

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/150037/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Zheng, Guowei, Zhang, Yu, Zhao, Ziyang, Wang, Yin, Liu, Xia, Shang, Yingying, Cong, Zhaoyang, Dimitriadis, Stavros I. , Yao, Zhijun and Hu, Bin 2022. A transformer-based multi-features fusion model for prediction of conversion in mild cognitive impairment. *Methods* 205 , pp. 241-248. 10.1016/j.ymeth.2022.04.015

Publishers page: <http://dx.doi.org/10.1016/j.ymeth.2022.04.015>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



A Transformer-based Multi-features Fusion Model for Prediction of Conversion in Mild Cognitive Impairment

Guowei Zheng^a, Yu Zhang^a, Ziyang Zhao^a, Yin Wang^a, Xia Liu^e, Yingying Shang^a, Zhaoyang Cong^a, Stavros Dimitriadis^{f, g, h, i, j, k, l, m, n, *}, Zhijun Yao^{a, *} and Bin Hu^{a, b, c, d, *}

^a Gansu Provincial Key Laboratory of Wearable Computing, School of Information Science and Engineering, Lanzhou University, Lanzhou, China

^b CAS Center for Excellence in Brain Science and Intelligence Technology, Shanghai Institutes for Biological Sciences, Chinese Academy of Sciences, Shanghai, China

^c Joint Research Center for Cognitive Neurosensor Technology of Lanzhou University & Institute of Semiconductors, Chinese Academy of Sciences, Lanzhou, China

^d Engineering Research Center of Open Source Software and Real-Time System (Lanzhou University), Ministry of Education, Lanzhou, China

^e School of Computer Science, Qinghai Normal University, Xining, China.

^f Integrative Neuroimaging Lab, 55133, Thessaloniki (Macedonia), Greece.

^g Neuroinformatics Group, Cardiff University Brain Research Imaging Centre (CUBRIC), School of Psychology, -College of Biomedical and Life Sciences, Cardiff, Wales, United Kingdom.

^h 1st Department of Neurology, G.H. "AHEPA " School of Medicine, Faculty of Health Sciences, Aristotle University of Thessaloniki (AUTH), Thessaloniki, Greece.

ⁱ Greek Association of Alzheimer's Disease and Related Disorders, Thessaloniki, Macedonia, Greece.

^j Cardiff University Brain Research Imaging Centre (CUBRIC), School of Psychology, College of Biomedical and Life Sciences, Cardiff University, Cardiff, Wales, United Kingdom.

^k Division of Psychological Medicine and Clinical Neurosciences, School of Medicine, College of Biomedical and Life Sciences, Cardiff University, Cardiff, Wales, United Kingdom.

^l School of Psychology, College of Biomedical and Life Sciences, Cardiff University, Cardiff, Wales, United Kingdom.

^m Neuroscience and Mental Health Research Institute, School of Medicine, College of Biomedical and Life Sciences, Cardiff University, Cardiff, Wales, United Kingdom.

29 ⁿ MRC Centre for Neuropsychiatric Genetics and Genomics, School of Medicine, College of
30 Biomedical and Life Sciences, Cardiff University, Cardiff, Wales, United Kingdom.
31 * Correspondence: DimitriadisS@cardiff.ac.uk (Stavros Dimitriadis), yaozj@lzu.edu.cn (Zhijun
32 Yao), bh@lzu.edu.cn (Bin Hu)
33 Email: zhenggw20@lzu.edu.cn (Guowei Zheng), yzhang20@lzu.edu.cn (Yu Zhang),
34 zhaozy2021@lzu.edu.cn (Ziyang Zhao), wangyin20@lzu.edu.cn (Yin Wang), liux2016@lzu.edu.cn
35 (Xia Liu), shangyy20@lzu.edu.cn (Yingying Shang), congchy19@lzu.edu.cn (Zhaoyang Cong),
36 DimitriadisS@cardiff.ac.uk (Stavros Dimitriadis), yaozj@lzu.edu.cn (Zhijun Yao), bh@lzu.edu.cn
37 (Bin Hu)

38 Abstract

39 Mild cognitive impairment (MCI) is usually considered the early stage of Alzheimer’s disease
40 (AD). Therefore, the accurate identification of MCI individuals with high risk in converting to AD is
41 essential for the potential prevention and treatment of AD. Recently, the great success of deep
42 learning has sparked interest in applying deep learning to neuroimaging field. However, deep
43 learning techniques are prone to overfitting since available neuroimaging datasets are not
44 sufficiently large. Therefore, we proposed a deep learning model fusing cortical features to address
45 the issue of fusion and classification blocks. To validate the effectiveness of the proposed model,
46 we compared seven different models on the same dataset in the literature. The results show that
47 our proposed model outperformed the competing models in the prediction of MCI conversion with
48 an accuracy of 83.3% in the testing dataset. Subsequently, we used deep learning to characterize
49 the contribution of brain regions and different cortical features to MCI progression. The results
50 revealed that the caudal anterior cingulate and pars orbitalis contributed most to the classification
51 task, and our model pays more attention to volume features and cortical thickness features.

52 **Keywords:** Mild cognitive impairment, Magnetic resonance imaging, Deep learning, Transformer

53 1 Introduction

54 Alzheimer’s disease (AD) is a common degenerative disease in aging populations. Cognitive

55 impairment and progressive memory loss are the fundamental characteristics of AD [1]. More than
56 30 million people worldwide are suffering from AD cause of the extending life expectancy, and this
57 number is estimated to be tripled by 2050 [2]. Despite the dramatic increase in the prevalence of
58 AD, no treatment can completely cure it currently. Thus, early diagnosis is crucial to developing
59 treatments for AD [3, 4]. Mild cognitive impairment (MCI) is generally considered a transitional
60 stage between normal aging and AD [5]. Studies have shown that approximately 5% to 15% of
61 persons with MCI will progress to AD each year [6, 7]. MCI can be divided into two subtypes,
62 progressive mild cognitive impairment (pMCI) and stable mild cognitive impairment (sMCI).
63 Subjects classified as pMCI were those with a higher risk of conversion to AD in a short period,
64 while subjects in the sMCI group remained stable for a certain period and had a lower risk of
65 progression to AD than the former [8]. Therefore, classifying the two different types of MCI can
66 predict the conversion from MCI to AD as early as possible, which is beneficial for AD prevention
67 and therapy.

68 Neuroimaging is widely used to understand the pathology of MCI and AD [9]. In previous
69 studies on the mechanism of AD, structural magnetic resonance imaging (MRI) is one of the most
70 extensively utilized imaging modalities in AD detection and prediction for its wide practicality, non-
71 invasion, high resolution, and moderate cost [10]. Applying machine learning techniques to
72 neuroimaging diagnosis is a developing field. In terms of MCI conversion prediction, numerous
73 studies are using different methods, including network features constructed based on graph theory
74 [11, 12], voxel-based morphometry (VBM) based on the segmentation of grey matter [13, 14],
75 multiple methods of hippocampal segmentation [15], etc. However, research using traditional
76 machine learning methods still suffers from inadequacies. The performance of traditional machine
77 learning methods largely depends on data representation [16], and it is challenging to learn high-
78 level information from poorly hand-picked features.

79 Recently, with the development of deep learning technology, many researchers have achieved
80 outstanding achievements in neuroscience [17-19]. Deep learning network models also progressed
81 in predicting AD conversion in advance from MCI [20-23]. Nevertheless, most of these studies used
82 3D subject-level features as input to deep learning network models, which suffer from overfitting
83 issues, since the sample size of available neuroimaging data sets is not significant compared with
84 millions of features in each image [24, 25]. Freesurfer is a powerful tool to reliably extract cortical

85 features such as volume, surface area, cortical thickness, and curvature index [26-33] through an
86 automated pipeline without any user interaction. The dimension of cortical features is significantly
87 lower than the original neuroimage but contains rich ROI-level brain morphological information,
88 which can effectively alleviate the overfitting problem. In 2017, the transformer was first proposed
89 by Vaswani et al. [34] and successfully applied to natural language processing (NLP) tasks.
90 Researchers have recently extended it to other tasks such as computer vision (CV) with great
91 success [35]. Its strong global perception capability makes it possible to find differences in brain
92 morphology of the cortex between pMCI and sMCI from fused cortical features for classification.

93 Based on the above considerations, in this work, we proposed a transformer-based multi-
94 features fusion model to predict conversion in MCI by using MRI. Specifically, our architecture was
95 designed to fuse the multiple cortical features and automatically learn high-level information from
96 the fused features. To validate the proposed model, we perform the classification on the MRI
97 datasets from the Alzheimer’s Disease Neuroimaging Initiative (ADNI, <http://adni.loni.usc.edu/>)
98 [36], and achieved better performance over other models. Furthermore, with occlusion analysis,
99 we investigated the contribution of different brain regions and different cortical features to the
100 classifying progression and stability of MCI.

101 The rest of the paper has been organized as follows. In Section 2, we mainly introduce the
102 architecture of the proposed model and the details of its construction and validation and the
103 implementation of occlusion analysis. Section 3 gives the analysis of the results followed by further
104 discussion in Section 4. Finally, Section 5 summarizes the full text.

105 **2 Materials and methods**

106 **2.1 Experimental Data**

107 Data used in our study were obtained from the ADNI database. The ADNI is an ongoing and
108 multicenter study that aims to develop imaging, clinical, genetic, and biochemical biomarkers for
109 AD's early detection and tracking [37]. 249 MCI participants with baseline T1-weighted structural
110 MRI were selected from ADNI in this work. All MCI subjects were divided into two groups: (1) stable
111 mild cognitive impairment (sMCI) who did not convert to AD within three years. In addition, the

112 subjects who were diagnosed as MCI at least twice, but reverse to a standard control, at last, are
113 also considered as sMCI [23]; (2) progressive mild cognitive impairment (pMCI) who were
114 diagnosed as MCI at the first visit, but converted to AD at longitudinal visits within three years. The
115 detailed demographic information is given in Table 1.

116 **TABLE 1** The demographic information.

	sMCI	pMCI
Subjects' number	104	145
Age range	55-88	55-88
Males/Females	67/37	90/55

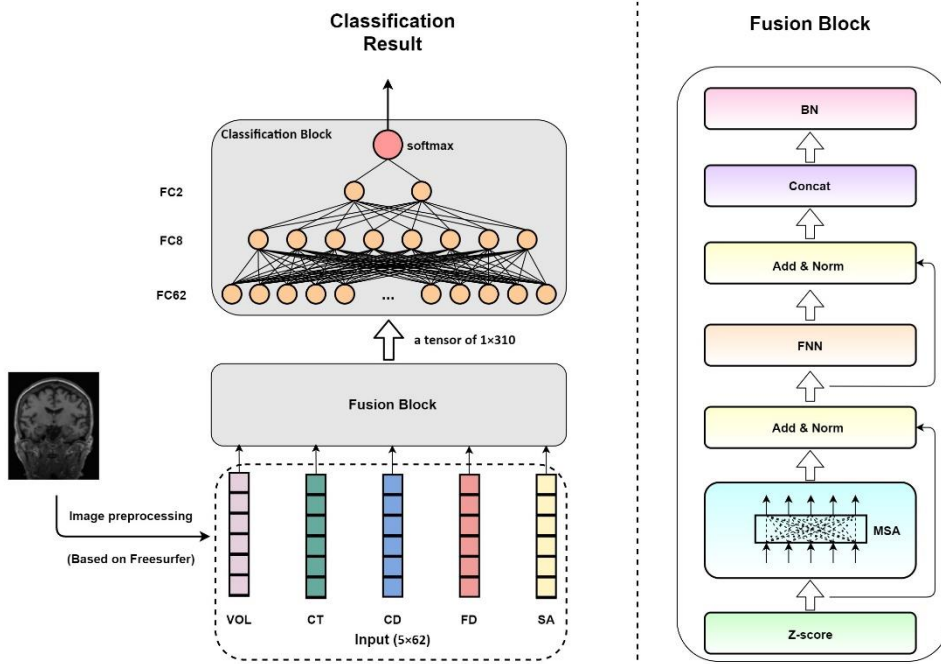
117 Abbreviations: pMCI = progressive mild cognitive impairment, sMCI = stable mild cognitive
118 impairment.

119 **2.2 Image Pre-processing**

120 T1-weighted structural images were processed using the Freesurfer software (v6.0;
121 <http://surfer.nmr.mgh.harvard.edu/>) [38]. The preprocessing steps are described below. Firstly, the
122 correction for non-uniformity artifacts was performed on the images [39], followed by the
123 coordinate transformation [40] and the brain tissue segmentation (including gray matter, white
124 matter, cerebrospinal fluid, and other background categories). Subsequently, the surface of
125 white/gray matter boundaries was reconstructed [40]. After completing the construction of
126 boundary models, surface expansion and registration were performed [30, 38]. Finally, we
127 extracted multiple cortical measurements including volume (VOL), cortical thickness (CT),
128 curvature index (CD), folding index (FD), and surface area (SA) for 62 brain regions (31 regions in
129 each hemisphere of the brain) using the Desikan-Killiany-Tourville (DKT) atlas [41].

130 **2.3 The Transformer-based Multi-features Fusion Model**

131 Here we proposed a multi-features fusion model to predict conversion in MCI, which is based
132 on the transformer model [34]. Our model was designed to input a cortical feature matrix
133 (extracted from the preprocessed image) and output the classification result. The model consists
134 of a fusion block and a classification block. For an overview, refer to Fig.1.



135

136 **Fig. 1.** Illustration of proposed deep learning model. Abbreviations: VOL = volume, CT = cortical
 137 thickness, CD = curvature index, FD = folding index, SA = surface area, MSA = multi-head self-
 138 attention, FFN = feed-forward network, Concat = concatenate, BN = Batch Normalization, FC62 =
 139 62-units fully connected layer, FC8 = 8-units fully connected layer, FC2 = 2-units fully connected
 140 layer.

141 The fusion block consists of five different sub-layers. Firstly, the Z-score method was applied
 142 to the input features to remove the effect of different feature sizes (Eq. (1)). The second is a multi-
 143 head self-attention (MSA) [34] and the third is a feed-forward network (FFN), the residual
 144 connections [42] were employed after the MSA and FFN, followed by layer normalization (LN) [43]
 145 (Eqs. (2) and (3)). To improve the model's efficiency, we set the number of heads in the MSA to 2,
 146 which could reduce the number of model parameters. The dimension of outputs for MSA and FFN
 147 is 62, which matches the model's input and enables these residual connections. The fourth is a
 148 concatenate (Concat) layer (Eq. (4)), which reshapes the input data (cortical feature matrix, 5×62)
 149 to a tensor of 1×310 for later classification. Finally, the Batch Normalization (BN) layers were
 150 applied to accelerate convergence. The output of fusion block f is calculated using Eq. (5) (x is
 151 the input of the model). Then, took f as the input to the classification block.

152
$$l_1 = Z - score(x) \tag{1}$$

153
$$l_2 = LN(l_1 + MSA(l_1)) \tag{2}$$

154
$$l_3 = LN(l_2 + FNN(l_2)) \tag{3}$$

155
$$l_4 = \text{Concat}(l_3) \quad (4)$$

156
$$f = \text{BN}(l_4) \quad (5)$$

157 The classification block consists of three fully connected (FC) layers with 62, 8, and 2 units
158 respectively. Later, the softmax activation function was used to predict the results.

159 **2.4 Implementation**

160 We implemented our model with Pytorch 1.8.0. Model training and testing were performed
161 on the Ubuntu 18.04 operating system. During training, we used the Binary Cross-Entropy (BCE)
162 loss function and set the number of epochs to 15, with a mini-batch size of 32. The optimizer was
163 Adam [44] with a learning rate of 1e-4 and weight decay of 1e-8.

164 **2.5 Validation Framework**

165 To validate the efficacy of the proposed model, we split our 249 subjects randomly into three
166 groups, including the training dataset (n=200), the validation dataset (n=25), and the testing
167 dataset (n=24). The training dataset was used for training models, while the validation dataset was
168 used for parameter tuning and the testing dataset for evaluating model performance.

169 **2.6 Model Comparison**

170 The proposed model was compared with four traditional machine learning methods: support
171 vector machine [45], decision tree [46], random forest [47], and logistic regression [48]. Compared
172 to traditional machine learning methods with feature engineering, deep learning models aim to
173 extract features automatically. Therefore, deep learning methods including Recurrent Neural
174 Network (RNN) [49], Long Short-Term Memory (LSTM) [50], and Gated Recurrent Unit (GRU) [51]
175 were also employed in this study for comparison with the proposed model.

176 To verify the performance of the above models, we randomly divide the entire dataset into a
177 training dataset, a validation dataset, and a testing dataset in a ratio of 8:1:1 (the division of the
178 dataset is the same as our proposed model). All the traditional machine learning models were
179 implemented using sklearn (<https://scikit-learn.org/stable/>) library in python3 (used the default
180 settings) and were trained on the training dataset then tested on the testing dataset. In addition,
181 all deep learning models were implemented with Pytorch 1.8.0, and for these deep learning

182 models, we defined the optimal hyperparameters of the classifiers by using the training and
 183 validating datasets. Subsequently, when the model achieved the best performance (the optimized
 184 hyperparameters were listed in *Supplementary Materials* Table S1) in the validation dataset, the
 185 model was validated using the testing dataset. We also performed 10-fold cross-validation on the
 186 entire dataset to compare the performance between our proposed and other models to ensure
 187 generalizability.

188 Furthermore, to validate the performance of all the models, we employed four
 189 measurements including classification accuracy (ACC), sensitivity (SEN), specificity (SPE) (shown
 190 in Eq. (6) to Eq. (8)), and the area under the receiver operating characteristic (ROC) curve (AUC).
 191 For these measurements, higher values demonstrate better performance.

$$192 \quad ACC = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

$$193 \quad SEN = \frac{TP}{TP + FN} \quad (7)$$

$$194 \quad SPE = \frac{TN}{TN + FP} \quad (8)$$

195 Where TP, TN, FP, and FN are abbreviations for True Positive, True Negative, False Positive, and
 196 False Negative, respectively.

197 **2.7 Implementation of Occlusion Analysis**

198 Occlusion analysis was employed to investigate the contribution of each brain region and each
 199 cortical feature to the performance of the proposed model. First, we set the value of five cortical
 200 features of each brain region (both left and right) to 0 from the cortical feature matrix of the test
 201 stage and retested the trained proposed model. The input corresponding to brain region m is x_m :

$$202 \quad BrainRegionOcc_n = \begin{cases} 0 & \text{if } m=n \\ x_m & \text{otherwise} \end{cases} \quad (9)$$

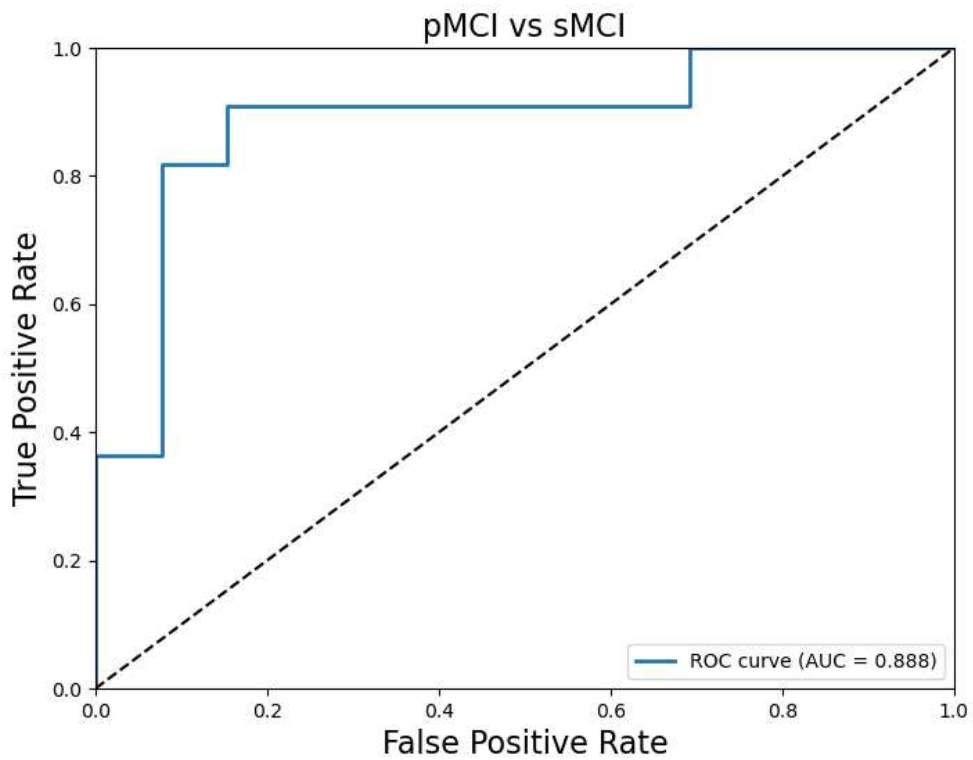
203 Where $BrainRegionOcc_n$ represents the occlusion of brain region n . If the m is equal to n ,
 204 the value of x_m is set to 0. Then, we masked different features to explore their impact on the
 205 model. See Eq. (10) for details (the $FeatureOcc_i$ means the occlusion of i -th feature and x_j
 206 means the input corresponding j -th feature).

$$207 \quad FeatureOcc_i = \begin{cases} 0 & \text{if } j=i \\ x_j & \text{otherwise} \end{cases} \quad (10)$$

208 **3 Results**

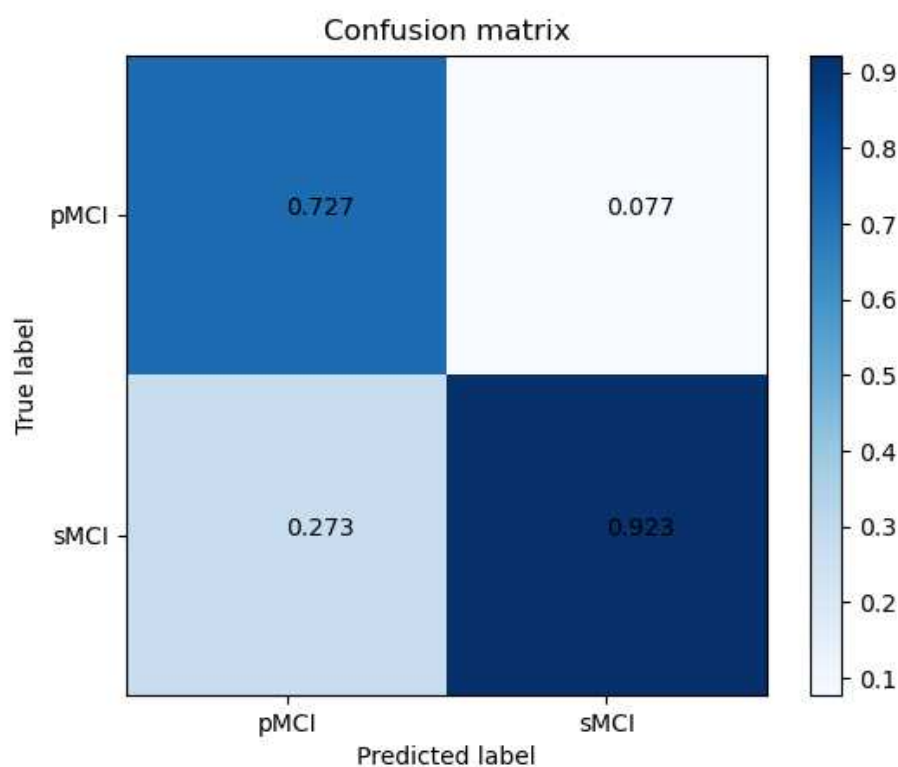
209 **3.1 Classification Performance**

210 We test the trained model on the testing dataset to verify the proposed model. The
211 proposed model achieved 83.3% accuracy and an AUC of 0.888 (Fig.2), with a sensitivity of 0.727
212 and a specificity of 0.923(Fig.3).



213

214 **Fig. 2.** ROC curve for classifying pMCI versus sMCI. Abbreviations: pMCI = progressive mild cognitive
215 impairment, sMCI = stable mild cognitive impairment, ROC = receiver operating characteristic, AUC
216 = area under the curve.



217

218 **Fig. 3.** Confusion matrix, evaluating the SEN and SPE obtained in pMCI versus sMCI. The matrix
 219 values were rescaled to the scope of [0,1]. Abbreviations: pMCI = progressive mild cognitive
 220 impairment, sMCI = stable mild cognitive impairment.

221 In addition, we compared our proposed method with three different deep learning methods
 222 (RNN, LSTM, GRU) and four different machine learning methods (random forest, decision tree,
 223 logistic regression, and support vector machine). As shown in Table 2, our proposed model
 224 showed better performance than the other models. The results show that GRU (AUC = 0.853)
 225 performs better than LSTM (AUC = 0.839) and RNN (AUC = 0.790) among the three deep learning
 226 models. Furthermore, the random forest has the best performance (AUC=0.678) among four
 227 machine learning models. See Table 2 for more detailed information. The results of 10-fold cross-
 228 validation also showed that our model can predict MCI conversion more accurately (see
 229 *Supplementary Materials Table S2*).

230 **TABLE 2** The performance of different models.

Model	ACC	SEN	SPE	AUC
Proposed model	83.3%	0.727	0.923	0.888
RNN	66.7%	0.818	0.538	0.790
LSTM	70.8%	0.727	0.692	0.839
GRU	70.8%	0.727	0.692	0.853

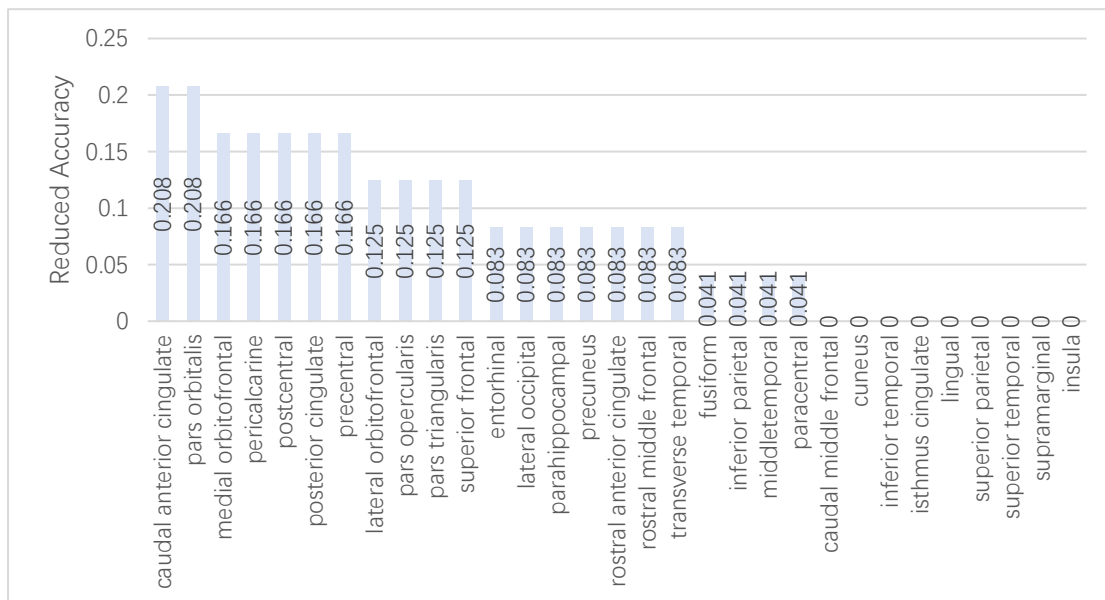
Random Forest	66.7%	0.818	0.538	0.678
Decision Tree	54.2%	0.545	0.538	0.542
Logistic Regression	66.7%	0.727	0.615	0.671
Support vector machine	54.2%	0.818	0.308	0.563

231 The best results for each column are shown in boldface. Abbreviations: RNN = Recurrent Neural
 232 Network, LSTM = Long Short-Term Memory, GRU = Gated Recurrent Unit, ACC = accuracy, SEN =
 233 sensitivity, SPE = specificity, AUC = area under the receiver operating characteristic (ROC) curve.

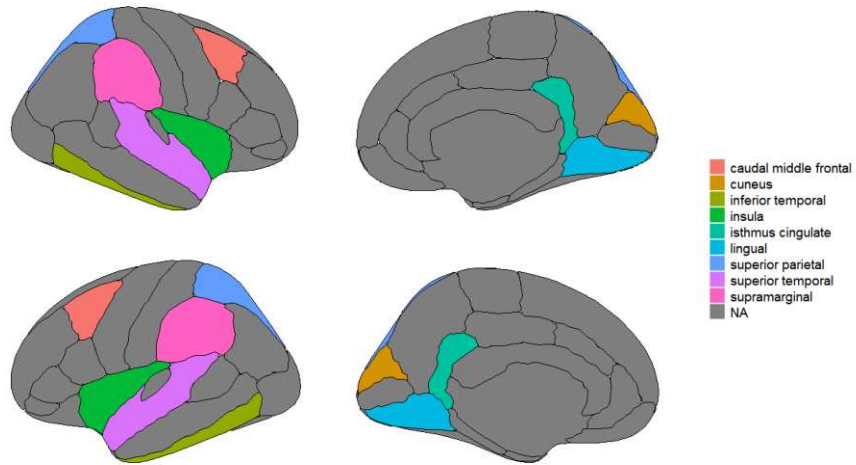
234 3.2 Occlusion Analysis

235 After extracting five features based on the DKT template for each subject, the relevant
 236 contribution of different brain regions and different features to the classification performance
 237 was evaluated using computer vision's commonly used occlusion analysis method [52].

238 As shown from Fig.4, the masking of most brain regions causes a decrease in model
 239 accuracy, and the masking of a small number of brain regions does not affect model accuracy
 240 (Fig.5). Notably, masking of the caudal anterior cingulate and the pars orbitalis (Fig.6) resulted in
 241 a dramatic decrease in model performance. Then, we performed occlusion analysis for different
 242 features. It can be seen from Fig.7 that the occlusion of different features all caused a significant
 243 decrease in model accuracy.

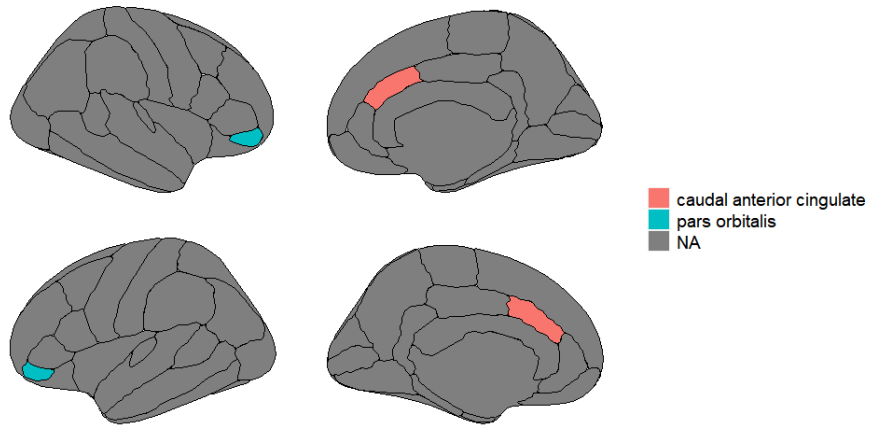


244
 245 **Fig. 4.** The reduced accuracies with each brain region occluded compared to the original intact
 246 model.



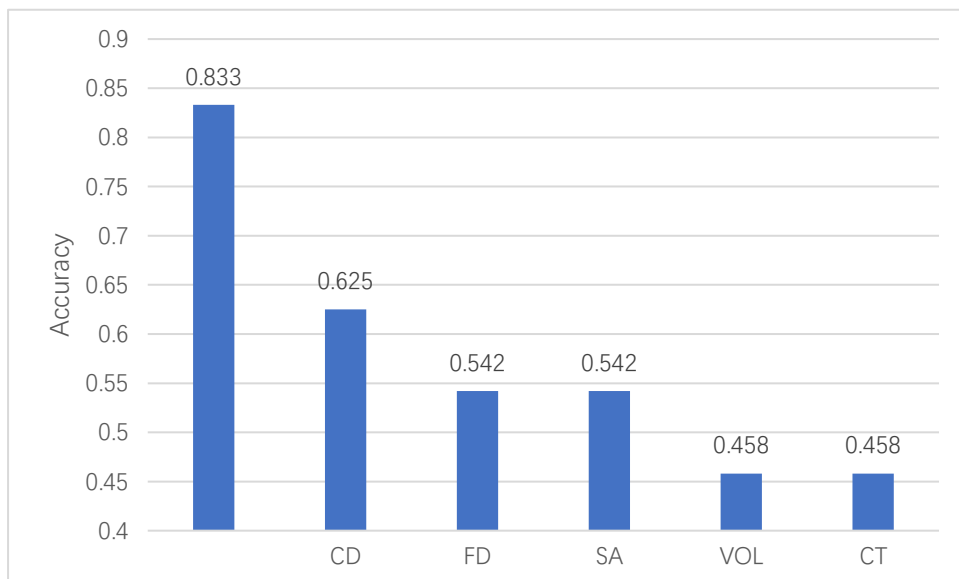
247
248
249

Fig. 5. The location of brain regions that do not affect classification accuracy.



250
251
252

Fig. 6. The location of the brain regions that contributed most to the classification task.



253

254 **Fig. 7.** Results for each feature occlusion (the first column is the model's accuracy with all
255 features input). Abbreviations: CD = curvature index, FD = folding index, SA = surface area, VOL =
256 volume, CT = cortical thickness.

257 **4 Discussion**

258 Patients with MCI show a strong variable trajectory of symptoms, with some individuals finally
259 diagnosed with AD, while others show a more stable cognitive ability pattern for a certain period.
260 Identifying these two different types of MCI is crucial and essential to preventing and treating AD.
261 Therefore, many researchers are committed to developing computer-aided systems to diagnose
262 AD early. To solve the overfitting problem of most previous methods, we proposed a transformer-
263 based model that predicts conversion in MCI using multiple ROI-level cortical features and achieved
264 an accuracy of 83.3% on the testing dataset.

265 The model comparison results demonstrated that the proposed model performs better than
266 other traditional machine learning models (random forest, decision tree, logistic regression, and
267 support vector machine) and deep learning models (RNN, LSTM, and GRU). The traditional machine
268 learning methods rely on the manual selection of features. For features that have not been
269 carefully selected, the traditional machine learning methods are challenging to thoroughly learn
270 sufficient information in cortical features. In addition, compared with other deep learning methods,
271 the proposed model includes a fusion block with MSA, which takes into account the features
272 themselves and fully considers the correlation between different cortical features to achieve better
273 performance. Furthermore, the classification performance of the proposed model also
274 outperformed previously developed deep learning models for classifying pMCI versus sMCI based
275 on MRI data [22, 53-55], which ranged from 73.95% to 78.79%.

276 The occlusion analysis results both extend and support prior reports by describing the
277 contribution of different brain regions and different cortical features to the progression of MCI. On
278 the one hand, the results revealed significant differences between the brain regions differentiating
279 pMCI from sMCI. Notably, the results have shown that the caudal anterior cingulate and pars
280 orbitalis (Fig.6) were most important for the classification task than any other brain region.
281 Previous studies have shown that neuronal loss in the caudal anterior cingulate begins in the early
282 stage of AD [56], and this timing may need to be advanced. This brain region contributed the most

283 to the model, possibly indicating that some neurons have been lost in pMCI. Given that the caudal
284 anterior cingulate is important to cognitive control of behavior [57], it suggests that pMCI may
285 show more severe cognitive impairment than sMCI. In addition, a previous study found that the
286 pars orbitalis, as well as some other brain regions, contributed to good classification performance
287 in this task [58], but the central role of the pars orbitalis should be highlighted. On the other hand,
288 the occlusion of different features all caused a significant decrease in model accuracy, this finding
289 demonstrates the existence of important complementary information in all five features.
290 Furthermore, the occlusion analysis on different features showed that VOL and CT had the
291 strongest impact on model performance, this may be related to the different volumes and atrophy
292 rates between sMCI and pMCI [59] and the significantly thinner cortical thicknesses in many brain
293 regions in pMCI [60]. The results also indicated that VOL and CT were more distinct in sMCI and
294 pMCI brains than CD, FD, and SA and were more reliable biomarkers in the progression of MCI.

295 Our study has some limitations. Firstly, our work only used MRI images, while researchers
296 have continuously disclosed the strength of multimodal features in computer-aided diagnosis
297 models [61-63]. Therefore, the model performance is expected to be improved by incorporating
298 data from multiple modalities, such as functional MRI. In addition, the cross-sectional nature is
299 another limitation of our study. Therefore, longitudinal data should be employed in our future
300 research.

301 **5 Conclusion**

302 This study proposed a transformer-based multi-features fusion model to predict the MCI-to-
303 AD conversion only using MRI data. Results show that our model can fuse the cortical features
304 extracted by Freesurfer. Compared with other models in the literature, our proposed model
305 achieves higher accuracy and AUC. In addition, our study reveals the contribution of brain regions
306 in differentiating between pMCI and sMCI, highlighting the central role of the caudal anterior
307 cingulate and pars orbitalis. Finally, the occlusion analysis results demonstrate that VOL and CT
308 may be more reliable biomarkers in MCI progression.

309

310 **Authors' Contributions:**

311 In this paper, Zhijun Yao and Bin Hu conceived the project. Guowei Zheng completed all
312 experiments in this work. Guowei Zheng wrote the manuscript, Stavros Dimitriadis, Yu Zhang,
313 Ziyang Zhao, Yin Wang, Xia Liu, Yingying Shang, and Zhaoyang Cong revised the manuscript.

314 Acknowledgments

315 This work was supported in part by the National Key Research and Development Program of
316 China (Grant No. 2019YFA0706200), in part by the National Natural Science Foundation of China
317 (Grant No.61632014, No.61627808), in part by the National Basic Research Program of China (973
318 Program, Grant No. 2014CB744600), in part by the Natural Science Foundation of Gansu Province
319 of China (Grant No.20JR5RA292), in part by the Fundamental Research Funds for the Central
320 Universities (Grant No. lzujbky-2018-125), and in part by the Department of education of Gansu
321 Province: "Innovation Star" project for excellent postgraduates (2021CXZX-121).

322 References

- 323 [1] J. Wen, E. Thibeau-Sutre, M. Diaz-Melo, J. Samper-González, A. Routier, S. Bottani, D. Dormont, S.
324 Durrleman, N. Burgos, O. Colliot, Convolutional neural networks for classification of Alzheimer's disease:
325 Overview and reproducible evaluation, *Medical image analysis* 63 (2020) 101694.
- 326 [2] D.E. Barnes, K. Yaffe, The projected effect of risk factor reduction on Alzheimer's disease prevalence,
327 *The Lancet Neurology* 10(9) (2011) 819-828.
- 328 [3] C. Samaey, A. Schreurs, S. Stroobants, D. Balschun, Early cognitive and behavioral deficits in mouse
329 models for tauopathy and Alzheimer's disease, *Frontiers in aging neuroscience* 11 (2019) 335.
- 330 [4] W. Shao, S. Xiang, Z. Zhang, K. Huang, J. Zhang, Hyper-graph based sparse canonical correlation
331 analysis for the diagnosis of Alzheimer's disease from multi-dimensional genomic data, *Methods* 189
332 (2021) 86-94.
- 333 [5] R.C. Petersen, Mild cognitive impairment as a diagnostic entity, *Journal of internal medicine* 256(3)
334 (2004) 183-194.
- 335 [6] A.J. Mitchell, M. Shiri-Feshki, Rate of progression of mild cognitive impairment to dementia—meta-
336 analysis of 41 robust inception cohort studies, *Acta psychiatrica scandinavica* 119(4) (2009) 252-265.
- 337 [7] R. Roberts, D.S. Knopman, Classification and epidemiology of MCI, *Clinics in geriatric medicine* 29(4)
338 (2013) 753-772.
- 339 [8] A.M. Anter, Y. Wei, J. Su, Y. Yuan, B. Lei, G. Duan, W. Mai, X. Nong, B. Yu, C. Li, A robust swarm
340 intelligence-based feature selection model for neuro-fuzzy recognition of mild cognitive impairment
341 from resting-state fMRI, *Information Sciences* 503 (2019) 670-687.
- 342 [9] Z. Li, H.-I. Suk, D. Shen, L. Li, Sparse multi-response tensor regression for Alzheimer's disease study

343 with multivariate clinical assessments, *IEEE transactions on medical imaging* 35(8) (2016) 1927-1936.

344 [10] T. Tong, Q. Gao, R. Guerrero, C. Ledig, L. Chen, D. Rueckert, A.s.D.N. Initiative, A novel grading
345 biomarker for the prediction of conversion from mild cognitive impairment to Alzheimer's disease, *IEEE*
346 *Transactions on Biomedical Engineering* 64(1) (2016) 155-165.

347 [11] S.H. Hojjati, A. Ebrahimzadeh, A. Khazaei, A. Babajani-Feremi, A.s.D.N. Initiative, Predicting
348 conversion from MCI to AD using resting-state fMRI, graph theoretical approach and SVM, *Journal of*
349 *neuroscience methods* 282 (2017) 69-80.

350 [12] R. Wei, C. Li, N. Fogelson, L. Li, Prediction of conversion from mild cognitive impairment to
351 Alzheimer's Disease using MRI and structural network features, *Frontiers in aging neuroscience* 8 (2016)
352 76.

353 [13] E. Moradi, A. Pepe, C. Gaser, H. Huttunen, J. Tohka, A.s.D.N. Initiative, Machine learning framework
354 for early MRI-based Alzheimer's conversion prediction in MCI subjects, *Neuroimage* 104 (2015) 398-412.

355 [14] G. Chételat, B. Landeau, F. Eustache, F. Mézenge, F. Viader, V. de La Sayette, B. Desgranges, J.-C.
356 Baron, Using voxel-based morphometry to map the structural changes associated with rapid conversion
357 in MCI: a longitudinal MRI study, *Neuroimage* 27(4) (2005) 934-946.

358 [15] C. Platero, M.C. Tobar, A fast approach for hippocampal segmentation from T1-MRI for predicting
359 progression in Alzheimer's disease from elderly controls, *Journal of neuroscience methods* 270 (2016)
360 61-75.

361 [16] H. Li, L. Chen, Z. Huang, X. Luo, H. Li, J. Ren, Y. Xie, DeepOMe: a web server for the prediction of 2'
362 -O-Me sites based on the hybrid CNN and BLSTM architecture, *Frontiers in cell and developmental*
363 *biology* 9 (2021) 1244.

364 [17] F. Li, L. Tran, K.-H. Thung, S. Ji, D. Shen, J. Li, A robust deep model for improved classification of
365 AD/MCI patients, *IEEE journal of biomedical and health informatics* 19(5) (2015) 1610-1616.

366 [18] N. Zeng, Z. Wang, H. Zhang, W. Liu, F.E. Alsaadi, Deep belief networks for quantitative analysis of a
367 gold immunochromatographic strip, *Cognitive Computation* 8(4) (2016) 684-692.

368 [19] N. Zeng, H. Zhang, B. Song, W. Liu, Y. Li, A.M. Dobaie, Facial expression recognition via learning deep
369 sparse autoencoders, *Neurocomputing* 273 (2018) 643-649.

370 [20] M. Liu, J. Zhang, E. Adeli, D. Shen, Joint classification and regression via deep multi-task multi-
371 channel learning for Alzheimer's disease diagnosis, *IEEE Transactions on Biomedical Engineering* 66(5)
372 (2018) 1195-1206.

373 [21] E. Ocasio, T.Q. Duong, Deep learning prediction of mild cognitive impairment conversion to
374 Alzheimer's disease at 3 years after diagnosis using longitudinal and whole-brain 3D MRI, *PeerJ*
375 *Computer Science* 7 (2021) e560.

376 [22] J. Zhang, B. Zheng, A. Gao, X. Feng, D. Liang, X. Long, A 3D densely connected convolution neural
377 network with connection-wise attention mechanism for Alzheimer's disease classification, *Magnetic*
378 *Resonance Imaging* 78 (2021) 119-126.

379 [23] W. Lin, T. Tong, Q. Gao, D. Guo, X. Du, Y. Yang, G. Guo, M. Xiao, M. Du, X. Qu, Convolutional neural
380 networks-based MRI image analysis for the Alzheimer's disease prediction from mild cognitive
381 impairment, *Frontiers in neuroscience* 12 (2018) 777.

382 [24] N. Goenka, S. Tiwari, Deep learning for Alzheimer prediction using brain biomarkers, *Artificial*
383 *Intelligence Review* (2021) 1-45.

384 [25] X. Zhao, X.-M. Zhao, Deep learning of brain magnetic resonance images: A brief review, *Methods*
385 192 (2021) 131-140.

386 [26] A.M. Dale, B. Fischl, M.I. Sereno, Cortical surface-based analysis: I. Segmentation and surface

387 reconstruction, *Neuroimage* 9(2) (1999) 179-194.

388 [27] B. Fischl, A.M. Dale, Measuring the thickness of the human cerebral cortex from magnetic
389 resonance images, *Proceedings of the National Academy of Sciences* 97(20) (2000) 11050-11055.

390 [28] B. Fischl, A. Liu, A.M. Dale, Automated manifold surgery: constructing geometrically accurate and
391 topologically correct models of the human cerebral cortex, *IEEE transactions on medical imaging* 20(1)
392 (2001) 70-80.

393 [29] B. Fischl, D.H. Salat, E. Busa, M. Albert, M. Dieterich, C. Haselgrove, A. Van Der Kouwe, R. Killiany,
394 D. Kennedy, S. Klaveness, Whole brain segmentation: automated labeling of neuroanatomical structures
395 in the human brain, *Neuron* 33(3) (2002) 341-355.

396 [30] B. Fischl, M.I. Sereno, A.M. Dale, Cortical surface-based analysis: II: inflation, flattening, and a
397 surface-based coordinate system, *Neuroimage* 9(2) (1999) 195-207.

398 [31] B. Fischl, M.I. Sereno, R.B. Tootell, A.M. Dale, High - resolution intersubject averaging and a
399 coordinate system for the cortical surface, *Human brain mapping* 8(4) (1999) 272-284.

400 [32] B. Fischl, D.H. Salat, A.J. Van Der Kouwe, N. Makris, F. Ségonne, B.T. Quinn, A.M. Dale, Sequence-
401 independent segmentation of magnetic resonance images, *Neuroimage* 23 (2004) S69-S84.

402 [33] B. Fischl, A. Van Der Kouwe, C. Destrieux, E. Halgren, F. Ségonne, D.H. Salat, E. Busa, L.J. Seidman,
403 J. Goldstein, D. Kennedy, Automatically parcellating the human cerebral cortex, *Cerebral cortex* 14(1)
404 (2004) 11-22.

405 [34] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A.N. Gomez, Ł. Kaiser, I. Polosukhin,
406 Attention is all you need, *Advances in neural information processing systems*, 2017, pp. 5998-6008.

407 [35] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Deghani, M.
408 Minderer, G. Heigold, S. Gelly, An image is worth 16x16 words: Transformers for image recognition at
409 scale, *arXiv preprint arXiv:2010.11929* (2020).

410 [36] B.T. Wyman, D.J. Harvey, K. Crawford, M.A. Bernstein, O. Carmichael, P.E. Cole, P.K. Crane, C. DeCarli,
411 N.C. Fox, J.L. Gunter, Standardization of analysis sets for reporting results from ADNI MRI data,
412 *Alzheimer's & Dementia* 9(3) (2013) 332-337.

413 [37] V.J. Lowe, P.J. Peller, S.D. Weigand, C.M. Quintero, N. Tosakulwong, P. Vemuri, M.L. Senjem, L. Jordan,
414 C.R. Jack, D. Knopman, Application of the National Institute on Aging–Alzheimer’s Association AD
415 criteria to ADNI, *Neurology* 80(23) (2013) 2130-2137.

416 [38] B. Fischl, *FreeSurfer*, *Neuroimage* 62(2) (2012) 774-781.

417 [39] J.G. Sled, A.P. Zijdenbos, A.C. Evans, A nonparametric method for automatic correction of intensity
418 nonuniformity in MRI data, *IEEE transactions on medical imaging* 17(1) (1998) 87-97.

419 [40] W. Zheng, Z. Yao, Y. Xie, J. Fan, B. Hu, Identification of Alzheimer’s disease and mild cognitive
420 impairment using networks constructed based on multiple morphological brain features, *Biological*
421 *Psychiatry: Cognitive Neuroscience and Neuroimaging* 3(10) (2018) 887-897.

422 [41] A. Klein, J. Tourville, 101 labeled brain images and a consistent human cortical labeling protocol,
423 *Frontiers in neuroscience* 6 (2012) 171.

424 [42] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, *Proceedings of the IEEE*
425 *conference on computer vision and pattern recognition*, 2016, pp. 770-778.

426 [43] J.L. Ba, J.R. Kiros, G.E. Hinton, Layer normalization, *arXiv preprint arXiv:1607.06450* (2016).

427 [44] D.P. Kingma, J. Ba, Adam: A method for stochastic optimization, *arXiv preprint arXiv:1412.6980*
428 (2014).

429 [45] V. Vapnik, *The nature of statistical learning theory*, Springer science & business media 1999.

430 [46] Y.-Y. Song, L. Ying, *Decision tree methods: applications for classification and prediction*, Shanghai

431 archives of psychiatry 27(2) (2015) 130.

432 [47] A. Liaw, M. Wiener, Classification and regression by randomForest, R news 2(3) (2002) 18-22.

433 [48] D.R. Cox, The regression analysis of binary sequences, Journal of the Royal Statistical Society: Series

434 B (Methodological) 20(2) (1958) 215-232.

435 [49] W. Zaremba, I. Sutskever, O. Vinyals, Recurrent neural network regularization, arXiv preprint

436 arXiv:1409.2329 (2014).

437 [50] S. Hochreiter, J. Schmidhuber, Long short-term memory, Neural computation 9(8) (1997) 1735-1780.

438 [51] J. Chung, C. Gulcehre, K. Cho, Y. Bengio, Empirical evaluation of gated recurrent neural networks on

439 sequence modeling, arXiv preprint arXiv:1412.3555 (2014).

440 [52] J. Sadr, I. Jarudi, P. Sinha, The role of eyebrows in face recognition, Perception 32(3) (2003) 285-293.

441 [53] H.-I. Suk, S.-W. Lee, D. Shen, A.s.D.N. Initiative, Hierarchical feature representation and multimodal

442 fusion with deep learning for AD/MCI diagnosis, NeuroImage 101 (2014) 569-582.

443 [54] K. Oh, Y.-C. Chung, K.W. Kim, W.-S. Kim, I.-S. Oh, Classification and visualization of Alzheimer's

444 disease using volumetric convolutional neural network and transfer learning, Scientific Reports 9(1)

445 (2019) 1-16.

446 [55] K. Kwak, M. Niethammer, K.S. Giovanello, M. Styner, E. Dayan, A.s.D.N. Initiative, Differential Role

447 for Hippocampal Subfields in Alzheimer's Disease Progression Revealed with Deep Learning, Cerebral

448 Cortex (2021).

449 [56] R.J. Killiany, T. Gomez-Isla, M. Moss, R. Kikinis, T. Sandor, F. Jolesz, R. Tanzi, K. Jones, B.T. Hyman,

450 M.S. Albert, Use of structural magnetic resonance imaging to predict who will get Alzheimer's disease,

451 Annals of Neurology: Official Journal of the American Neurological Association and the Child Neurology

452 Society 47(4) (2000) 430-439.

453 [57] J.R. Gray, T.S. Braver, Personality predicts working-memory—related activation in the caudal

454 anterior cingulate cortex, Cognitive, Affective, & Behavioral Neuroscience 2(1) (2002) 64-75.

455 [58] C.Y. Wee, P.T. Yap, D. Shen, A.s.D.N. Initiative, Prediction of Alzheimer's disease and mild cognitive

456 impairment using cortical morphological patterns, Human brain mapping 34(12) (2013) 3411-3425.

457 [59] H. Tabatabaei-Jafari, M.E. Shaw, E. Walsh, N. Cherbuin, A.s.D.N. Initiative, Regional brain atrophy

458 predicts time to conversion to Alzheimer's disease, dependent on baseline volume, Neurobiology of

459 aging 83 (2019) 86-94.

460 [60] V. Julkunen, E. Niskanen, S. Muehlboeck, M. Pihlajamäki, M. Könönen, M. Hallikainen, M. Kivipelto,

461 S. Tervo, R. Vanninen, A. Evans, Cortical thickness analysis to detect progressive mild cognitive

462 impairment: a reference to Alzheimer's disease, Dementia and geriatric cognitive disorders 28(5) (2009)

463 389-397.

464 [61] T. Zhang, M. Shi, Multi-modal neuroimaging feature fusion for diagnosis of Alzheimer's disease,

465 Journal of Neuroscience Methods 341 (2020) 108795.

466 [62] L. Xu, X. Wu, R. Li, K. Chen, Z. Long, J. Zhang, X. Guo, L. Yao, A.s.D.N. Initiative, Prediction of

467 progressive mild cognitive impairment by multi-modal neuroimaging biomarkers, Journal of Alzheimer's

468 Disease 51(4) (2016) 1045-1056.

469 [63] F. Liu, C.-Y. Wee, H. Chen, D. Shen, Inter-modality relationship constrained multi-modality multi-task

470 feature selection for Alzheimer's Disease and mild cognitive impairment identification, NeuroImage 84

471 (2014) 466-475.

472 **Supplementary Materials**

473 **TABLE S1** The hyperparameters of RNN, LSTM, and GRU.

Hyperparameter	RNN	LSTM	GRU
Epoch	100	50	50
Batch size	64	32	64
Optimizer	Adam	Adam	Adam
Learning rate	1e-3	1e-3	1e-3
Weight decay	1e-8	1e-8	1e-8

474 Abbreviations: RNN = Recurrent Neural Network, LSTM = Long Short-Term Memory, GRU = Gated
 475 Recurrent Unit.

476 **TABLE S2** Results of the 10-fold cross-validation. All metrics are reported as mean \pm SD across folds.

Model	ACC	SEN	SPE	AUC
Proposed model	0.719 \pm 0.084	0.797 \pm 0.109	0.545 \pm 0.214	0.668 \pm 0.092
RNN	0.603 \pm 0.108	0.866 \pm 0.095	0.272 \pm 0.182	0.600 \pm 0.105
LSTM	0.615 \pm 0.105	0.928 \pm 0.101	0.203 \pm 0.139	0.624 \pm 0.072
GRU	0.611 \pm 0.106	0.898 \pm 0.110	0.264 \pm 0.207	0.603 \pm 0.100
Random Forest	0.562 \pm 0.063	0.638 \pm 0.150	0.478 \pm 0.128	0.621 \pm 0.108
Decision Tree	0.618 \pm 0.119	0.693 \pm 0.143	0.533 \pm 0.157	0.613 \pm 0.122
Logistic Regression	0.619 \pm 0.062	0.682 \pm 0.126	0.529 \pm 0.164	0.619 \pm 0.071
Support vector machine	0.518 \pm 0.117	0.583 \pm 0.159	0.449 \pm 0.227	0.516 \pm 0.109

477 The best results for each column are shown in boldface. Abbreviations: SD = standard deviation,
 478 RNN = Recurrent Neural Network, LSTM = Long Short-Term Memory, GRU = Gated Recurrent Unit,
 479 ACC = accuracy, SEN = sensitivity, SPE = specificity, AUC = area under the receiver operating
 480 characteristic (ROC) curve.