

## **A Practical Alzheimer Disease Classifier via Brain Imaging-Based Deep Learning on 85,721 Samples**

Bin Lu<sup>1,2</sup>, Hui-Xian Li<sup>1,2</sup>, Zhi-Kai Chang<sup>1,2</sup>, Le Li<sup>3</sup>, Ning-Xuan Chen<sup>1,2</sup>, Zhi-Chen Zhu<sup>1,2</sup>, Hui-Xia Zhou<sup>1,2</sup>, Xue-Ying Li<sup>1,4,5</sup>, Yu-Wei Wang<sup>1,2</sup>, Shi-Xian Cui<sup>1,4,5</sup>, Zhao-Yu Deng<sup>1,2</sup>, Zhen Fan<sup>6</sup>, Hong Yang<sup>7</sup>, Xiao Chen<sup>1,2</sup>, Paul M. Thompson<sup>8</sup>, Francisco Xavier Castellanos<sup>9,10</sup>, Chao-Gan Yan<sup>1,2,11,12\*</sup>, for the Alzheimer's Disease Neuroimaging Initiative \*\*

<sup>1</sup>CAS Key Laboratory of Behavioral Science, Institute of Psychology, Beijing, China; <sup>2</sup>Department of Psychology, University of Chinese Academy of Sciences, Beijing, China; <sup>3</sup>Center for Cognitive Science of Language, Beijing Language and Culture University, Beijing, China; <sup>4</sup>Sino-Danish College, University of Chinese Academy of Science, Beijing, China; <sup>5</sup>Sino-Danish Center for Education and Research, Beijing, China; <sup>6</sup>Department of Neurosurgery, Huashan Hospital, Fudan University, Shanghai Neurosurgical Clinical Center, Shanghai, China; <sup>7</sup>Department of Radiology, The First Affiliated Hospital, College of Medicine, Zhejiang University, Hangzhou, Zhejiang, China; <sup>8</sup>Imaging Genetics Center, Mark & Mary Stevens Institute for Neuroimaging & Informatics, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA; <sup>9</sup>Department of Child and Adolescent Psychiatry, NYU Grossman School of Medicine, New York, NY, USA; <sup>10</sup>Nathan Kline Institute for Psychiatric Research, Orangeburg, NY, USA; <sup>11</sup>International Big-Data Research Center for Depression (IBRCD), Institute of Psychology, Chinese Academy of Sciences, Beijing, China; <sup>12</sup>Magnetic Resonance Imaging Research Center, Institute of Psychology, Chinese Academy of Sciences, Beijing, China. \*e-mail: [ycg.yan@gmail.com](mailto:ycg.yan@gmail.com). \*\*Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database ([adni.loni.usc.edu](http://adni.loni.usc.edu)). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: [http://adni.loni.usc.edu/wp-content/uploads/how\\_to\\_apply/ADNI\\_Acknowledgement\\_List.pdf](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)

1     **Abstract**

2     Beyond detecting brain lesions or tumors, comparatively little success has been attained in  
3     identifying brain disorders such as Alzheimer’s disease (AD), based on magnetic resonance  
4     imaging (MRI). Many machine learning algorithms to detect AD have been trained using  
5     limited training data, meaning they often generalize poorly when applied to scans from  
6     previously unseen populations. Here we aimed to build a practical brain imaging-based AD  
7     diagnostic classifier using deep learning/transfer learning on dataset of unprecedented size  
8     and diversity. We pooled MRI data from more than 217 sites/scanners to constitute the largest  
9     brain MRI sample to date (85,721 scans from 50,876 participants). Next, we applied a  
10    state-of-the-art deep convolutional neural network, Inception-ResNet-V2, to build a sex  
11    classifier with high generalization capability. The sex classifier achieved 94.9% accuracy and  
12    served as a base model in transfer learning for the objective diagnosis of AD. After transfer  
13    learning, the model fine-tuned for AD classification achieved 91.3% accuracy in  
14    leave-sites-out cross-validation on the Alzheimer’s Disease Neuroimaging Initiative (ADNI)  
15    dataset and 94.2%/87.9% accuracy for direct tests on two unseen independent datasets  
16    (AIBL/OASIS). When this AD classifier was tested on brain images from unseen mild  
17    cognitive impairment (MCI) patients, MCI patients who finally converted to AD were 3 times  
18    more likely to be predicted as AD than MCI patients who did not convert (65.2% vs 20.6%).  
19    Predicted scores from the AD classifier showed significant correlations with illness severity.  
20    In sum, the proposed AD classifier could offer a medical-grade biomarker that could be  
21    integrated into AD diagnostic practice. Our trained model, code and preprocessed data are  
22    freely available to the research community.

23

24    **Keywords**

25    Alzheimer’s disease, convolutional neural network, magnetic resonance brain imaging, sex  
26    difference, transfer learning

## 27 **1. Introduction**

28 Magnetic resonance imaging (MRI) is widely used in neuroradiology to detect brain lesions  
29 including stroke, vascular disease, and tumor tissue. Even so, MRI has been less useful in  
30 definitively identifying degenerative diseases including Alzheimer's disease (AD), mainly  
31 because signatures of the disease are diffusely found in the images and hard to distinguish  
32 from normal aging. Machine learning and deep learning methods have been trained on  
33 relatively small datasets, but the limited training data often leads to poor generalization  
34 performance on new datasets not used to train the algorithms. In the current study, we aim to  
35 create a practical brain imaging-based AD classifier with high generalization capability via  
36 learning/transfer learning on a diverse range of large-scale datasets.

37

38 In recently updated AD diagnostic criteria, such as those proposed by IWG-2 and NIA-AA,  
39 biomarkers such as amyloid measures from cerebrospinal fluid (CSF) and amyloid-sensitive  
40 positron emission tomography (PET) have been integrated to improve the specificity in  
41 diagnosing AD<sup>1,2</sup>. However, the invasive nature and the low sensitivity of these biomarkers  
42 hampers their application in routine clinical settings. For example, the IWG-1 criteria only  
43 achieved 68% sensitivity and the NIA-AA criteria achieved only 65.6% sensitivity<sup>3,4</sup>. A novel  
44 non-invasive biomarker with both high sensitivity and specificity is needed for diagnosing  
45 AD. Structural MRI is a promising candidate considering its non-invasive nature and wider  
46 availability than PET. In addition, there are well-developed MRI data preprocessing pipelines  
47 that make it feasible to integrate MRI biomarkers into automatic end-to-end deep learning  
48 algorithms. Deep learning has already been successfully deployed in real-world scenarios  
49 such as extreme weather condition prediction<sup>5</sup>, aftershock pattern prediction<sup>6</sup> and automatic  
50 speech recognition<sup>7</sup>. In clinical scenarios, convolutional neural networks (CNN) – a  
51 widely-used architecture that is well-suited for image-based deep learning has been  
52 successfully used for objective diagnosis of retinal diseases<sup>8</sup> and skin cancer<sup>9</sup>, and for breast  
53 cancer screening<sup>10</sup>. Given the success of recent CNN algorithms in pattern recognition,  
54 similar MRI-based diagnostic classifiers are likely to be highly valuable, if they can be  
55 integrated into routine neurological practice.

56

57 Even so, prior attempts at MRI-based AD diagnosis have yet to reach clinical utility. A major  
58 challenge for brain MRI-based algorithms, especially if they are trained on limited data, is  
59 their failure to generalize. Brain imaging data varies depending on scanner characteristics  
60 such as scanner vendor and magnetic field strength, head coil hardware, the pulse sequence,  
61 applied gradient fields, reconstruction methods, scanning parameters, voxel size, field of view,  
62 etc. Participants also differ in sex, age, race and education, and robust methods need to work  
63 well on diverse populations. These variations in the scans - and in the populations studied -  
64 make it hard for a brain imaging-based classifier trained on data from a single site (or a few  
65 sites) to generalize to data from unseen sites/scanners. This has prevented brain  
66 imaging-based classifiers from becoming practically useful, in clinical settings. For instance,  
67 Qiu and colleagues built a deep-learning classifier for AD with an average accuracy of 82.2%  
68 using brain imaging data from four datasets<sup>11</sup>. However, when tested on the FHS  
69 (Framingham Heart Study) dataset, the accuracy and specificity of the AD classifier dropped  
70 to 76.6% and 71.2%, with a relatively high sensitivity (90.1%). On the contrary, when tested  
71 on AIBL dataset (from the Australian Imaging, Biomarker and Lifestyle Study of Ageing),  
72 the same classifier achieved relatively high accuracy and specificity (87.0% and 92.4%), but  
73 the sensitivity was poor (59.4%). The variable accuracy and inconsistent tradeoff between  
74 sensitivity and specificity in data from different medical centers hampers these proposed  
75 methods from being deployed across multiple clinical institutions. To alleviate the  
76 unsatisfactory generalization performance, Bashyam et al. used a more heterogeneous sample  
77 to build a brain age prediction model that would be more generalizable to data from unseen  
78 sites/scanners<sup>12</sup>. However, when transfer learning to AD, they only used random  
79 cross-validation on the ADNI dataset with an accuracy of 86%, and did not implement  
80 independent dataset validation. Reviews of brain imaging-based AD classifiers suggest that  
81 most machine learning methods have been trained on scans from only a few hundreds of  
82 participants, which makes them unable to achieve stable performance when validated  
83 independently<sup>13</sup>.

84

85 Therefore, one bottleneck in developing a practical brain imaging-based classifier is the  
86 variety and comprehensiveness of training datasets. As most publicly available AD datasets  
87 only contain several hundred patients and corresponding healthy controls, directly training  
88 models on these datasets may result in overfitting with poor generalization to unseen test  
89 data<sup>13</sup>. In the current paper, we propose a transfer learning framework to solve this problem,  
90 by training a model on a certain characteristic for which there are abundant samples available,  
91 and fine-tuning it to another characteristic in smaller samples<sup>14</sup>, following successful  
92 examples in diagnosing retinal disease<sup>8</sup> and skin cancer<sup>9</sup>. In the brain imaging field, the  
93 scientific community has shared hundreds of thousands of brain images from hundreds of  
94 sites/scanners all over the world. Nonetheless, no studies have fully implemented this  
95 abundant resource to promote the generalizability of an AD classifier. Thus, in the current  
96 study, we used the largest and most diverse sample to date ( $n = 85,721$  from more than 217  
97 sites/scanners, see Supplementary Table 1) to pre-train a brain imaging based classifier, in  
98 order to ensure high generalization capability. After that, the pre-trained model was  
99 fine-tuned for AD classification and was validated through leave-sites-out cross-validation  
100 and independent validation. Mild cognitive impairment (MCI) is a syndrome defined as  
101 relative cognitive decline without symptoms interfering with daily life, but more than half of  
102 MCI patients progress to dementia within 5 years<sup>15</sup>. Discrimination between MCI patients  
103 who will progress to AD (pMCI) and MCI patients who will not progress to AD over a given  
104 time interval (stable MCI, or sMCI) would facilitate the early treatment of pMCI. Therefore,  
105 we used the present AD classifier to predict progression in MCI patients to further evaluate  
106 its generalization capability.

107

108 The goal of the present study was to build a practical AD classifier with high generalization  
109 capability. We incorporated three design features to improve the clinical utility of the method.  
110 First, we trained and tested the algorithm on a datasets of unprecedented size and diversity -  
111 from more than 217 sites/scanners - the variety of training samples critical for improving the  
112 generalization capability of the models. Secondly, a rigorous leave-datasets/sites-out  
113 cross-validation and independent validation was implemented to make sure that the classifier

114 accuracy would be robust to site/scanner variability. Thirdly, compared to 2D modules  
115 (feature detectors) typically used in CNNs for natural images, fully 3D convolution filters in  
116 the present study were used capture more sophisticated and distributed spatial features for  
117 diagnostic classification. We also openly share our preprocessed data, trained model, code  
118 and framework, and have built an online predicting website (<http://brainimagenet.org>) for  
119 anyone interested in testing our classifier with brain imaging data from any research  
120 participants and any scanner.

121

## 122 **2. Materials and methods**

123 **Data acquisition.** We submitted data access applications to nearly all the open-access brain  
124 imaging data archives and received permissions from the administrators of 34 datasets. The  
125 full dataset list is shown in **Table S1**. Deidentified data were contributed from datasets  
126 approved by local Institutional Review Boards. The reanalysis of these data was approved by  
127 the Institutional Review Board of Institute of Psychology, Chinese Academy of Sciences. All  
128 participants had provided written informed consent at their local institution. All 50,876  
129 participants (contributing 85,721 samples) had at least one session with a T1-weighted  
130 structural brain image and information on their sex and age. For participants with multiple  
131 sessions of structural images, each image was considered as an independent sample for data  
132 augmentation in training. Importantly, scans from the same person were never split into  
133 training and testing sets, as that could artifactually inflate performance.

134

135 **MRI preprocessing.** We did not feed raw data into the classifier for training, but used  
136 accepted pre-processing pipelines that are known to generate valuable features from the brain  
137 scans. The brain structural data were segmented and normalized to acquire grey matter  
138 density (GMD) and grey matter volume (GMV) maps. Specifically, we used the voxel-based  
139 morphometry (VBM) analysis module within Data Processing Assistant for Resting-State  
140 fMRI (DPARSF)<sup>16</sup>, which is based on SPM<sup>17</sup>, to segment individual T1-weighted images  
141 into grey matter, white matter and cerebrospinal fluid (CSF). Then, the segmented images  
142 were transformed from individual native space to MNI space (a coordinate system created by

143 Montreal Neurological Institute) using the Diffeomorphic Anatomical Registration Through  
144 Exponentiated Lie algebra (DARTTEL) tool<sup>18</sup>. Two voxel-based structural metrics, GMD and  
145 GMV were fed into the deep learning classifier as two features for each participant. GMV  
146 was based on modulated GMD images by using the Jacobian determinants derived from the  
147 spatial normalization in the VBM analysis<sup>19</sup>.

148

149 **Quality control.** Poor quality raw structural images would produce distorted GMD and  
150 GMV maps during segmentation and normalization. To remove such participants from  
151 affecting the training of the classifiers, we excluded participants in each dataset with a spatial  
152 correlation exceeding the threshold defined by (mean - 2SD) of the Pearson's correlation  
153 between each participant's GMV map and the grand mean GMV template (See **Fig. S6** for  
154 the distribution of correlations for each dataset). The grand mean GMV template was  
155 generated by randomly selecting 10 participants from each dataset and averaging the GMV  
156 maps of all these 340 (from 34 datasets) participants. All these participants were visually  
157 checked for image quality. After quality control, 83,735 samples were retained for classifier  
158 training.

159

160 **Deep learning: classifier training and testing for sex.** We trained a 3-dimensional  
161 Inception-ResNet-v2<sup>20</sup> model adopted from its 2-dimensional version in the Keras built-in  
162 application (see **Fig. 1A** for its structure). This is a state-of-the-art model in pattern  
163 recognition, and it integrates two classical series of CNN models, Inception and ResNet. We  
164 replaced the convolution, pooling and normalization modules with their 3-dimensional  
165 versions and adjusted the number of layers and convolutional kernels to make them suitable  
166 for 3-dimensional MRI inputs (e.g., GMD and GMV as different input channels). The present  
167 model consists of one stem module, three groups of convolutional modules  
168 (Inception-ResNet-A/B/C) and two reduction modules (Reduction-A/B). The model can take  
169 advantage of convolutional kernels with different shapes and sizes, and can extract features of  
170 different sizes. The model also can mitigate vanishing gradients and exploding gradients by  
171 adding residual modules. We utilized the Keras built-in stochastic gradient descent optimizer

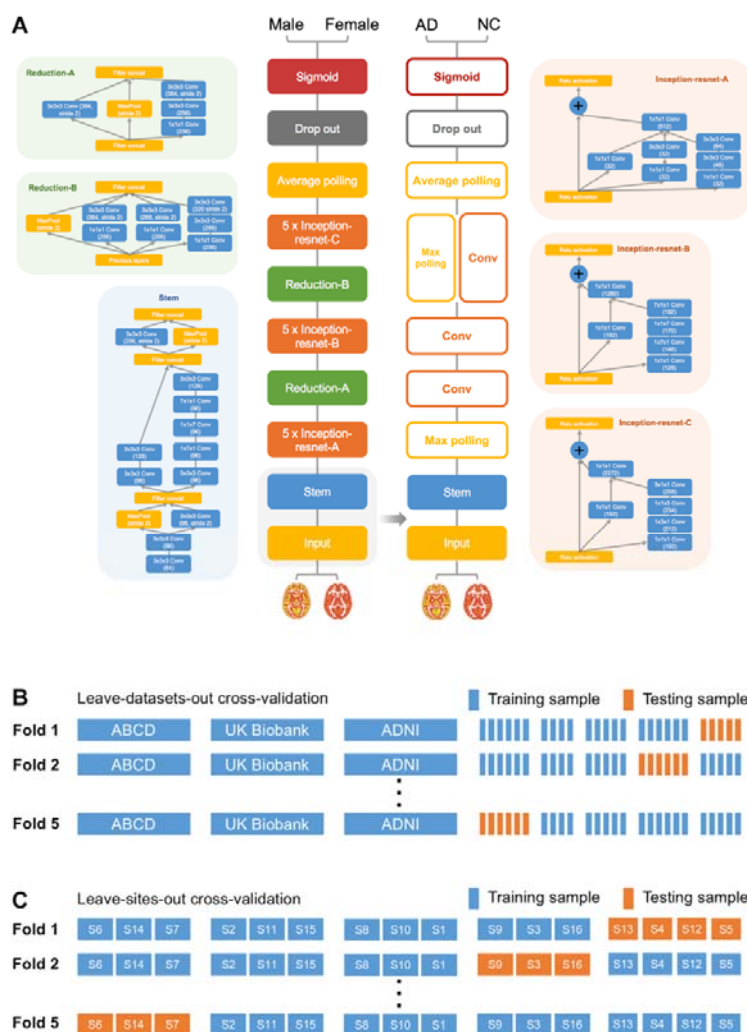
172 with learning rate = 0.01, Nesterov momentum = 0.9, decay = 0.003 (e.g., learn rate = learn  
173  $\text{rate}_0 \times (1 / (1 + \text{decay} \times \text{batch}))$ ). The loss function was set to binary cross-entropy. The batch  
174 size was set to 24 and the training procedure lasted 10 epochs for each fold. To avoid  
175 potential overfitting, we randomly split 600 samples out of the training sample as a validation  
176 sample and set a checking point at the end of every epoch. We saved the model in which the  
177 epoch classifier showed the lowest validation loss. Thereafter, the testing sample was fed into  
178 this model to test the classifier.

179

180 While training the sex classifier, random cross-validation may share participants from the  
181 same sites between training and testing samples, so the model may not generalize well to  
182 datasets from unseen sites due to the site information leakage during training. To ensure  
183 generalizability, we used cross-dataset validation. In the testing phase, all the data from a  
184 given dataset would never be seen during the classifier training phase. This also ensured the  
185 data from a given site (and thus a given scanner) were unseen by the classifier during training  
186 (see **Fig. 1B** for an illustration). This strict setting can limit classifier performance, but it  
187 makes it feasible to generalize to any participant at any site (scanner). Five-fold cross-dataset  
188 validation was used to assess classifier accuracy. Of note, 3 datasets were always kept in the  
189 training sample due to the massive number of samples: Adolescent Brain Cognition  
190 Development (ABCD) (n = 31,176), UK Biobank (n = 20,124) and the Alzheimer's Disease  
191 Neuroimaging Initiative (ADNI) (n = 16,596). The remaining 31 datasets were randomly  
192 allocated to the training and testing samples. The allocating schemas were the solution that  
193 balanced the sample size of 5 folds the best from 10,000 random allocating procedures.

194





195

196

197

198

199

200

201

202

203

204

205

206

**Fig. 1 | Flow diagram for the Alzheimer disease (AD) transfer learning framework and cross-validation procedure.** (A) Schema for the 3D Inception-ResNet-V2 model and the transfer learning framework for the Alzheimer disease classifier. (B) Schematic diagram for the leave-datasets-out 5-fold cross-validation for the sex classifier. (C) Schematic diagram for the leave-sites-out 5-fold cross-validation for the AD classifier.

**Transfer learning: classifier training and testing for AD.** After obtaining a highly robust and accurate brain imaging-based sex classifier as a base model, we used transfer learning to further fine-tune the AD classifier. Rather than retaining the intact sophisticated structure of the base model (Inception-ResNet-V2), we only leveraged the pre-trained weights in the **stem** module and simplified the upper layers (e.g., replacing Inception-ResNet modules with

207 ordinary convolutional layers). The retained bottom structure of the model works as a feature  
208 extractor and can take advantage of the massive training of sex classifier. And the pruned  
209 upper structure of the AD model can avoid potential overfitting and promote generalizability  
210 by reducing the number of parameters (10 million parameters for the AD classifier vs. 54  
211 million parameters for the sex classifier). This derived AD classifier was fine-tuned on the  
212 ADNI dataset (2,186 samples from 380 AD patients and 4,671 samples from 698 normal  
213 controls (NCs)). ADNI was launched in 2003 (Principal Investigator: Michael W. Weiner,  
214 MD) to investigate biological markers of the progression of MCI and early AD (see  
215 [www.adni-info.org](http://www.adni-info.org)). We used the Keras built-in stochastic gradient descent optimizer with  
216 learning rate = 0.0003, Nesterov momentum = 0.9, decay = 0.002. The loss function was set  
217 to binary cross-entropy. The batch size was set to 24 and the training procedure lasted 10  
218 epochs for each fold. Similar to the cross-dataset validation for the sex classifier training,  
219 five-fold cross-site validation was used to assess classifier accuracy (see **Fig. 1C** for an  
220 illustration). By ensuring that the data from a given site (and thus a given scanner) were  
221 unseen by the classifier during training, this strict strategy made the classifier generalizable  
222 with non-inflated accuracy, thus better simulating realistic clinical applications than  
223 traditional five-fold cross-validation.

224

225 Furthermore, to test the generalizability of the AD classifier, we directly tested the classifier  
226 on another unseen independent AD sample, i.e., Australian Imaging, Biomarker and Lifestyle  
227 Flagship Study of Ageing (AIBL)<sup>21</sup> and Open Access Series of Imaging Studies (OASIS)<sup>22,23</sup>.  
228 We used the averaged output of 5 AD classifiers in the five-fold cross-validation trained on  
229 ADNI as the final output for a participant. We used diagnoses provided by the AIBL dataset  
230 as the labels of samples (101 samples from 82 AD patients and 523 samples from 324 NCs).  
231 As OASIS did not specify the criteria for an AD diagnosis, we adopted 2 criteria from  
232 ADNI-1 to define AD patients, i.e., 1) mini-mental state examination score between 20 and  
233 26 (inclusive) and 2) clinical dementia rating score = 0.5 or 1.0. Thus, we tested on 277 AD  
234 samples and 995 normal control samples who met the ADNI-1 criteria for AD and NC in the  
235 OASIS dataset. Of note, AIBL and OASIS scanning conditions and recruitment criteria

236 differed much more than variations among different ADNI sites (where scanning and  
237 recruitment was deliberately coordinated), so we expected the AD classifier to achieve lower  
238 performance.

239

240 We further investigated whether the AD classifier could predict disease progression in people  
241 with MCI. MCI is a diagnosis defined as cognitive decline without impairment in everyday  
242 activities<sup>15</sup>. People with the amnesic subtype of MCI have a high risk of converting to AD.  
243 We screened imaging records of the MCI patients who converted to AD later in the ADNI  
244 1/2/'GO' phases, and collected 2,371 images from 243 participants labeled as 'pMCI' (i.e.,  
245 their early scans before entering the AD phase; the images labeled 'Conversion: MCI to AD'  
246 and images labeled as 'AD' after conversion were not included). We also assembled 4,018  
247 samples from 524 participants labeled 'sMCI' without later progression for contrast. We  
248 directly fed all these MCI images into the AD classifier without further fine-tuning, thus  
249 evaluating the performance of the AD classifier on unseen MCI information.

250

### 251 **Interpretation of the deep learning classifiers.**

252 To better understand the brain imaging-based deep learning classifier, we calculated  
253 occlusion maps for the classifiers. We repeatedly tested the images in testing sample using  
254 the model with the highest accuracy within the 5 folds, while successively masking brain  
255 areas (volume = 18mm\*18mm\*18mm, step = 9mm) of all input images. The accuracy  
256 achieved on "intact" samples by the classifier minus accuracy achieved on "defective"  
257 samples indicated the "importance" of the occluded brain area for the classifier. The  
258 occlusion maps were calculated for both sex and AD classifiers. To investigate the clinical  
259 significance of the output of the AD classifier, we calculated the Spearman's correlation  
260 coefficient between the predicted scores and mini-mental state examination (MMSE) scores  
261 of AD, NC and MCI samples. We also used general linear models (GLM) to verify whether  
262 the predicted scores (or MMSE score) showed a group difference between people with sMCI  
263 and pMCI. The age and sex information of MCI participants was included in this GLM as  
264 covariates. We selected the T1-weighted images from the first visit for each MCI subject and

265 finally collected data from 243 pMCI patients and 524 sMCI patients.

266

### 267 **3. Results**

#### 268 **3.1. Large-Scale Brain imaging Data**

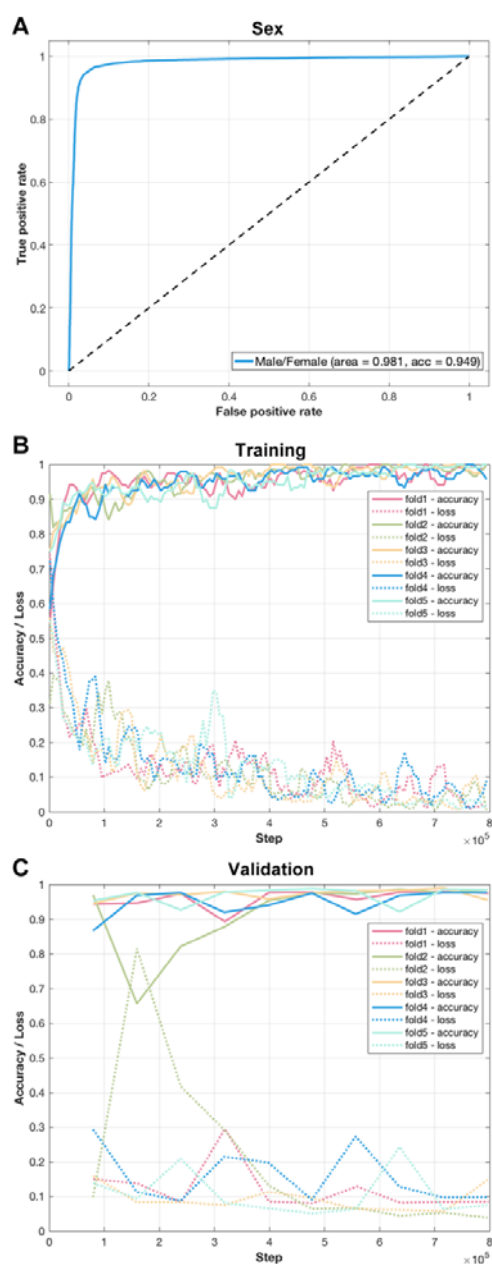
269 Only brain imaging data with enough size and variety can make deep learning accurate and  
270 robust enough to build a practical classifier. We received permissions from the administrators  
271 of 34 datasets (85,721 samples of 50,876 participants from more than 217 sites/scanners, see  
272 **Table S1**; there were no application requirements for some datasets). Data for each  
273 participant contained at least one session with a T1-weighted brain structural image and  
274 information on participant sex. After quality control, all these samples were used to pre-train  
275 the stem module to achieve better generalization for further AD classifier training. For the  
276 further fine-tuning of the AD classifier, ADNI (16,596 samples from 2,212 participants),  
277 AIBL (624 samples from 406 participants) and OASIS (3,150 samples from 1,664  
278 participants) were selected to train and test the model.

279

#### 280 **3.2. Performance of the sex classifier**

281 We trained a 3-dimensional Inception-ResNet-v2 model adapted from its 2-dimensional  
282 version in the Keras built-in application (see Fig. 1A for structure). As noted in the Methods,  
283 we did not feed raw data into the classifier for training, but used prior knowledge regarding  
284 helpful analytic pipelines. The brain structural data were segmented and normalized to yield  
285 grey matter density (GMD) and grey matter volume (GMV) maps (i.e., GMD and GMV  
286 maps were treated as different input channels). To ensure generalizability, five-fold  
287 cross-dataset validation was used to assess classifier accuracy. The five-fold cross-dataset  
288 validation accuracies were: 94.8%, 94.0%, 94.8%, 95.7% and 95.8%. Taken together,  
289 accuracy was 94.9% in testing samples when pooling results across the five folds. The area  
290 under the curve (AUC) of the receiver operating characteristic (ROC) curve reached 0.981  
291 (see **Fig. 2**). In short, our model can classify the sex of a participant based on brain structural  
292 imaging data from anyone and any scanner with an accuracy of about 95%. Interested readers  
293 can test this model on our online prediction website (<http://brainimagenet.org>).

294



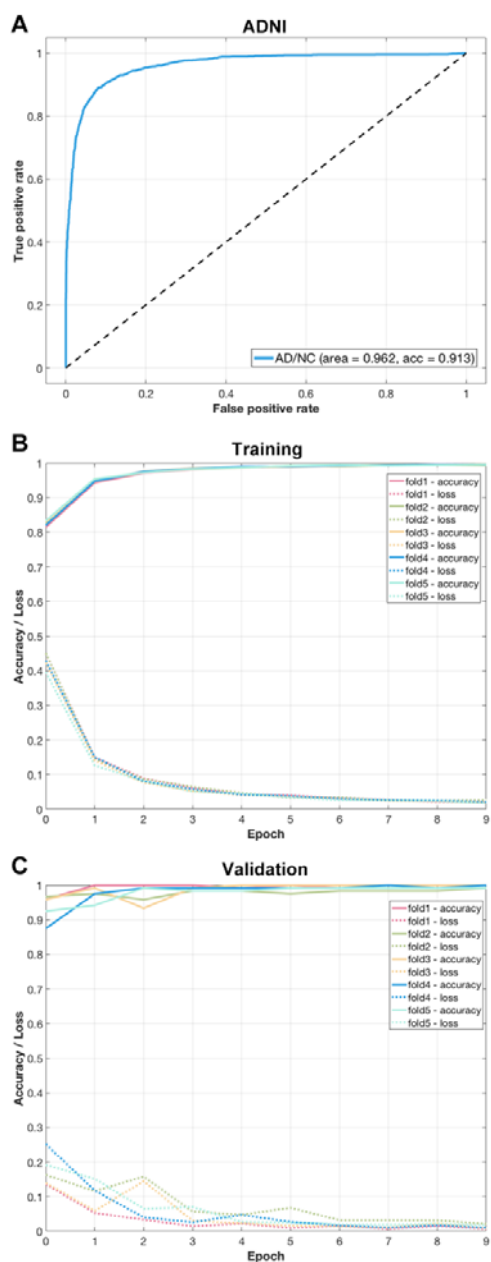
295

296 **Fig. 2 | Performance of the sex classifier.** (A) The receiver operating characteristic curve of  
297 the sex classifier. (B) The tensorboard monitor graph of the sex classifier in the training  
298 sample. The curve was smoothed for better visualization. (C) The tensorboard monitor graph  
299 of the sex classifier in the validation sample.

300

301 **3.3. Performance of the AD classifier**

302 After creating a practical brain imaging-based classifier for sex with high cross-dataset  
303 accuracy, we used transfer learning to see if we could classify patients with AD. The AD  
304 classifier achieved an average accuracy of 91.3% (accuracy = 93.2%, 90.3%, 92.0%, 94.4%  
305 and 86.7% in 5 cross-site folds) in the test samples. Average sensitivity and specificity were  
306 0.848 and 0.943, respectively. The ROC AUC reached 0.962 when results from the 5 testing  
307 samples were taken together (see **Fig. 3** and **Table 1**). The AD classifier achieved an average  
308 accuracy of 91.4% in 3T field strength MR testing samples and achieved an average accuracy  
309 of 91.1% in 1.5T MR testing samples. The accuracy in 3T MR testing sample was not  
310 significantly different from that of 1.5T MR testing sample ( $p = 0.316$ , statistical examined  
311 by permutation test of randomly allocating the testing samples into 1.5T group or 3T group  
312 and calculated the accuracy difference between the two groups for 100,000 times, see **Fig.**  
313 **S7**).



314

315 **Fig. 3 | Performance of the Alzheimer's disease (AD) classifier. (A)** The receiver  
316 operating characteristic curve of the AD classifier. **(B)** The tensorboard monitor panel of the  
317 AD classifier in the training sample. **(C)** The tensorboard monitor panel of the AD classifier  
318 in the validation sample.

319

320 To test the generalizability of the AD classifier, we applied it to unseen independent AD  
321 datasets, i.e., AIBL and OASIS 1 and 2. The AD classifier achieved 94.2% accuracy in AIBL

322 with 0.97 AUC (see **Fig. 4A**). Sensitivity and specificity were 0.881 and 0.954, respectively.  
323 The AD classifier achieved 87.9% accuracy in OASIS with 0.936 AUC (see **Fig. 4B**).  
324 Sensitivity and specificity were 0.796 and 0.902, respectively.

325

326 **Table 1 | Performance of the Alzheimer's disease classifier**

Dataset	n (AD)	n (NC)	Accuracy	AUC	Sensitivity	Specificity
ADNI	2,186	4,671	0.913	0.962	0.848	0.943
AIBL	101	523	0.942	0.97	0.881	0.954
OASIS	277	995	0.879	0.936	0.796	0.902

327 **AD = Alzheimer's disease; NC = normal control. The sample sizes shown here are the numbers of T1-weighted**  
328 **brain MRI scans.**

329

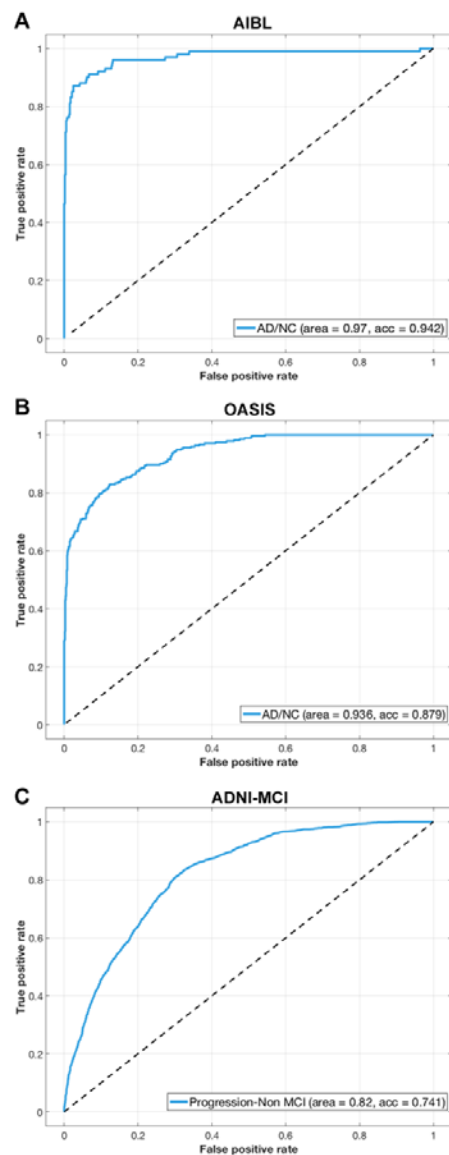
330 Importantly, although the AD classifier is agnostic to brain imaging data of MCI, we directly  
331 tested it on the MCI dataset in ADNI to see if it has the potential to predict the progression of  
332 MCI to AD. The idea behind this test is that even though people with MCI do not yet have  
333 AD, their scans may appear closer to the AD class learned by the deep learning model. In the  
334 end, 65.2% of pMCI patients were predicted as closer to the AD class but only 20.4% of  
335 sMCI patients were predicted as having AD by the AD classifier. If the percentage of pMCI  
336 patients who were predicted as AD was considered as sensitivity and the percentage of sMCI  
337 patients who were predicted as AD was considered as 1-specificity, the AUC of ROC curve  
338 for AD classifier reached 0.82. These results suggest that the classifier is practical for  
339 screening MCI patients who have a higher risk of progression to AD. In sum, we believe our  
340 AD classifier can provide important insights relevant to computer-aided diagnosis and  
341 prediction of AD, and we have freely provided it on the website <http://brainimagenet.org>.  
342 Importantly, classification results by the online classifier should be interpreted with caution,  
343 as they cannot replace diagnosis by licensed clinicians.

344

345



346



347

348 **Fig. 4 | Receiver operating characteristic (ROC) curves for the Alzheimer disease (AD)**  
349 **classifier when tested on independent AD samples and a mild cognitive impairment**  
350 **sample. (A) The ROC curve of AD classifier tested on the AIBL sample. (B) The ROC curve**  
351 **of AD classifier tested on the OASIS sample. (C) The ROC curve of AD classifier tested on**  
352 **MCI sample in ADNI. The images of MCI subjects with future conversion to AD were**  
353 **labeled as “AD”, and the images of MCI subjects who had not shown conversion to AD were**  
354 **labeled as “NC”.**

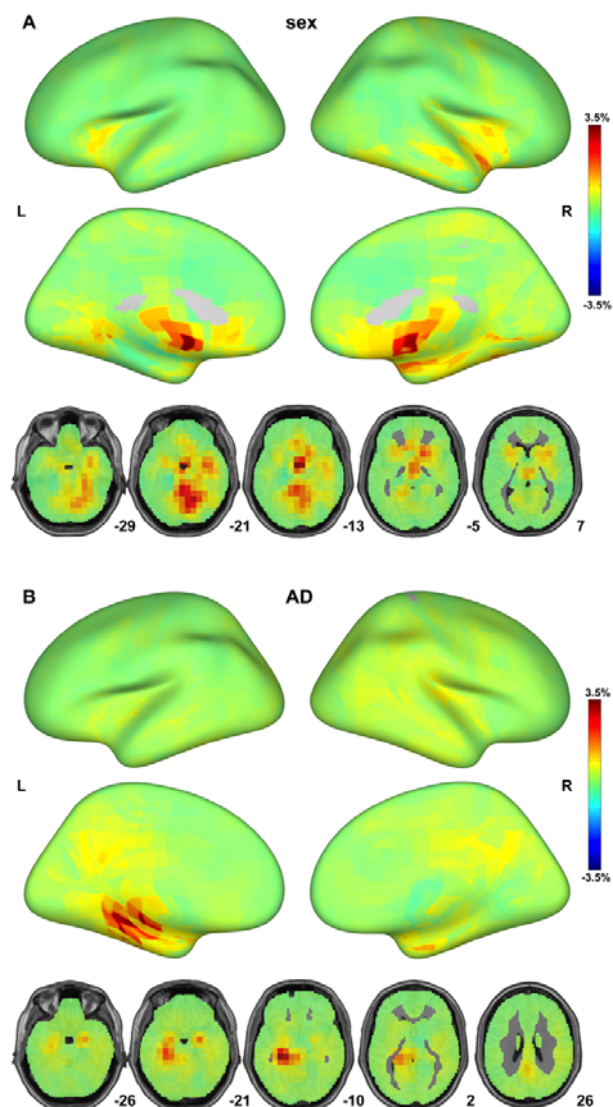
355

356 As a supplementary analysis, we also trained the AD classifier that kept the intact structure of  
357 the base model in transfer learning. The performance of the proposed model was  
358 comprehensively inferior to the optimized AD classifier. The “intact” AD classifier achieved  
359 an average accuracy of 88.4% with 0.938 AUC in the ADNI test samples (see **Fig. S2A**).  
360 Average sensitivity and specificity were 0.814 and 0.917, respectively. When tested on  
361 independent samples, the AD classifier achieved 91.2% accuracy in AIBL with 0.948 AUC  
362 (see **Fig. S3A**). Sensitivity and specificity were 0.851 and 0.924, respectively. The AD  
363 classifier achieved 86.1% accuracy in OASIS with 0.921 AUC (see Fig. S3B). Sensitivity and  
364 specificity were 0.789 and 0.881, respectively. When tested on MCI samples, 63.2% of pMCI  
365 patients were predicted as having AD and only 22.1% of sMCI patients was predicted as  
366 having AD by the AD classifier (see **Fig. S3C**).

367

### 368 **3.5. Interpretation of the deep learning classifiers**

369 To better understand the brain imaging-based deep learning classifier, we calculated  
370 occlusion maps for the classifiers. In brief, we continuously set a cubic brain area of every  
371 input image to zeros, and made the classifier attempt classification based on the defective  
372 samples. The occlusion map showed that hypothalamus, superior vermis, thalamus, amygdala,  
373 putamen, accumbens, hippocampus and parahippocampal gyrus played critical roles in  
374 predicting sex (see **Fig. 5A**). The occlusion map for the AD classifier highlighted that the  
375 hippocampus and parahippocampal gyrus - especially in the left hemisphere - played unique  
376 roles in predicting AD (see **Fig. 5B**).



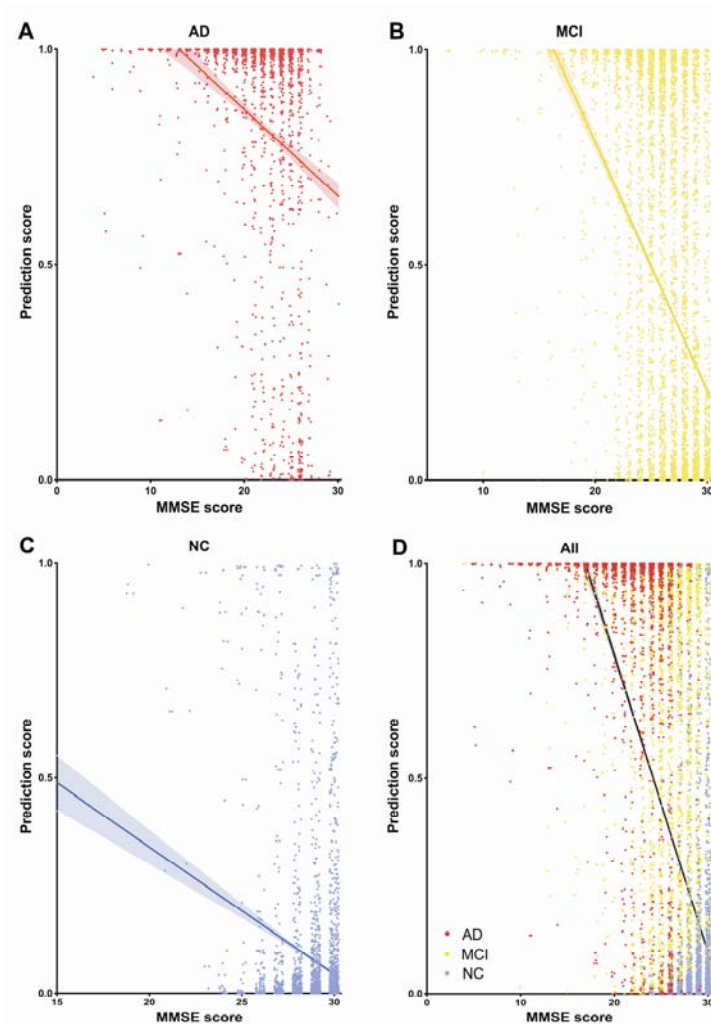
377

378 **Fig. 5 | Interpretation of the deep learning classifiers with occlusion maps.** Classifier  
379 performance dropped considerably when the brain areas rendered in red were masked out of  
380 the model input. **(A)** The occlusion maps for the sex classifier. **(B)** The occlusion maps for  
381 Alzheimer disease classifier.

382

383 To investigate the clinical significance of the output of the AD classifier, we calculated the  
384 Spearman's correlation coefficient between the predicted scores by the classifier and  
385 mini-mental state examination (MMSE, provided by ADNI datasets) scores in AD, NC and  
386 MCI samples, although the classifier had not been trained for MMSE scores before. This  
387 analysis confirmed significant negative correlations between the predicted scores and MMSE

388 scores for AD ( $r = -0.37, p < 1 \times 10^{-55}$ ), NC ( $r = -0.11, p < 1 \times 10^{-11}$ ), MCI ( $r = -0.52, p < 1 \times$   
389  $10^{-307}$ ) and the overall samples ( $r = -0.64, p < 1 \times 10^{-307}$ ) (See Fig. 6). As lower MMSE scores  
390 indicate more severe cognitive impairment for AD and MCI patients, we confirmed that the  
391 more severe the disease, the higher the predicted score by the classifier. In addition, both the  
392 predicted scores and MMSE scores showed significant differences between pMCI and sMCI  
393 (predicted scores:  $t = 13.88, p < 0.001$ , Cohen's  $d = 1.08$ ; MMSE scores:  $t = -9.42, p < 0.01$ ,  
394 Cohen's  $d = -0.73$ , See **Fig. S5**). Importantly, the effect size of the predicted scores by the  
395 classifier is much larger than the behavioral measure (MMSE scores).



396

397 **Fig. 6 | Correlations between the output of the Alzheimer's disease (AD) classifier and**  
398 **the severity of illness.** The predicted scores from the AD classifier showed significant  
399 negative correlations with the mini-mental state examination (MMSE) scores of AD, normal

400 control (NC) and mild cognitive impairment (MCI) samples. **(A)** Correlation between the  
401 predicted scores from AD classifier and the MMSE scores of AD samples. **(B)** Correlation  
402 between the predicted scores from the AD classifier and the MMSE scores of MCI samples.  
403 **(C)** Correlation between the predicted scores from AD classifier and the MMSE scores of NC  
404 samples. **(D)** Correlation between the predicted scores from AD classifier and the MMSE  
405 scores of AD, NC, and MCI samples.

406

#### 407 **4. Discussion**

408 Using an unprecedentedly diverse brain imaging sample, we pre-trained an industrial-grade  
409 sex classifier with about 95% accuracy which served as a base-model for transfer learning to  
410 promote model generalization capability. After transfer learning, the model fine-tuned to AD  
411 achieved 91.3% accuracy in stringent leave-sites-out cross-validation and 94.2%/87.9%  
412 accuracy for direct tests on unseen independent datasets. Predicted scores from the AD  
413 classifier showed significant negative correlations with the severity of illness. The AD  
414 classifier also showed the potential to predict the prognosis of MCI patients.

415

416 The industrial-grade high accuracy and generalization capability of our deep neural network  
417 classifiers demonstrate that brain imaging did have practical utility for auxiliary diagnosis.  
418 The current prototype may facilitate future research to apply brain imaging in many practical  
419 application fields. Of note, the output of the deep neural network model is a continuous  
420 variable, so the threshold can be adjusted to balance sensitivity and specificity. For example,  
421 when tested on the independent sample (OASIS), sensitivity and specificity results were  
422 0.796 and 0.902, respectively, as the default threshold was set at 0.5. However, for screening,  
423 the false-negative rate should be minimized even at the cost of higher false-positive rates. If  
424 we lower the threshold (e.g., to 0.2), sensitivity can be improved to 0.881 at a cost of  
425 decreasing specificity to 0.796. Thus, in our freely available AD prediction website, users can  
426 obtain continuous outputs and adjust the threshold by themselves. This adjustable  
427 characteristic of the present model makes itself more suitable for integration into the current

428 diagnostic criteria as a diagnostic biomarker. The proposed MRI-based biomarker has a high  
429 sensitivity, which may address the lack of sensitivity of other biomarkers.

430

431 Except for the feasibility of being integrated into diagnostic criteria, the present AD model  
432 also showed outstanding characteristics as a progression biomarker. First, the present model  
433 was able to quantify key disease milestones by predicting disease progression in MCI patients.  
434 In fact, people with pMCI were 3 times more likely to be classified as AD than sMCI (65.2%  
435 vs 20.4%). Recently, a critical review about predicting the progression of MCI noted that  
436 about 40% of studies had methodological issues, such as lack of a test dataset, data-leakage in  
437 feature selection or parameter tuning, and leave-one-out validation performance bias<sup>24</sup>. The  
438 present AD classifier was only trained on AD/NC samples and was not fine-tuned using MCI  
439 data, so data leakage was avoided. The estimated true AUC of current published state-of-art  
440 classifiers for predicting progression of MCI is about 0.75<sup>13,24</sup>. The proposed AD classifier  
441 here outperformed the benchmark considerably (AUC = 0.82). Considering the discouraging  
442 clinical trial failures for AD treatments, early identification of people with MCI with potential  
443 to progress would help in evaluation of early treatments<sup>25</sup>. Clinicians can use the present AD  
444 classifier as auxiliary decision support system. Second, the output of the deep neural network  
445 can indicate the clinical severity for AD patients or people with MCI, as the predicted scores  
446 showed significant negative correlations with MMSE scores. Considering the “greedy”  
447 characteristic of CNN for reducing training loss, the prediction scores for AD and NC were  
448 overstated, so the magnitude of negative correlations might be even underestimated. Third,  
449 when directly comparing the predicted scores (or MMSE scores) between sMCI and pMCI  
450 groups, the predicted scores showed much higher effect size than MMSE scores (Cohen’s  
451  $d_{\text{prediction}} = 1.08$  vs Cohen’s  $d_{\text{MMSE}} = -0.73$ ), indicating that predicted scores may offer better  
452 prompting or ‘warning’ effects for the physician to differentiate MCI patients.

453

454 Although deep-learning algorithms have often been described as “black boxes” for their poor  
455 interpretability, our subsequent analyses showed that the current MRI-based AD biomarker  
456 was in line with former pathological findings and clinical practices. For example, AD induced

457 brain structural changes have been frequently reported by MRI studies. Among all the  
458 structural findings, hippocampal atrophy is the most prominent change and is used in imaging  
459 assisted diagnosis<sup>26</sup>. Neurobiological changes in the hippocampus typically precede  
460 progressive neocortical damage and AD symptoms. The convergence of our deep learning  
461 system and human physicians on alterations in hippocampal structure for classifying AD  
462 patients is in line with the crucial role of the hippocampus in AD. Furthermore, brain atrophy  
463 in AD has been frequently reported as left lateralized<sup>27,28</sup>. Compared to the un-optimized AD  
464 classifier, a slight left hemisphere preference for input features may help explain the  
465 improved performance of the optimized AD classifier (see **Fig. 5** and **Fig. S4** to compare the  
466 occlusion maps).

467

468 Rather than indiscriminately imitate the structure of the base model in transfer learning, the  
469 present AD classifier significantly simplified the model before the fine-tuning procedure. In  
470 fact, the performance of the unoptimized AD classifier was far poorer than that of the  
471 optimized AD classifier in accuracy, sensitivity, specificity, and in independent validation  
472 performance (see **Fig. S2-3**). There is some reported evidence that truncating or pruning  
473 models before transfer learning may facilitate the performance of the transferred models<sup>29,30</sup>.  
474 As the sample for training the AD classifier is considerably smaller than that used to train the  
475 sex classifier, the simplified model structure may help to avoid overfitting and improve  
476 generalizability.

477

478 By precisely predicting the sex of people, the present study also advances our understanding  
479 of sex differences in human brain. Daphna and colleagues extracted hundreds of VBM  
480 features from structural MRI and concluded that “the so-called male/female brain” does not  
481 exist as no individual structural feature supports a sexually dimorphic view of human brains<sup>31</sup>.  
482 However, human brains may embody sexually dimorphic features in a multivariate manner.  
483 The high accuracy and generalizability of the present sex classifier demonstrate that sex is  
484 separable in a 1,981,440-dimension (96\*120\*86\*2) feature space. Among those 1,981,440  
485 features, hypothalamus played the most critical role in predicting sex. The hypothalamus



486 regulates testosterone secretion through the hypothalamic-pituitary-gonadal axis and thus  
487 plays a critical role in brain masculinization<sup>32</sup>. Men have significantly larger hypothalamus  
488 than women relative to cerebrum size<sup>33</sup>. Taken together, our machine learning evidence  
489 shows that the “male/female brain” does exist, in the sense that accurate classification is  
490 possible.

491

492 In the deep learning field, the appearance of ImageNet tremendously accelerated the  
493 evolution of computer vision<sup>34</sup>. It provided large amounts of well-labeled image data for  
494 researchers to pre-train their models. Studies have shown that pre-trained models can  
495 facilitate the performance and robustness of subsequently fine-tuned models<sup>35</sup>. The present  
496 study confirms that the “pre-train + fine-tuning” paradigm does work for MRI-based  
497 auxiliary diagnosis. Unfortunately, no such well-preprocessed dataset exists in brain imaging  
498 domain. As data organization and preprocessing of MRI data require tremendous time,  
499 manpower and computational load, these constraints impede scientists from other fields  
500 entering brain imaging. Open access to large amounts of preprocessed brain imaging data is  
501 fundamental to facilitate the participation of a broader range of researchers. Beyond building  
502 and sharing a practical brain imaging-based deep learning classifier, we would openly share  
503 all sharable preprocessed data to invite researchers (especially computer scientists) to join the  
504 efforts to create predictive models using brain images (Link\_To\_Be\_Added upon publication,  
505 preprocessed data of some datasets could not be shared as the raw data owners do not allow  
506 sharing of data derivatives). We anticipate that this dataset may boost the clinical utility of  
507 brain imaging as ImageNet has done in computer vision research. We openly share our  
508 models to allow other researchers to deploy them  
509 (<https://github.com/Chaogan-Yan/BrainImageNet>). Our code is also openly shared as well  
510 (<https://github.com/Chaogan-Yan/BrainImageNet>), allowing other researchers to replicate the  
511 present results and further develop brain imaging-based classifiers based on our existing  
512 work. Finally, we have also built a demonstration website for classifying sex and AD  
513 (<http://brainimagenet.org>). Users can upload their own raw T1-weighted or preprocessed  
514 GMD and GMV data to make predictions of sex or AD labels in real-time.



515

516 Some limitations of the current study should be acknowledged. Considering the lower  
517 reproducibility of functional MRI compared to structural MRI, only structural MRI derived  
518 images were used in the present deep learning model. Even so, functional measures of  
519 physiology and activation may further improve the performance of sex and brain disorder  
520 classifiers. In future studies, functional MRI, especially resting-state functional MRI, may  
521 provide additional information for model training. Furthermore, with advances in software  
522 such as FreeSurfer<sup>36</sup>, fmriprep<sup>37</sup> and DPABISurf, surface-based algorithms have shown their  
523 superiority when compared with traditional volume-based algorithms<sup>38</sup>. Surface-based  
524 algorithms are more time consuming to run in terms of computation load, but can provide  
525 more precise brain registration and reproducibility. Future studies should take surface-based  
526 images as inputs for deep learning models. In addition, the present AD classification model  
527 was built based on labels provided by ADNI database. Further study may also benefit by  
528 using post-mortem neuropathological data as a gold standard for AD to further advance the  
529 clinical value of MRI-based biomarkers.

530

531 In summary, we pooled MRI data from more than 217 sites/scanners to constitute the largest  
532 brain MRI sample to date, and applied a state-of-the-art architecture deep convolutional  
533 neural network, Inception-ResNet-V2, to pre-train an industrial-grade brain image-based  
534 classifier. The AD classifier obtained via transfer learning reached high accuracy and  
535 sufficient generalizability to be of practical use, demonstrating the feasibility of transfer  
536 learning in brain disorder applications. Future work is needed to deploy such a framework for  
537 assessment of psychiatric disorders, to predict treatment response, and other aspects of  
538 individual differences.

539

#### 540 **Data and code availability**

541 The imaging, phenotype and clinical data used for the training, validation and test sets were  
542 applied from the administrators of 34 datasets. The preprocessed brain imaging data will be  
543 available on (Link\_To\_Be\_Added upon publication, preprocessed data of some datasets

544 could not be shared as the raw data owners do not allow sharing data derivatives). The code  
545 for training and testing the model are openly shared at  
546 <https://github.com/Chaogan-Yan/BrainImageNet>.

547

548

549 **Acknowledgement**

550 Data used in the preparation of this article for training and testing the sex classifier was  
551 obtained from the Adolescent Brain Cognitive Development (ABCD) Study  
552 (<https://abcdstudy.org>), held in the NIMH Data Archive (NDA). This is a multisite,  
553 longitudinal study designed to recruit more than 10,000 children age 9-10 and follow them  
554 over 10 years into early adulthood. The ABCD Study is supported by the National Institutes  
555 of Health and additional federal partners under award numbers U01DA041048,  
556 U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037,  
557 U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028,  
558 U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025,  
559 U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089. A full list  
560 of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of  
561 participating sites and a complete listing of the study investigators can be found at  
562 <https://abcdstudy.org/scientists/workgroups/>. ABCD consortium investigators designed and  
563 implemented the study and/or provided data but did not necessarily participate in analysis or  
564 writing of this report. This manuscript reflects the views of the authors and may not reflect  
565 the opinions or views of the NIH or ABCD consortium investigators. This research has been  
566 conducted using the UK Biobank Resource. Data collection and sharing for the training and  
567 testing the sex and AD classifier were funded by the Alzheimer's Disease Neuroimaging  
568 Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI  
569 (Department of Defense award number W81XWH-12-2-0012). ADNI is funded by the  
570 National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering,  
571 and through generous contributions from the following: AbbVie, Alzheimer's Association;  
572 Alzheimer's Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen;  
573 Bristol-Myers Squibb Company; CereSpir, Inc.; Cogstate; Eisai Inc.; Elan Pharmaceuticals,  
574 Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated  
575 company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer  
576 Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical  
577 Research & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale

578 Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals  
579 Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and  
580 Transition Therapeutics. The Canadian Institutes of Health Research is providing funds to  
581 support ADNI clinical sites in Canada. Private sector contributions are facilitated by the  
582 Foundation for the National Institutes of Health ([www.fnih.org](http://www.fnih.org)). The grantee organization is  
583 the Northern California Institute for Research and Education, and the study is coordinated by  
584 the Alzheimer's Therapeutic Research Institute at the University of Southern California.  
585 ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of  
586 Southern California.

587

### 588 **Funding**

589 This work was supported by the National Key R&D Program of China (grant number:  
590 2017YFC1309902), the National Natural Science Foundation of China (grant number:  
591 81671774, 81630031), the 13th Five-year Informatization Plan of Chinese Academy of  
592 Sciences (grant number: XXH13505), the Key Research Program of the Chinese Academy of  
593 Sciences (grant NO. ZDBS-SSW-JSC006), Beijing Nova Program of Science and  
594 Technology (grant number: Z191100001119104).

595

### 596 **Competing interests**

597 The authors declare no competing interests.

598

### 599 **Supplementary material**

600 Supplementary material is available online.

601

### 602 **References**

- 603 1 Dubois, B. *et al.* Advancing research diagnostic criteria for Alzheimer's disease: the IWG-2 criteria. *The*  
604 *Lancet Neurology* **13**, 614-629, (2014).
- 605 2 Jack Jr, C. R. *et al.* Introduction to the recommendations from the National Institute on  
606 Aging-Alzheimer's Association workgroups on diagnostic guidelines for Alzheimer's disease.  
607 *Alzheimer's & dementia* **7**, 257-262, (2011).

608 3 de Jager, C. A., Honey, T. E., Birks, J. & Wilcock, G. K. Retrospective evaluation of revised criteria for the  
609 diagnosis of Alzheimer's disease using a cohort with post-mortem diagnosis. *Int. J. Geriatr. Psychiatry*  
610 **25**, 988-997, (2010).

611 4 Harris, J. M. *et al.* Do NIA-AA criteria distinguish Alzheimer's disease from frontotemporal dementia?  
612 *Alzheimer's & Dementia* **11**, 207-215, (2015).

613 5 Ham, Y.-G., Kim, J.-H. & Luo, J.-J. Deep learning for multi-year ENSO forecasts. *Nature* **573**, 568-572,  
614 (2019).

615 6 DeVries, P. M. R., Viegas, F., Wattenberg, M. & Meade, B. J. Deep learning of aftershock patterns  
616 following large earthquakes. *Nature* **560**, 632-634, (2018).

617 7 Liu, W. *et al.* A survey of deep neural network architectures and their applications. *Neurocomputing*  
618 **234**, 11-26, (2017).

619 8 Kermany, D. S. *et al.* Identifying medical diagnoses and treatable diseases by image-based deep  
620 learning. *Cell* **172**, 1122-1131, (2018).

621 9 Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural networks. *Nature*  
622 **542**, 115-118, (2017).

623 10 McKinney, S. M. *et al.* International evaluation of an AI system for breast cancer screening. *Nature* **577**,  
624 89-94, (2020).

625 11 Qiu, S. *et al.* Development and validation of an interpretable deep learning framework for Alzheimer's  
626 disease classification. *Brain* **143**, 1920-1933, (2020).

627 12 Bashyam, V. M. *et al.* MRI signatures of brain age and disease over the lifespan based on a deep brain  
628 network and 14 468 individuals worldwide. *Brain* **143**, 2312-2324, (2020).

629 13 Rathore, S., Habes, M., Iftikhar, M. A., Shacklett, A. & Davatzikos, C. A review on neuroimaging-based  
630 classification studies and associated feature extraction methods for Alzheimer's disease and its  
631 prodromal stages. *NeuroImage* **155**, 530-548, (2017).

632 14 Yosinski, J., Clune, J., Bengio, Y. & Lipson, H. in *Adv. Neural Inf. Process. Syst.* 3320-3328.  
633 15 Gauthier, S. *et al.* Mild cognitive impairment. *The lancet* **367**, 1262-1270, (2006).

634 16 Yan, C. G. & Zang, Y. F. DPARSF: A MATLAB Toolbox for "Pipeline" Data Analysis of Resting-State fMRI.  
635 *Front. Syst. Neurosci.* **4**, 13, (2010).

636 17 Friston, K. J. *et al.* Statistical parametric maps in functional imaging: a general linear approach. *Hum.*  
637 *Brain Mapp.* **2**, 189-210, (1994).

638 18 Goto, M. *et al.* Diffeomorphic Anatomical Registration Through Exponentiated Lie Algebra provides  
639 reduced effect of scanner for cortex volumetry with atlas-based method in healthy subjects.  
640 *Neuroradiology* **55**, 869-875, (2013).

641 19 Good, C. D. *et al.* A voxel-based morphometric study of ageing in 465 normal adult human brains.  
642 *NeuroImage* **14**, 21-36, (2001).

643 20 Szegedy, C., Ioffe, S., Vanhoucke, V. & Alemi, A. A. in *National Conference on Artificial Intelligence.*  
644 4278-4284.

645 21 Ellis, K. A. *et al.* Addressing population aging and Alzheimer's disease through the Australian Imaging  
646 Biomarkers and Lifestyle study: Collaboration with the Alzheimer's Disease Neuroimaging Initiative.  
647 *Alzheimer's & dementia* **6**, 291-296, (2010).

648 22 Marcus, D. S. *et al.* Open Access Series of Imaging Studies (OASIS): cross-sectional MRI data in young,  
649 middle aged, nondemented, and demented older adults. *J. Cogn. Neurosci.* **19**, 1498-1507, (2007).

650 23 Marcus, D. S., Fotenos, A. F., Csernansky, J. G., Morris, J. C. & Buckner, R. L. Open access series of

651 imaging studies: longitudinal MRI data in nondemented and demented older adults. *J. Cogn. Neurosci.*  
652 **22**, 2677-2684, (2010).

653 24 Ansart, M. *et al.* Predicting the progression of mild cognitive impairment using machine learning: A  
654 systematic, quantitative and critical review. *Med. Image Anal.* **67**, 101848, (2021).

655 25 Selkoe, D. J. Preventing Alzheimer's disease. *Science* **337**, 1488-1492, (2012).

656 26 Frisoni, G. B., Fox, N. C., Jack, C. R., Jr., Scheltens, P. & Thompson, P. M. The clinical use of structural  
657 MRI in Alzheimer disease. *Nat. Rev. Neurol.* **6**, 67-77, (2010).

658 27 Wachinger, C., Salat, D. H., Weiner, M., Reuter, M. & Initiative, A. s. D. N. Whole-brain analysis reveals  
659 increased neuroanatomical asymmetries in dementia for hippocampus and amygdala. *Brain* **139**,  
660 3253-3266, (2016).

661 28 Derflinger, S. *et al.* Grey-matter atrophy in Alzheimer's disease is asymmetric but not lateralized.  
662 *Journal of Alzheimer's Disease* **25**, 347-357, (2011).

663 29 Liu, J., Wang, Y. & Qiao, Y. in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence.*  
664 2245-2251.

665 30 Ke, A., Ellsworth, W., Banerjee, O., Ng, A. Y. & Rajpurkar, P. CheXtransfer: Performance and Parameter  
666 Efficiency of ImageNet Models for Chest X-Ray Interpretation. *arXiv preprint arXiv:2101.06871*, (2021).

667 31 Joel, D. *et al.* Sex beyond the genitalia: The human brain mosaic. *Proc. Natl. Acad. Sci. U. S. A.* **112**,  
668 15468-15473, (2015).

669 32 Forest, M. G., Peretti, E. D. & Bertrand, J. Hypothalamic-pituitary-gonadal relationships in man from  
670 birth to puberty. *Clin. Endocrinol. (Oxf.)* **5**, 551-569, (1976).

671 33 Makris, N. *et al.* Volumetric parcellation methodology of the human hypothalamus in neuroimaging:  
672 Normative data and sex differences. *NeuroImage* **69**, 1-10, (2013).

673 34 Deng, J. *et al.* in *2009 IEEE conference on computer vision and pattern recognition.* 248-255 (Ieee).

674 35 Hendrycks, D., Lee, K. & Mazeika, M. Using pre-training can improve model robustness and uncertainty.  
675 *arXiv preprint arXiv:1901.09960*, (2019).

676 36 Fischl, B. FreeSurfer. *NeuroImage* **62**, 774-781, (2012).

677 37 Esteban, O. *et al.* fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nat. Med.* **16**, 111-116,  
678 (2019).

679 38 Coalson, T. S., Van Essen, D. C. & Glasser, M. F. The impact of traditional neuroimaging methods on the  
680 spatial localization of cortical areas. *Proc. Natl. Acad. Sci. U. S. A.* **115**, e6356-e6365, (2018).

681