UNINETT Sigma2 - the  
1115 National Infrastructure for High Performance Computing and Data Storage in Norway – with  
1116 project no. (nn9769k/ns9769k).

## 1117 Competing interests:

1118 KP report work with Roche BN29553 and Novo Nordisk NN6535-4730 trials; All other authors  
1119 declare that they have no competing interests.

## 1120 Data availability

1121 The data used in this study were gathered from various sources, an overview including  
1122 acknowledgements of their respective funding sources is provided in Supplementary Table 1.  
1123 Among others, data used in the preparation of this article was obtained from the Alzheimer's  
1124 Disease Neuroimaging Initiative (ADNI, see [adni.loni.usc.edu](http://adni.loni.usc.edu) for further details), the Australian  
1125 Imaging Biomarkers and Lifestyle flagship study of ageing (AIBL, [www.aibl.csiro.au](http://www.aibl.csiro.au)) the

1126 AddNeuroMed consortium, and MIRIAD ([www.nitrc.org/projects/miriad](http://www.nitrc.org/projects/miriad)). The investigators  
1127 within these studies contributed to the design and implementation of the data collection process  
1128 but did not participate in the analysis or writing of this report, and this publication is solely the  
1129 responsibility of the authors.

### 1130 Code availability:

1131 The trained model and explainable pipeline and the underlying code will be made available at  
1132 <https://github.com/estehnl/pyment-public> upon publication. Generic code for generating  
1133 explanations for 3D CNNs is available at <https://github.com/estehnl/keras-explainability>.