

IT'S YOUR TIME BE PRECISE

STANDARDIMAGING



Every day we spend **our time**
optimizing ways to make
QA easy and reliable.

Ask us how our solutions
can benefit you.

WWW.STANDARDIMAGING.COM

Ensemble of Convolutional Neural Networks and Multilayer Perceptron for the Diagnosis of Mild Cognitive Impairment and Alzheimer's Disease

Minglei Li¹, Yuchen Jiang¹, Xiang Li¹, Shen Yin^{2,*}, Hao Luo^{1,*}

¹Department of Control Science and Engineering, Harbin Institute of Technology, Harbin 150001, Heilongjiang, China

²Department of Mechanical and Industrial Engineering, Norwegian University of Science and Technology, Trondheim 7033, Norway

Shen Yin. Email: shen.yin@ntnu.no

Hao Luo. Email: hao.luo@hit.edu.cn

Abstract

Background: Structural magnetic resonance imaging (sMRI) can provide morphological information about the structure and function of the brain in the same scanning process. It has been widely used in the diagnosis of Alzheimer's disease (AD) and mild cognitive impairment (MCI).

Purpose: To capture the anatomical changes in the brain caused by AD/MCI, deep learning-based MRI image analysis methods have been proposed in recent years. However, it is observed that the performance of most existing methods is limited as they only construct a single type of deep network and ignore the significance of other clinical information.

Methods: To make up for these defects, an ensemble framework that incorporates three types of dedicatedly-designed convolutional neural networks (CNNs) and a multilayer perceptron (MLP) network is proposed, where three CNNs with entropy-based multi-instance learning pooling layers have more reliable feature selection abilities. The dedicatedly-designed base classifiers can make use of the heterogeneous data, and empower the framework with enhanced diversity and robustness. In particular, to consider the interactions among the base classifiers, a novel multi-head self-attention voting scheme is designed. Moreover, considering the chance that MCI can be transformed to AD, the proposed framework is designed to diagnose AD and predict MCI conversion simultaneously, with the aid of the transfer learning technique.

Results: For performance evaluation and comparison, extensive experiments are conducted on the public dataset of the Alzheimer's Disease Neuroimaging Initiative (ADNI). The results show that the proposed ensemble framework provides superior performance under most of the evaluation metrics. Especially, the proposed framework achieves state-of-the-art diagnostic accuracy (98.61% for the AD diagnosis task, and 84.49% for the MCI conversion prediction task).

Conclusions: These promising results demonstrate the proposed ensemble framework can accurately diagnose AD patients and predict the conversion of MCI patients, which has the potential of clinical practice for diagnosing AD and MCI.

Keywords: Alzheimer’s disease, magnetic resonance imaging, computer-aided diagnosis, ensemble learning, multiple instance learning

Contents

I. INTRODUCTION	1
II. MATERIAL AND METHODS	3
II.A. Data acquisition and image pre-processing	3
II.B. Overall ensemble learning architecture	5
II.C. Base classifiers in the ensemble framework	6
II.C.1. Entropy-based MIL pooling	6
II.C.2. Base classifier 1	8
II.C.3. Base classifier 2	9
II.C.4. Base classifier 3	10
II.C.5. Base classifier 4	10
II.D. Ensemble approaches for classifiers	11
II.E. Implementations	13
II.E.1. Data split	13
II.E.2. Training strategy	14
II.E.3. Evaluation metrics	14
III. RESULTS	15
III.A. Comparison with other methods	15
III.B. Effectiveness of entropy-based MIL pooling	16
III.C. Effectiveness of MHSA voting	16
III.D. Influence of transfer learning	18
III.E. Influence of clinical information	18
III.F. Indispensability of four base classifiers	20
IV. DISCUSSIONS	22
V. CONCLUSION	26

67	VI. ACKNOWLEDGMENTS	26
68	VII.CONFLICT OF INTEREST DISCLOSURE	27
69	References	27

Accepted Article

1. INTRODUCTION

Alzheimer's disease (AD) is a chronic neurodegenerative disease and contributes to 60-80% of dementias, over 30 million people around the world are diagnosed with AD^{1,2}. As the most common form of dementia, AD can cause irreversible damage or destruction of neurons in brain regions over time, and gradually has a serious impact on the life of patients. Mild cognitive impairment (MCI) is often seen as a preclinical stage of AD, the predominant symptom of MCI is mild memory loss which has less impact on a person than AD³. Around 10% of the MCI patients worldwide develop to AD per year, while a majority of them stay stable or even revert to the normal state⁴. Those MCI patients who develop to AD are medically known as progressive MCI (pMCI), in contrast, patients who stay stable are stable MCI (sMCI). Therefore, distinguishing sMCI from pMCI has been typically considered as an early prediction of AD dementia. In particular, because there is no effective treatment to cure AD, reliable early diagnosis is crucial for the control of AD. And early diagnosis will help for the better targeted selection of individuals with MCI, thus allowing early implementation of treatment strategies and altering the course of this disease⁵.

Various biomarkers (e.g., positron emission tomography (PET)⁶ and MRI⁷) and biospecimens (e.g., cerebrospinal fluid (CSF)⁸) measured in vivo constitute dominant features in the diagnosis of AD. These biomarkers and biospecimens are typically employed for evaluating the development of AD, which have been well validated in many clinical settings⁹. For example, structural MRI can noninvasively capture cerebral atrophy caused by loss of neurons and dendritic pruning¹⁰, which provides a powerful auxiliary pattern for brain research and clinical diagnosis. In addition, the clinical information of individuals can be used to partially indicate disease status, which typically includes demographic information and cognitive and neuropsychological measures. Many cognitive and neuropsychological measures, such as the Mini Mental State Examination (MMSE)¹¹, Clinical Dementia Rating Scale (CDRSB)¹², Alzheimer's Disease Assessment Scale (ADAS)¹³, and Ray Auditory Verbal Learning Test (RAVLT)¹⁴, etc., can reflect the cognitive level of an individual and reveal the disease progression.

Computer-aided methods have been a growing interest in the assessment and treatment of serious brain diseases, such as brain tumors¹⁵, autism¹⁶, and Parkinson's disease¹⁷. AD as one of the serious brain diseases also receives much attention. To achieve the reliable diagnosis of AD and MCI, machine learning- (ML) or deep learning- (DL) based methods have been developed in many studies based on sMRI. These existing methods include at least two main components: 1) extraction of imaging features and 2) construction of classification

104 models. According to the scale of feature extraction, these methods are usually categorized
105 into 1) subject-level, 2) region-level, 3) patch-level and 4) slice-level¹⁸. The subject-level
106 methods^{19–22} extract features from voxel intensities directly, while the extracted features
107 are high dimensional and these methods are susceptible to overfitting due to the small
108 number of samples. The region-level methods^{23–26} focus on pre-determined brain regions
109 of structure or function, and extract representative features from these regions. Although
110 region-level features have lower dimensions than subject-level features, they may not cover
111 all possible pathological parts of the whole brain and miss some subtle changes in pathology.
112 The patch-level methods^{27–29} combine the above two methods, attempting to capture the
113 disease-related pathologies in the local brain. The key step of patch-level methods is to
114 select patches and combine them to obtain information about the brain. The slice-level
115 methods^{30,31} are closer to the diagnosis modes of physicians, which utilize 2D slice images
116 from sMRI to extract features and then count each slice-level result to obtain a subject-level
117 diagnosis. ML-based methods usually need to extract features manually and then construct
118 a conventional classifier to complete diagnosis, such as support vector machine (SVM). While
119 DL-based methods perform feature extraction and classification only by convolutional neural
120 networks (CNNs), which have been demonstrated more powerful than ML-based methods.

121 In the above methods, the requirements of slice-level methods for computing resources
122 are much lower than the use of regions, patches, or subjects. And the architectures of clas-
123 sifiers in slice-level methods are also simpler than other methods. In addition, the superior
124 performance of DL-based methods often depends on numerous learnable parameters of net-
125 works. Many existing DL-based studies have been limited to using a single CNN for AD
126 diagnosis or MCI conversion prediction. However, due to the scarcity of medical data, it is
127 challenging for an individual CNN to achieve reliable classification with the small number
128 of available training data.

129 To overcome this limitation, ensemble learning methods have been applied to the disease
130 diagnosis, and effectively combined with the CNN³². There are very few works used CNN-
131 based ensemble classifiers for AD diagnosis in recent years^{33–36}. Ensemble learning is the
132 algorithm that constructs a set of classifiers and then performs classification by aggregating
133 their predictions³⁷. And the ensemble learning methods have been proved that can enhance
134 the reliability of diagnosis, while the main drawback of these works is that each classifier is
135 assigned the same weight when the final results are obtained by the majority- and average-
136 voting. These fusion methods do not perform adaptive fusion based on each classifier and
137 may be affected by the weaker classifier in the ensemble.

I. INTRODUCTION

138 In this work, the target is to propose an ensemble framework that can conduct the
139 reliable diagnosis of AD and MCI simultaneously. For clarity, the following two research
140 tasks are defined:

141 1) *Task 1 (AD vs. CN)*: Distinguish between whether a subject (a patient) is cognitively
142 normal (CN) or with AD.

143 2) *Task 2 (pMCI vs. sMCI)*: Distinguish between whether an MCI patient belongs to
144 pMCI or sMCI.

145 The contributions of this work can be summarized as follows:

- 146 • A robust ensemble learning framework is proposed to make use of the multi-modal
147 information/heterogeneous data. Three types of dedicatedly-designed CNNs are in-
148 corporated to exploit information from sMRI, and a shallow network (i.e., MLP) is
149 employed to exploit the clinical information.
- 150 • A multi-head self-attention voting scheme is proposed as an ensemble approach for
151 base classifiers. The interactions among the classifiers are considered, and the defect
152 that common voting approaches ignore the relationships among classifiers is overcome.
- 153 • Multi-instance learning (MIL) is incorporated into base CNN classifiers. The entropy-
154 based MIL pooling layer can reasonably consider the expressive abilities of different
155 slices and integrate slice-level features.

156 II. MATERIAL AND METHODS

157 II.A. Data acquisition and image pre-processing

158 We consider a dataset obtained from ADNI-1 and ADNI-2 in the Alzheimer’s Disease Neu-
159 roimaging Initiative (<http://www.loni.ucla.edu/ADNI>)³⁸. The ADNI database is the largest
160 publicly available Alzheimer’s disease dataset and has been used in quite a few studies.
161 Specifically, the baseline dataset contains T1-weighted MRI obtained from 771 subjects,
162 which consists of 244 CN, 299 MCI, and 228 AD subjects. Depending on whether the MCI
163 subjects progressed to the AD stage within 36 months after baseline assessment, they can be
164 further divided into 170 sMCI and 129 pMCI subjects. The demographic information (age,
165 gender, and education years), cognitive and neuropsychological measures (CDRSB, ADAS,
166 MMSE, RAVLT) as well as the ApoE4 genotyping of the subjects are shown in Table 1.

Table 1: Information summary of the studied dataset extracted from ANDI

Gender (M/F)	Age	Education (years)	APoE4 level			CDRSB	ADAS			MMSE	RAVLT				
			0	1	2		ADAS11	ADAS13	ADASQ4		immediate	learning	forgetting	%forgetting	
CN	118/126	74.2±6.0	16.5±2.6	17861	5	0.2±0.1	5.6±2.7	8.6±4.0	2.7±1.7	29.1±1.1	45.4±10.0	5.9±2.2	3.7±2.6	34.9±26.7	
sMCI	99/71	71.8±7.4	16.2±2.9	10453	13	1.2±0.7	8.7±3.8	13.9±5.6	4.7±2.2	28.1±1.7	37.9±11.1	4.9±2.6	4.3±2.4	50.9±30.7	
pMCI	73/56	73.8±7.1	15.9±2.8	42	57	30	2.0±1.0	13.0±4.0	21.4±5.2	7.4±1.9	26.6±1.7	28.0±6.9	3.1±2.0	5.2±2.3	77.4±27.8
AD	124/104	74.9±7.8	15.2±2.9	71	11542	4.5±1.6	19.9±6.6	30.1±7.8	8.6±1.5	23.1±2.0	22.9±7.1	2.0±1.6	4.5±1.7	88.8±21.4	

*The data are presented as mean ± standard deviation (std).

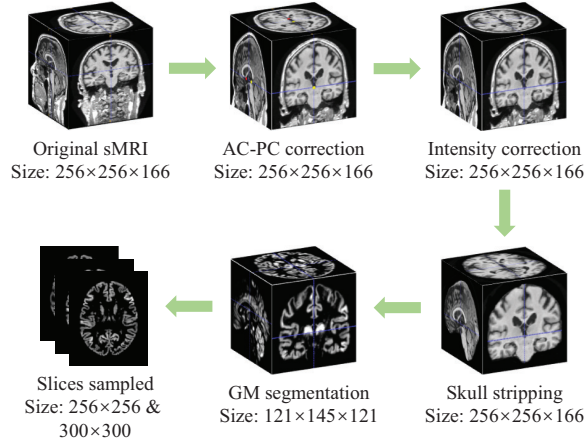


Figure 1: The preprocessing pipeline of sMRI. The pipeline includes AC-PC correction, intensity correction, skull stripping, tissue segmentation, and slice selection. Taking an sMRI with the size of $256 \times 256 \times 166$ voxels as an example, the image size after each processing step is shown.

As shown in Fig. 1, the sMRI data go through a standard pipeline preprocessing procedure, including anterior commissure (AC)-posterior commissure (PC) correction, intensity correction, skull stripping, tissue segmentation, and slice selection. Specifically, we use the MIPAV software (<https://mipav.cit.nih.gov/clickwrap.php>) for AC-PC correction and adopt N3 algorithm³⁹ for intensity correction. Skull stripping and tissue segmentation are performed by using the CAT12 toolbox (<http://dbm.neuro.uni-jena.de/cat/>) via SPM12 software (<http://www.fil.ion.ucl.ac.uk/spm/software/spm>). Following skull stripping, the quality of the preprocessed images is checked manually. And the qualified images are then segmented to obtain the gray matter (GM) tissues, which are aligned to Montreal Neurological Institute T1 Template⁴⁰. The GM images are smoothed with a 3.0 mm full width at half maximum (FWHM) isotropic Gaussian kernel. As a result, the sizes of obtained GM tissues are $121 \times 145 \times 121$ voxels, and the spatial resolutions are $1.5 \times 1.5 \times 1.5 \text{ mm}^3$. Considering that GM is the most notably affected tissue by AD, it is used for feature extraction. Then, the 3D volumetric data is sectioned along the axial direction, and the slices are sampled from the central slice to the edges of the 3D volumetric data. The edge slices largely cover cross-sections of the brain stem, cerebellum, and cerebral cortex, which are the anatomic

183 areas less relevant to AD pathology. Therefore, the middle two-thirds of the slices (80 slices)
 184 are selected and resized to 256×256 and 300×300 pixels. The selected slices cover areas
 185 including ventricle, inferior temporal, and middle temporal cortices. And these areas have
 186 been reported as the regions correlated with AD pathology, which can provide rich tissue
 187 information⁴¹.

188 For the clinical information, numerical normalization (i.e., Min-Max normalization) is
 189 employed to normalize the values of each separate clinical factor to the range of $[0, 1]$.

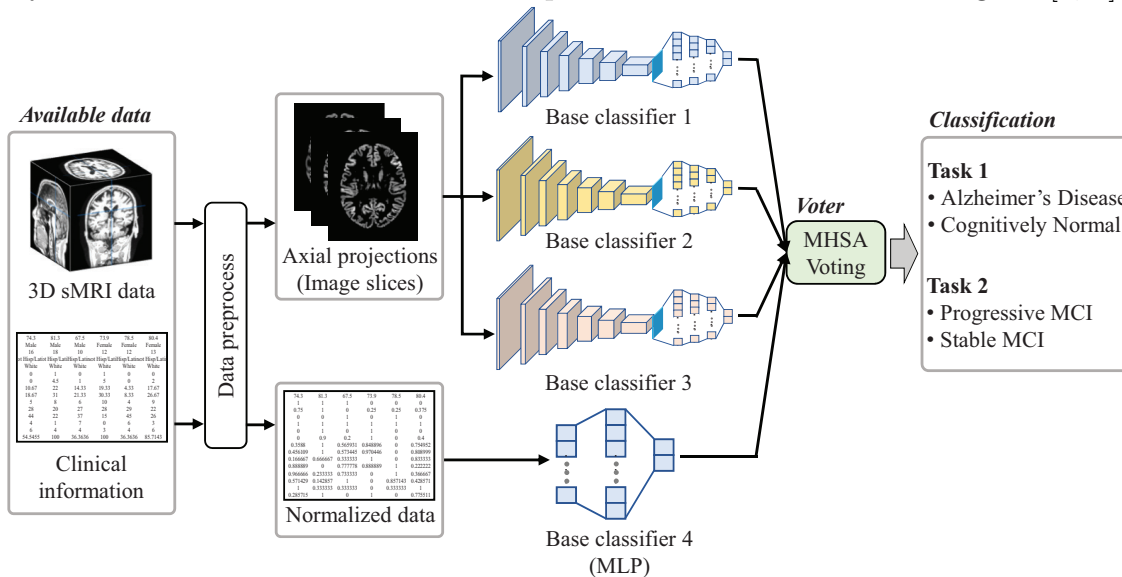


Figure 2: Illustration of the proposed ensemble framework for AD diagnosis and MCI conversion prediction. Raw 3D sMRI and corresponding clinical information of each individual are first preprocessed, multiple 2D slices are sampled from each sMRI, and the clinical information is normalized. The processed data are then fed into four different base classifiers, and an MHS voting scheme aggregates the outputs of each base classifier for the final prediction.

190 II.B. Overall ensemble learning architecture

191 The proposed ensemble framework is illustrated in Fig. 2, where the inputs are the 3D sMRI
 192 data and clinical information, and the output is the AD diagnosis (i.e., AD or CN) or MCI
 193 conversion prediction (i.e., pMCI or sMCI). Specifically, 3D sMRI and clinical information of
 194 each individual are processed via several preprocessing steps. After that, the multiple slices
 195 sampled from 3D sMRI and normalized clinical information are as the inputs of different base
 196 classifiers. The base classifiers are designed to have different architectures, each base classifier

197 can play an important part in this ensemble framework. Base classifier 1, base classifier 2,
198 and base classifier 3 are used to extract the features of images and give the initial predictions
199 based on sMRI data, where the entropy-based multi-instance learning (MIL) pooling layer
200 is designed to consider different information densities of slices and further improve their
201 expression abilities. Base classifier 4 is designed as an MLP to make use of the clinical
202 information, which can introduce different patient information than the sMRI modal. Then,
203 four base classifiers are fused via MHSA voting to obtain the classification results for two
204 classification tasks (i.e., AD *vs.* CN and pMCI *vs.* sMCI).

205 II.C. Base classifiers in the ensemble framework

206 In this section, the detailed architecture of each base classifier and their mentalities of de-
207 signing are introduced, including three CNNs (base classifier 1, 2, and 3) and an MLP model
208 (base classifier 4).

209 The architectures of base classifiers are shown in Fig. 3, all base CNN classifiers (Base
210 classifier1, 2, and 3) have feature extraction, entropy-based MIL pooling, and classification
211 layer three parts. Scaling up the dimension of network width, depth, and resolution has been
212 widely used to improve the performance of networks. However, scaling up a CNN in all three
213 dimensions of width, depth, and resolution will greatly increase the number of parameters.
214 Considering the consumption of computing resources and the efficiency of ensemble learning,
215 it is not necessary to design an overly complex model as one of the base classifiers. Thus, the
216 three base CNN classifiers are scaled up in width, depth and resolution, respectively. The
217 number of layers and the number of parameters in these classifiers are controlled. As a result,
218 the average number of three CNNs parameters is less than that of ResNet34⁴², and the layers
219 of them are less than 19 layers. Specifically, three base CNN classifiers have different scales
220 of network width, depth and resolution, respectively. Base classifier 1 has higher resolutions
221 than the other two, which means that it can potentially capture more fine-grained patterns.
222 Base classifier 2 only scales up in terms of network depth. Deeper networks can fit more
223 complex deep features. Base classifier 3 has a wider architecture and can focus on richer
224 features. More details of these base classifiers will be introduced as follows.

225 II.C.1. Entropy-based MIL pooling

226 AD-related pathological areas usually exist in some partial areas of the brain, and these areas
227 in sMRI images are unlabeled, namely, only the entire sMRI image is labeled as a certain

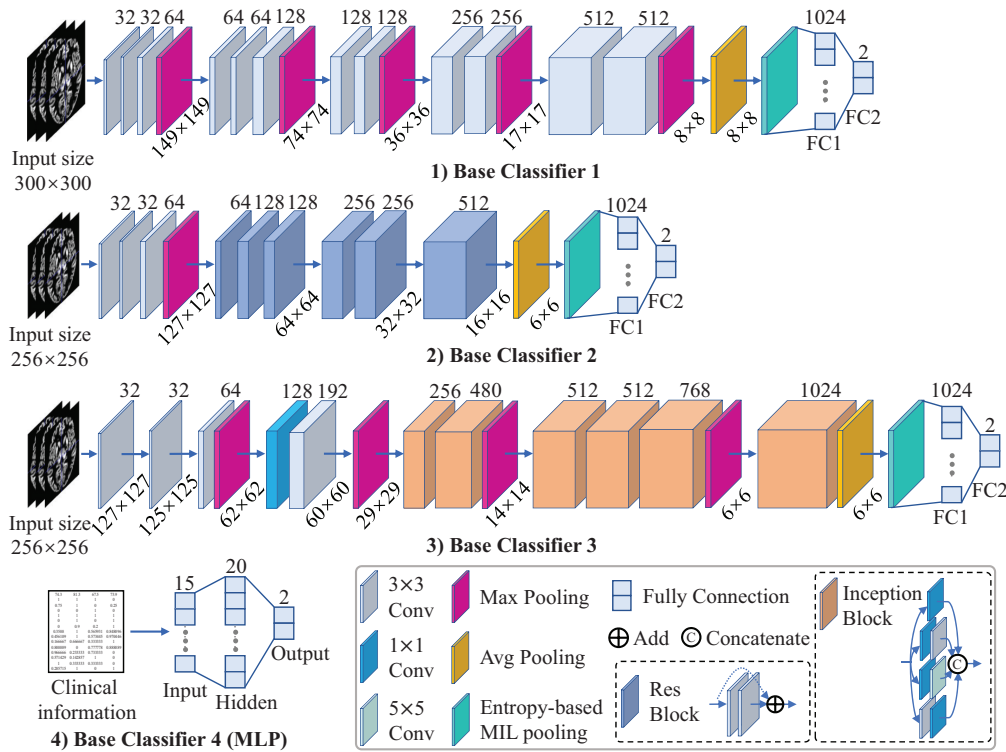


Figure 3: The architectures of base classifiers. Base classifier 1, base classifier 2, and base classifier 3 are CNN based classifiers with MRI images as inputs, which mainly consist of convolutional layers, designed special blocks (i.e., Res block, Inception block), and pooling layers. Base classifier 4 is an MLP with clinical information as inputs, and it consists of fully connected layers. The number of channels for each convolutional layer or special block is displayed above them. When the sizes of the feature maps change after passing through some layers, the sizes are shown below the convolutional layers or special blocks in the form of $H \times W$.

category. As described in Section II.A., the slices are sampled from 3D volumetric data along the axial direction and used as inputs of base CNN classifiers. These processes can be seen as the construction of bags in MIL. Considering the properties and preprocessing processes of sMRI images, both tasks in this work can be solved with the MIL strategy.

Let $X_i = \{x_{i1}, x_{i2}, \dots, x_{in_i}\}$ denotes the bag of the i -th sMRI, where $x_{kl} \in \mathbb{R}^d$ ($k = 1, 2, \dots, n_k$) represents the l -th slice of the k -th bag. Then, these slices are input into the feature extraction part of base CNN classifiers to obtain slice-level features $E_i = \{e_{i1}, e_{i2}, \dots, e_{in_i}\}$, followed by a proposed entropy-based MIL pooling layer to generate embedding-level features \mathcal{B}_i from slice-level features. The proposed entropy-based MIL pooling layer combines information entropy with MIL. The information entropy of an image is a statistical form of the features, which evaluates the information density of an image. In

239 general, the images with high entropy values have more information about target areas (e.g.,
 240 brain, lung, etc.). In the clinical environment, for medical images with explicit sequences,
 241 such as MRI and CT, physicians also focus on the slices with more abundant tissue infor-
 242 mation when diagnosing diseases. Entropy as a form of reflecting image information density,
 243 combining it with MIL can not only be closer to actual clinical diagnosis, but also further
 244 improve the performance of diagnosis. This is the motivation for us to design entropy-
 245 based MIL pooling. The entropy-based MIL pooling layer can be described by the following
 246 equations:

$$247 \quad \mathcal{B}_i = \text{Concat}_{l=1}^{n_i}(h_{il} \cdot e_{il}) \quad (1)$$

$$248 \quad h_{il} = \text{norm}\left(\frac{H_{il}}{\sum_{l=1}^{n_i} H_{il}}\right) \quad (2)$$

250 where h_{il} is normalized weight that can be calculated by Eq. (2), and H_{il} in Eq. (2)
 251 is the information entropy of the l -th slice of i -th sMRI. e_{il} corresponds to the l -th slice-
 252 level features of E_i . *Concat* is channel concatenation. In addition, mean MIL pooling and
 253 maximum MIL pooling are commonly used operators in MIL. Mean MIL pooling considers
 254 that all slices have the same ability to express the information of features, it generates
 255 embedding-level features by averaging slice-level features. Maximum MIL pooling depends
 256 on only one slice to determine the prediction of the individual. Different from these two
 257 pooling operators, entropy-based MIL pooling comprehensively considers the information
 258 entropy of different slices, which can utilize the information expression ability of these slices
 259 to achieve a more accurate diagnosis.

260 After obtaining embedding-level features \mathcal{B}_i , the classification layer is used to predict
 261 the category (i.e., AD, CN, pMCI, or sMCI) of each input sMRI.

$$262 \quad P(Y|X) = f_{cls}(\mathcal{B}_i) \quad (3)$$

263 where $P(Y|X)$ is the probability that the subject belongs to a specific class, Y denotes the
 264 true category, $f_{cls}(\cdot)$ denotes the mapping function of the classification layer.

265 II.C.2. Base classifier 1

266 The base classifier 1 is designed to have higher resolution, and it is constructed by stacking
 267 convolutional layers without adopting more complex modules. Specifically, base classifier 1
 268 contains twelve convolutional (Conv) layers, an entropy-based MIL pooling layer, and two
 269 fully connected (FC) layers. The number of channels for Conv layers is mainly 32, 64, 128,

270 256, and 512. Each Conv layer consists of one convolutional layer, batch normalization (BN),
271 and rectified linear unit (ReLU) activation, where the convolutional layer has 3×3 kernel size,
272 unit stride with unit zero padding. Several 3×3 max pooling layers and an adaptive average
273 pooling layer are inserted in the specific positions of the model, which can downsample the
274 number or depth of the intermediate feature maps. An entropy-based MIL pooling layer is
275 inserted between the average pooling layer and FC layers. At the end, two FC layers with
276 1024 and 2 nodes respectively as classification layer are adopted to map distributed features
277 into the sample label space. The input images of base classifier 1 have higher resolutions
278 than those of base classifier 2 and base classifier 3, and the intermediate feature maps also
279 have higher resolutions. With high resolutions, base classifier 1 tends to be more sensitive
280 to fine-grained patterns, which can better focus on subtle pathological changes in slices.

281 II.C.3. Base classifier 2

282 The base classifier 2 with deeper depth is designed to characterize complex nonlinearities.
283 Scaling up the depth of networks may bring gradient instability and network degradation,
284 therefore, base classifier 2 draws on the idea of residual learning, which adopts Conv layers
285 and residual (Res) blocks as main components. Specifically, it consists of three Conv layers,
286 six Res blocks, an entropy-based MIL pooling layer, and two FC layers. At the beginning of
287 the model, three Conv layers with the same composition as in base classifier 1 are used to
288 extract shallow feature maps, where the number of channels for Conv layers is 32, 32 and 64,
289 respectively. Then a max pooling layer merges the features and reduces their dimensions,
290 followed by six Res blocks. As shown in Fig. 3, each Res block contains two serial Conv
291 layers, and the output of the second Conv layer adds the input of the Res block through a
292 shortcut connection, the result of the addition is used as the output of the Res block. The
293 number of channels for Res blocks is 64, 128, 128, 256, 256 and 512, respectively. In order to
294 achieve the effect of downsampling, the stride of the first Conv layer in the third, fifth and
295 sixth Res block is respectively set to 2, other Conv layers in Res blocks have the same settings
296 as the Conv layers in the base classifier 1. After that, the average pooling layer, MIL pooling
297 layer, and classification layer that same as base classifier 1 are adopted. Base classifier 2
298 with deeper depth is designed to characterize complex nonlinearities. The Res blocks can
299 transfer shallow feature information extracted by three Conv layers to deeper layers, thereby
300 enhancing feature representations and strengthening their learning. Benefiting from network
301 depth, base classifier 2 has better nonlinear representation ability, which can learn to fit more
302 complex features and generalize well on diagnostic tasks.

303 II.C.4. Base classifier 3

304 The base classifier 3 is designed as a network with wider architecture. The suitable network
305 width can ensure that the layers learn rich features, such as texture features in different
306 frequencies and different directions. Base classifier 3 consists of five Conv layers, six Inception
307 blocks, an entropy-based MIL pooling layer, and two FC layers. To maintain the proper size
308 of feature maps, the stride of the first Conv layer is set to 2, followed by four Conv layers,
309 where 1×1 Conv layer allows the model to control the depth of the feature more flexibly as
310 needed. The number of channels for Conv layers is 32, 32, 64, 128 and 192, respectively.
311 After serial Conv layers, the Inception blocks further process the extracted features. As
312 shown in Fig. 3, each Inception block has four paths to perform convolution operations on
313 the input and concatenates to generate the output of the block, it contains several 1×1 , 3×3 ,
314 and 5×5 Conv layers. The number of input channels for Inception blocks is 256, 480, 512,
315 512, 768 and 1024, respectively. Similar to the base classifier 1, the 3×3 max pooling layers
316 and an adaptive average pooling layer are inserted in the specific positions to downsample
317 the feature maps, the MIL pooling layer and classification layer are inserted at the end of
318 the model. In base classifier 3, the maximum number of channels for blocks reaches 1024,
319 which is twice the maximum number of the other base CNN classifiers. More channels
320 characterize richer feature information of images, which can endow the model with better
321 representational ability. Thus, base classifier 3 with wide architecture can potentially better
322 learn and characterize rich tissue information in slices.

323 II.C.5. Base classifier 4

324 As summarized in Table 1, the clinical information data including age, gender, cognitive test,
325 etc. were collected from the subjects. Since these data are not as complicated as images,
326 shallow neural networks are enough to mine information in these clinical data. For this
327 reason, MLP is chosen as the base classifier 4. In more detail, the MLP is composed of three
328 layers, including an input layer, a hidden layer, and an output layer. The number of nodes
329 for three layers is 15, 20 and 2, respectively. All layers contain one FC layer, followed by BN
330 and ReLU activations. Since the MLP is simple in structure and with few parameters, it is
331 suitable for clinical information data analysis.

332 The loss function in the proposed base classifiers for classification can be formulated as:

$$333 \mathcal{L}(X, Y, P, \omega_c) = -\log(P(Y|X), Y) \quad (4)$$

334 where X denotes the input data of the base classifiers (i.e., sMRI for base CNN classifiers,
 335 clinical information data for MLP), Y denotes the corresponding true label, P denotes the
 336 predicted results, and ω_c is the learnable parameters of these classifiers.

337 II.D. Ensemble approaches for classifiers

338 The predictions from these trained base classifiers are combined by different ensemble ap-
 339 proaches. Specifically, common voting approaches (i.e., majority voting, weighted voting,
 340 SVM-based voting) and proposed multi-head self-attention (MHSA) voting have been per-
 341 formed on classifiers of ensemble framework and compared. In common voting approaches,
 342 the fixed weight is assigned to each classifier in the ensemble for the aggregation of the clas-
 343 sification results. The major drawback of these approaches is that the aggregation is not
 344 data-adaptive and ignores the interactions among base classifiers, which potentially brings
 345 bias to the final classification, especially in the presence of weak base classifiers.

346 Considering that common voting approaches ignore the interactions among base clas-
 347 sifiers and potentially introduce bias resulting in unreliable predictions, an MHSA voting
 348 scheme is proposed to aggregate the results of base classifiers, which can calculate and ex-
 349 ploit the interactions among base classifiers during their fusion. The MHSA voting is to
 350 calculate the correlation and importance among the base classifiers, and then use these in-
 351 teractions to aggregate the results and obtain the final classification results. It is defined as
 352 linear transformation, interaction calculation, and aggregation & final decision three stages.
 353 The proposed MHSA voting scheme is shown in Fig. 4

354 1) *Linear Transformation*: In this stage, the outputs of each base classifier are linearly
 355 transformed into three vectors q , k , and v , and the distribution spaces of these vectors
 356 are basically the same. Formally, an embedded representation is constructed to represent
 357 the outputs of all base classifiers. Denote the embedded representation as Φ , where $\Phi =$
 358 $[\phi_1, \phi_2, \dots, \phi_n, \dots, \phi_N]^T \in \mathbb{R}^{N \times C}$. Here, $\phi_n \in \mathbb{R}^{1 \times C}$ ($n = 1, 2, \dots, N$) indicates the outputs of
 359 the n -th base classifier, N and C are the number of base classifiers and the output dimension
 360 of each base classifier, respectively. Define Q , K and V as the set of q , k and v , respectively,
 361 where $Q = [q_1, q_2, \dots, q_N]^T = \Phi \cdot W^Q \in \mathbb{R}^{N \times C}$, $K = [k_1, k_2, \dots, k_N]^T = \Phi \cdot W^K \in \mathbb{R}^{N \times C}$, and
 362 $V = [v_1, v_2, \dots, v_N]^T = \Phi \cdot W^V \in \mathbb{R}^{N \times C}$. Here, $W^Q \in \mathbb{R}^{C \times C}$, $W^K \in \mathbb{R}^{C \times C}$, and $W^V \in \mathbb{R}^{C \times C}$
 363 are the weight of the linear transformation matrix.

364 2) *Interaction Calculation*: In the second stage, we need to score all the base classifiers
 365 based on the results of a certain classifier, and this score determines the degree of interactions

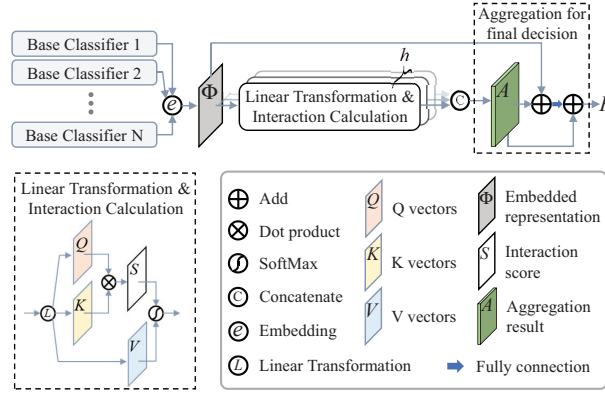


Figure 4: Illustration of the proposed MHSA voting scheme. It includes linear transformation, interaction calculation, and aggregation for the final decision three parts. The linear transformation part transforms the outputs of based classifiers into three vectors Q , K , and V . The interactions among base classifiers are calculated based on Q , K , and V by the interaction calculation part. Then, these interactions are adopted to enhance the representation and generate the final decision.

among this classifier and other base classifiers. The similarity between each pair of base classifiers is calculated by the dot product of K and Q , namely QK^T . Then a SoftMax function is used to normalize the similarity QK^T , and get an interaction score $S \in \mathbb{R}^{N \times N}$ which can reflect the interactions among the base classifiers.

$$S = \begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1N} \\ s_{21} & s_{22} & \cdots & s_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ s_{N1} & s_{N2} & \cdots & s_{NN} \end{bmatrix} = \text{SoftMax}\left(\frac{QK^T}{\sqrt{d}}\right) \quad (5)$$

where s_{ij} represents the interaction between the q_i and k_j , \sqrt{d} can make the MHSA voting scheme have a more stable gradient flow during the training process. After that, the V is multiplied by S , which means maintaining the relationship among the associated base classifiers and reducing the impact of the less-correlated classifiers.

3) *Aggregation & Final Decision*: To learn interaction information in different representation subspaces, the above two stages are performed several times, and the results of these times are concatenated and linearly transformed.

$$A = \text{Aggregation}(Q, K, V) = \text{Concat}(S_1V, \dots, S_hV) \cdot W^A \quad (6)$$

where A is the aggregation result, $S_k (k = 1, \dots, h)$ indicates interaction score in different representation subspaces, $W^A \in \mathbb{R}^{C \times C}$ is the linear transformation matrix, *Concat* is channel

381 concatenation. Then the aggregation result is passed through the residual connection and
382 the FC layer to enhance the representation, and get the final decision, which can be described
383 by the following equation:

$$384 \quad F = FC(A + \Phi) + A \quad (7)$$

385 where F is the final decision generated by the MHSA voting. The MHSA voting can achieve
386 the modeling of the interactions among the base classifiers and fuse the outputs of each base
387 classifier based on these interactions.

388 II.E. Implementations

389 The proposed ensemble framework is implemented based on the PyTorch deep learning
390 library. The framework is trained on a PC with an NVIDIA GTX 1080Ti graphics card.
391 The loss function in Eq. (4) is adopted to supervise the learning of the base classifiers
392 parameters, which are optimized by the Adam optimizer with a low learning rate of 0.0001.

393 To validate the proposed framework, a series of comparison and ablation experiments are
394 conducted. In the comparison experiments, several ML-based and DL-based methods were
395 compared with the proposed framework to demonstrate the superiority of our framework.
396 Since all results acquired by different methods are measured based on the same ADNI cohort,
397 and most of these methods have similar pre-processing pipeline and implementation details
398 to that in the proposed method, we compare our results with the reported results by the
399 compared methods. In the ablation experiments, the effectiveness of the entropy-based
400 MIL pooling layer and MHSA voting scheme, several studies are conducted to evaluate the
401 influence of transfer learning and clinical information, and the indispensability of four base
402 classifiers. More details about the implementations are as follows.

403 II.E.1. Data split

404 20% samples (154 samples) of the dataset are selected as the test samples and the remaining
405 80% samples (617 samples) as the training samples. A five-fold cross-validation strategy
406 is adopted to verify the reliability of the proposed framework, in which four folds of the
407 training samples are used for training and one fold for validation. To make sure that no
408 significant difference in the age and gender distributions among the training, validation, and
409 test samples, the Chi-square test is used to verify the distributions.

410 II.E.2. Training strategy

411 For task 1 (i.e., AD *vs.* CN), the base CNN classifiers are trained from scratch directly, and
 412 the parameters of them are initialized randomly. For task 2 (i.e., pMCI *vs.* sMCI), transfer
 413 learning is adopted to train the base CNN classifiers. MCI is a preclinical stage of AD,
 414 the structural changes of brains caused by MCI may be more subtle than those caused by
 415 AD, which means task 2 is more challenging than task 1. According to the development
 416 of AD, the two tasks are highly correlated, and the information learned from AD and CN
 417 subjects can be employed as a supplement to enrich the information for task 2^{28,29}. Thus,
 418 the parameters of base CNN classifiers trained on task 1 are transferred to initialize the
 419 training for task 2. Early stopping is applied for all training processes, the training process
 420 is terminated when the validation loss exceeds the lower threshold in 10 continuous epochs.

421 II.E.3. Evaluation metrics

422 In two classification tasks, four evaluation metrics, namely, classification accuracy (ACC),
 423 sensitivity (SEN), specificity (SPE), and the area under the receiver operating characteristic
 424 curve (AUC) are adopted to evaluate the classification performance. These metrics are
 425 respectively defined as:

$$426 \quad ACC = \frac{TP + TN}{TP + TN + FP + FN} \quad (8)$$

$$427 \quad SEN = \frac{TP}{TP + FN} \quad (9)$$

$$428 \quad SPE = \frac{TN}{TN + FP} \quad (10)$$

429 where TP denotes true positive, TN denotes true negative, FP denotes false positive, and FN
 430 denotes false negative. The ROC curve is generated according to the (SEN, 1−SPE) pairs.
 431 The AUC characterizes the classification performance of the methods, the performance is
 432 better when AUC is closer to 1.

III. RESULTS

III.A. Comparison with other methods

To demonstrate the superiority of the proposed ensemble framework, we compare the results on two tasks of our method and other methods. The classification results on ADNI dataset are summarized in Table 2.

Table 2: Comparison of the proposed method with the existing state-of-the-art methods reported in the literature.

Methods	Data	AD <i>vs.</i> CN				pMCI <i>vs.</i> sMCI				
		ACC	SEN	SPE	AUC	ACC	SEN	SPE	AUC	
ML-based	Moradi et al. ²⁰	sMRI + Clinical info	-	-	-	-	82.00%	87.00%	74.00%	90.00%
	Beheshti et al. ²¹	sMRI	93.01%	89.13%	96.80%	93.51%	75.00%	76.92%	73.23%	75.08%
	Calvini et al. ²³	sMRI	-	74.00%	85.00%	86.30%	-	-	-	-
	Koikkalainen et al. ²⁴	sMRI	86.00%	81.00%	91.00%	-	72.10%	77.00%	71.00%	-
	Liu et al. ²⁵	sMRI	93.06%	94.85%	90.49%	95.79%	79.25%	87.92%	75.54%	83.44%
	Shi et al. ²⁶	sMRI + PET + CSF	95.00%	95.30%	94.70%	93.20%	-	-	-	-
	Tong et al. ²⁷	sMRI	90.00%	86.00%	93.00%	-	72.00%	69.00%	74.00%	-
	Coupe et al. ²⁸	sMRI	91.00%	87.00%	94.00%	-	74.00%	73.00%	74.00%	-
DL-based	Suk et al. ¹⁸	sMRI + PET	95.35%	94.65%	95.22%	98.77%	75.92%	48.04%	95.23%	74.66%
	Shi et al. ⁴³	sMRI + PET	97.13%	95.93%	98.53%	97.20%	78.88%	68.04%	86.81%	80.10%
	Liu et al. ⁴⁴	sMRI + PET	91.40%	92.32%	90.42%	-	-	-	-	-
	Cui et al. ⁴⁵	sMRI	92.29%	90.63%	93.72%	96.95%	75.00%	73.33%	76.19%	79.70%
	Liu et al. ²⁹	sMRI	91.09%	88.05%	93.50%	95.86%	76.90%	42.11%	82.43%	77.64%
	Kang et al. ³⁴	sMRI	90.40%	-	-	-	66.70%	-	-	-
	Lian et al. ⁴⁶	sMRI	90.30%	82.40%	96.50%	95.10%	80.9%	52.60%	85.40%	78.10%
	Chen et al. ⁴⁷	sMRI	95.32%	91.18%	93.94%	-	77.60%	71.62%	75.85%	-
	Zhang et al. ⁴⁸	sMRI	93.20%	92.40%	94.00%	96.10%	82.90%	90.00%	75.70%	86.50%
	Basaia et al. ²²	sMRI + PET + CSF	93.20%	93.00%	93.30%	-	-	-	-	-
Proposed	sMRI + Clinical info	98.61%	98.54%	98.67%	99.08%	84.49%	83.50%	81.48%	85.69%	

In the task of AD *vs.* CN, the best ACC, SEN, SPE, and AUC values implemented by previous works are respectively 97.13%, 95.93%, 98.53%, and 98.77%, which are realized by the works of Shi et al.⁴³ and Suk et al.¹⁸ The proposed method has the ACC of 98.61%, the SEN of 98.54%, the SPE of 98.67%, and the AUC of 99.08%, which are respectively 1.48%, 2.61%, 0.14%, and 0.31% higher than the best metrics achieved by other methods. In the task of pMCI *vs.* sMCI, the values of ACC, SEN, SPE, and AUC obtained by the proposed framework are respectively 84.49%, 83.50%, 81.48%, and 85.69%. Our method achieves the best prediction accuracy, which is 1.59% higher than the best ACC obtained by Zhang et al.⁴⁸ These results show that the proposed framework can indeed yield a more accurate diagnosis, and have satisfactory performance on other evaluation metrics.

Last edited *Date* :

450 III.B. Effectiveness of entropy-based MIL pooling

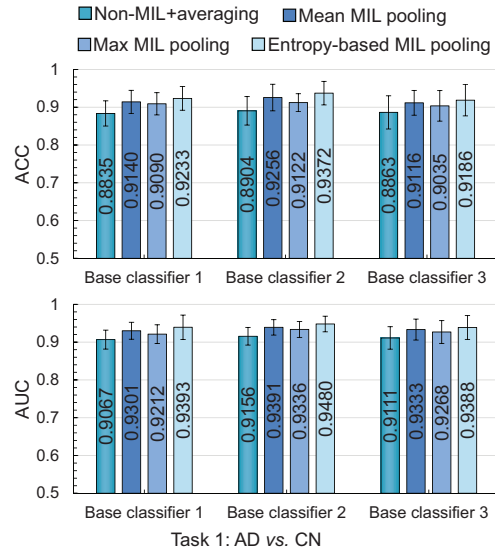
451 To evaluate the effectiveness of entropy-based MIL pooling, we compare the results of base
452 classifiers without MIL pooling and with different MIL pooling layers. The compared meth-
453 ods include non-MIL+averaging method, Mean MIL pooling method, and Maximum pooling
454 method. The non-MIL+averaging method has the same architectures as base CNN classi-
455 fiers except no MIL pooling, and performs classification through averaging the slice-level
456 results. Both mean MIL pooling method and maximum pooling method also have the same
457 architectures as base CNN classifiers, only replacing the entropy-based MIL pooling layer.
458 The classification results in terms of ACC and AUC for two tasks are shown in Fig. 5.

459 From Fig. 5, it can be learned that MIL methods (i.e., mean MIL pooling, max MIL
460 pooling, and entropy-based MIL pooling) yield better results in terms of ACC and AUC.
461 Taking the base classifier 1 as an example, the ACC and AUC achieved by MIL methods
462 are on average higher 0.0319 and 0.0247 than non-MIL method in task 1, and higher 0.0199
463 and 0.0249 in task 2. Compared with mean MIL pooling and max MIL pooling methods,
464 the proposed entropy-based MIL pooling achieves the best results on both tasks, which can
465 reach 0.9372 ACC and 0.9480 AUC on task 1 (achieved by base classifier 2), 0.7959 ACC
466 and 0.8081 AUC on task 2 (achieved by base classifier 1). The above results reflect that the
467 MIL methods can improve the classification performance than the non-MIL method, and
468 confirm that the entropy-based MIL pooling method is more effective than the normal MIL
469 methods, which shows the effectiveness of entropy-based MIL pooling.

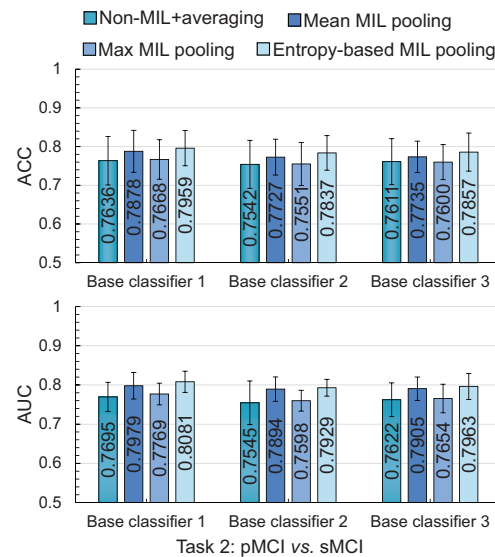
470 III.C. Effectiveness of MHSA voting

471 A key component of the proposed ensemble framework is the ensemble approaches to fuse
472 the base classifiers. We conduct the experiments to verify the effectiveness of MHSA voting.
473 Specifically, base classifier 1, base classifier 2, and base classifier 3 are fused via different
474 voting approaches including majority voting (MV), weighted voting (WV), SVM-based vot-
475 ing (SVM), and the proposed MHSA voting. Table 3 reports the corresponding results of
476 different ensemble approaches.

477 From Table 3, it can be observed that two learnable ensemble approaches (i.e., SVM-
478 based voting, and MHSA voting) yield better classification performance on two tasks than
479 unlearnable approaches (i.e., majority voting, and weighted voting). In the task of AD
480 *vs.* CN, the results obtained by majority voting and weighted voting are lower than the
481 maximum values of ACC and AUC (achieved by base classifier 2) before fusion. And the



(a) Classification results in terms of ACC and AUC for task 1



(b) Classification results in terms of ACC and AUC for task 2

Figure 5: Classification results in terms of ACC and AUC achieved by three base CNN classifiers with different MIL pooling layers for two tasks, i.e., AD vs. CN, and pMCI vs. sMCI. The error bars denote the standard deviations of the results.

482 results obtained by SVM-based voting are basically consistent with the maximum values
 483 before fusion. Only the MHSA voting achieves an improvement in results, with the ACC of
 484 0.9419, and the AUC of 0.9545, which is at least 0.0047 higher than the metrics generated
 485 by base classifiers. In the task of pMCI vs. sMCI, all ensemble approaches can obtain better
 486 results than that before fusion. The results obtained via majority voting have the minimum
 487 improvement, with the ACC of 0.8061, and the AUC of 0.8165. The maximum improvement

Table 3: Classification results of different ensemble approaches on two tasks.

Ensemble Members	Ensemble Approach	AD <i>vs.</i> CN		pMCI <i>vs.</i> sMCI	
		ACC	AUC	ACC	AUC
Base classifier 1	-	0.9233 \pm 0.0314	0.9393 \pm 0.0323	0.7959 \pm 0.0454	0.8081 \pm 0.0271
Base classifier 2	-	0.9372 \pm 0.0311	0.9480 \pm 0.0207	0.7837 \pm 0.0447	0.7929 \pm 0.0251
Base classifier 3	-	0.9186 \pm 0.0416	0.9388 \pm 0.0315	0.7857 \pm 0.0492	0.7963 \pm 0.0332
Base classifier 1, 2, 3	MV	0.9279 \pm 0.0283	0.9306 \pm 0.0301	0.8061 \pm 0.0366	0.8165 \pm 0.0318
	WV	0.9302 \pm 0.0245	0.9415 \pm 0.0202	0.8265 \pm 0.0409	0.8316 \pm 0.0431
	SVM	0.9349 \pm 0.0209	0.9478 \pm 0.0199	0.8286 \pm 0.0422	0.8367 \pm 0.0395
	MHSA	0.9419 \pm 0.0232	0.9545 \pm 0.0205	0.8408 \pm 0.0350	0.8535 \pm 0.0283

Data are mean \pm standard deviation.

MV: majority voting; WV: weighted voting; SVM: SVM-based voting; MHSA: MHSA voting.

488 on results is achieved by MHSA voting, which is at least 0.0449 higher than the metrics
 489 generated by base classifiers. Compared with these common ensemble approaches, MHSA
 490 voting can further improve the effects of fusion. These results confirm the effectiveness of
 491 using MHSA voting.

492 III.D. Influence of transfer learning

493 To demonstrate the impact of transfer learning, we compare the experimental results with
 494 and without transfer learning. In this group of experiments, we train base classifiers from
 495 scratch for task 2 without adopting transfer learning strategy, and compare their classification
 496 performance with that obtained by base classifiers trained with transfer learning strategy.
 497 Fig. 6 shows the classification results in terms of ACC and AUC for task 2.

498 As shown in Fig. 6, it can be seen that transfer learning strategy significantly improves
 499 the classification performance. Take base classifier 1, 2, 3, 4 fused via MHSA voting as an
 500 example, with the aid of transfer learning, it improves the ACC from 0.8106 to 0.8449, the
 501 AUC from 0.8196 to 0.8569, which has at least a 4.23% boost. Meanwhile, other methods
 502 trained with transfer learning have higher gain percentages, the ACC has an average gain
 503 of 5.65%, and the AUC has an average gain of 6.13%. These results indicate that the use of
 504 transfer learning strategy can indeed improve the classification performance on task 2.

505 III.E. Influence of clinical information

506 As introduced in Section II.C.5., base classifier 4 (i.e., MLP) is chosen for clinical information
 507 analysis. Base classifier 4 is fused with other base classifiers via MHSA voting to construct
 508 a multi-model ensemble framework. To investigate the influence of clinical information, we
 509 compare the classification performance achieved by *Only Clin info* (base classifier 4), *Without*

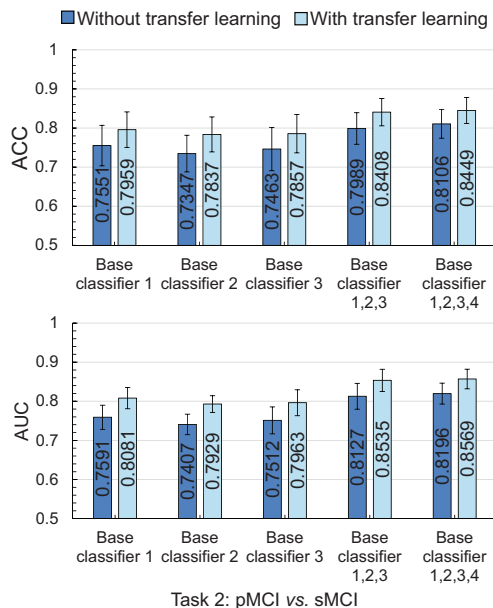


Figure 6: Classification results in terms of ACC and AUC achieved by base classifiers trained without and with transfer learning for task 2. Base classifier 1, 2, 3 and base classifier 1, 2, 3, 4 are fused via MHSA voting. The error bars denote the standard deviations of the results.

Clin info (base classifier 1, 2, 3 fused via MHSA voting), and *With Clin info* (base classifier 1, 2, 3, 4 fused via MHSA voting). The corresponding results are as demonstrated in Table 4.

Table 4: Classification results of different methods with and without clinical information.

Task	Method	ACC	SEN	SPE	AUC
AD vs. CN	<i>Only clin info.</i>	0.9837 ± 0.0204	0.9756 ± 0.0209	0.9778 ± 0.0298	0.9848 ± 0.0167
	<i>Without clin info.</i>	0.9419 ± 0.0232	0.9268 ± 0.0311	0.9556 ± 0.0199	0.9545 ± 0.0205
	<i>With clin info.</i>	0.9861 ± 0.0182	0.9854 ± 0.0233	0.9867 ± 0.0221	0.9908 ± 0.0143
pMCI vs. sMCI	<i>Only clin info.</i>	0.6980 ± 0.0870	0.7255 ± 0.0744	0.7037 ± 0.0661	0.6987 ± 0.0598
	<i>Without clin info.</i>	0.8408 ± 0.0350	0.8273 ± 0.0422	0.8234 ± 0.0406	0.8535 ± 0.0283
	<i>With clin info.</i>	0.8449 ± 0.0332	0.8350 ± 0.0405	0.8148 ± 0.0341	0.8569 ± 0.0214

Data are mean ± standard deviation.

As shown in Table 4, for the task of AD vs. CN, the use of clinical information can significantly improve the diagnosis performance. Compared with the results achieved by *Without Clin info*, *With Clin info* improves the ACC from 0.9419 to 0.9861, the SEN from 0.9268 to 0.9854, the SPE from 0.9556 to 0.9867, and the AUC from 0.9545 to 0.9908. And the quantification biases of ACC, SEN, and AUC obtained by *With Clin info* are smaller than that of *Without Clin info*. *Only Clin info* can obtain similar performance to *With Clin*

519 *info* in terms of ACC. However, the SEN, SPE, AUC achieved by *Only Clin info* are lower
 520 than that achieved by *With Clin info*, and the quantification biases of these metrics are
 521 also higher. For the task of pMCI *vs.* sMCI, the use of clinical information also improves
 522 diagnosis performance, but not as significantly as the task of AD *vs.* CN. *With Clin info*
 523 yields better results, with the ACC of 0.8449, the AUC of 0.8569, which are higher than that
 524 obtained by the other two methods. Though *Without Clin info* yields similar performance
 525 to *With Clin info*, the quantification biases of all metrics are higher than that obtained by
 526 *With Clin info*. The above results reveal that the use of clinical information can provide
 527 better classification performance, and reduce the quantification bias of diagnosis.

528 III.F. Indispensability of four base classifiers

529 To prove the indispensability of the four types of base classifiers, we summarize and compare
 530 the classification performance of fused different types of base classifiers. Specifically, base
 531 classifier 1, 2, 3, and 4 are randomly fused by MHSA voting. The corresponding results for
 532 task 1 and task 2 are reported in Table 5, and some of the ROC curves for the two tasks are
 533 respectively represented in Fig. 7.

Table 5: Classification results of fused different types of base classifiers.

No. of Cls	Members	AD <i>vs.</i> CN		pMCI <i>vs.</i> sMCI	
		ACC	AUC	ACC	AUC
1	Base classifier 1	0.9233 ± 0.0314	0.9393 ± 0.0323	0.7959 ± 0.0454	0.8081 ± 0.0271
	Base classifier 2	0.9372 ± 0.0311	0.9480 ± 0.0207	0.7837 ± 0.0447	0.7929 ± 0.0251
	Base classifier 3	0.9186 ± 0.0416	0.9388 ± 0.0315	0.7857 ± 0.0492	0.7963 ± 0.0332
	Base classifier 4	0.9837 ± 0.0204	0.9848 ± 0.0167	0.6980 ± 0.0870	0.6987 ± 0.0598
2	Base classifier 1, 2	0.9396 ± 0.0276	0.9539 ± 0.0191	0.7999 ± 0.0371	0.8102 ± 0.0298
	Base classifier 1, 3	0.9253 ± 0.0291	0.9405 ± 0.0198	0.8018 ± 0.0466	0.8143 ± 0.0248
	Base classifier 2, 3	0.9380 ± 0.0323	0.9485 ± 0.0212	0.7993 ± 0.0322	0.8036 ± 0.0231
	Base classifier 1, 4	0.9847 ± 0.0197	0.9863 ± 0.0155	0.7967 ± 0.0507	0.8098 ± 0.0336
	Base classifier 2, 4	0.9856 ± 0.0184	0.9902 ± 0.0152	0.7901 ± 0.0581	0.8003 ± 0.0364
	Base classifier 3, 4	0.9847 ± 0.0187	0.9866 ± 0.0152	0.7896 ± 0.0482	0.7998 ± 0.0342
3	Base classifier 1,2,3	0.9419 ± 0.0232	0.9545 ± 0.0205	0.8408 ± 0.0350	0.8535 ± 0.0283
	Base classifier 1,2,4	0.9855 ± 0.0265	0.9902 ± 0.0144	0.8059 ± 0.0382	0.8154 ± 0.0350
	Base classifier 1,3,4	0.9841 ± 0.0227	0.9862 ± 0.0156	0.8122 ± 0.0394	0.8205 ± 0.0312
	Base classifier 2,3,4	0.9852 ± 0.0197	0.9900 ± 0.0137	0.8041 ± 0.0435	0.8181 ± 0.0344
4	Base classifier 1, 2, 3, 4	0.9861 ± 0.0182	0.9908 ± 0.0143	0.8449 ± 0.0332	0.8569 ± 0.0214

Data are mean ± standard deviation.

534 From Table 5, when four base classifiers are fused, the best classification results can
 535 be obtained, the values of ACC for task 1 and task 2 are respectively 0.9861 and 0.8449,
 536 the values of AUC are respectively 0.9908 and 0.8569. And the quantification bias is also
 537 satisfactory. Base classifier 1, base classifier 2, and base classifier 3 have similar performance

538 on both tasks. Base classifier 4 (i.e., MLP) achieves great performance on task 1, while
 539 it performs not good on task 2. When two base CNN classifiers are randomly fused, the
 540 classification results are similar to that achieved by a single CNN classifier, and the quanti-
 541 fication biases are lower. Due to the influence of clinical information, any base CNN classifier
 542 (i.e., base classifier 1, 2, and 3) fused with base classifier 4 could further boost the diagnosis
 543 performance, especially in the task of AD *vs.* CN. Though one base CNN classifier fused
 544 with base classifier 4 can improve the ACC and AUC, the quantification biases of them are
 545 higher than that before fused with base classifier 4. When three base classifiers are randomly
 546 fused, the fusions that include base classifier 4 can yield satisfactory results in the task of
 547 AD *vs.* CN, which are better than the fusions only including base CNN classifiers. In the
 548 task of pMCI *vs.* sMCI, we can see that the fusions only including base CNN classifiers have
 549 the better performance than that fusions including base classifier 4.

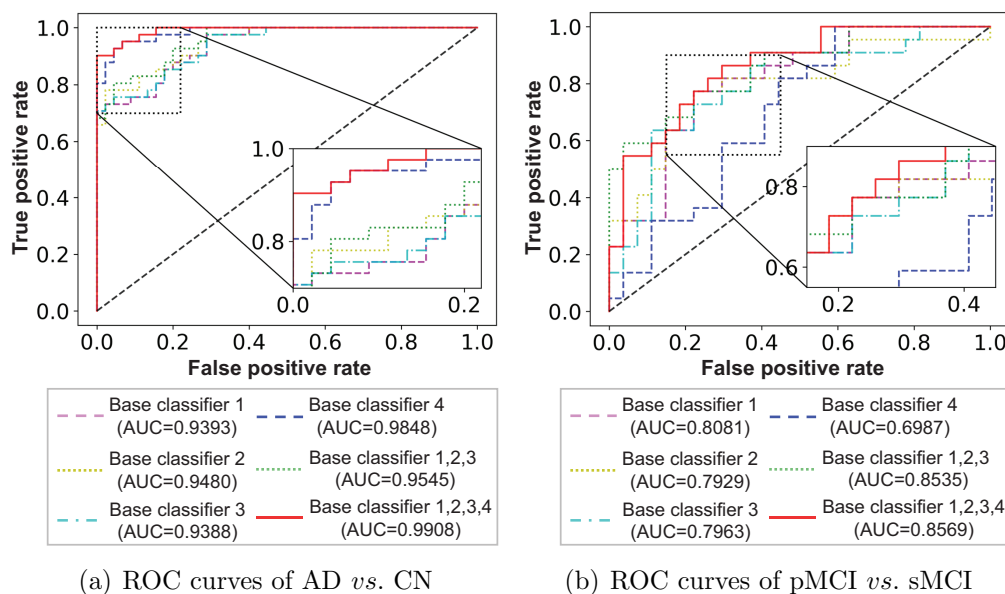


Figure 7: Comparison of the ROC curves. The ROC curves are obtained by base classifier 1, base classifier 2, base classifier 3, base classifier 4, the fusion of base classifier 1, 2, 3, and the fusion of base classifier 1, 2, 3, 4. The upper left area of ROC curve is zoomed for clarity.

550 From Fig. 7, it can be learned that the fusion of four base classifiers has better ROC
 551 curves than others. The results in Table 5 and Fig. 7 illustrate that the fusion of these base
 552 classifiers can achieve better diagnosis performance than a single classifier, and each base
 553 classifier could play an important part in the ensemble framework.

IV. DISCUSSIONS

This work presents a reliable ensemble framework to diagnose AD and MCI using neural networks. MHSA voting improves the fusion of base classifiers in the ensemble, and entropy-based MIL strategy could use more effective information contained in sMRI. Overall, the proposed method provides the reliable diagnosis of AD and prediction of MCI conversion. We built our method based on ensemble learning for several reasons. First, though DL-based methods have been shown to surpass human experts in predictive accuracy, they tend to exhibit higher variance, especially when only a single DL model is adopted. However, reliable diagnosis is needed in the clinic, high variance makes it hard for a single model to generate convincing judgments. In contrast to a single DL model, ensemble learning that combines the outputs of multiple DL models has been proven to achieve better outcomes and generalizability⁴⁹, which is more applicable in clinical settings. Second, because the characteristics of AD are concealed, slow, and non-lethal, the collection of samples is difficult, often resulting in the limited number of samples. The limited number of samples may lead to over-fitting or inadequate training of a model, and limit the identification of complex AD patterns. Ensemble learning has the power in dealing with these challenges³².

We compared the performance of the proposed method against several ML-based and DL-based methods. In all compared methods, MRI images were preprocessed through a similar pipeline to this work, including motion correction, intensity correction, skull stripping, and normalization. Following this basic pipeline, different methods then performed some specific operations (e.g., tissue segmentation and slices sampled in this work) to generate slices, regions, or patches of the brain according to the needs of these methods. In addition, cross-validation and corresponding data split were also adopted in most of the compared methods^{18,20,21,25–28,34,43–45,47}, and they took the average of the cross-validation results as the final performance. These means such as preprocessing procedures or cross-validation are a part of the compared methods and have no impact on demonstrating the effectiveness of the proposed method. As the results shown in Section III.A., our method significantly outperformed the compared methods in classification accuracy for both tasks (AD *vs.* CN, pMCI *vs.* sMCI). Noting that some compared methods^{18,20,24,25,27,29,43,46,48} had quite unbalanced SEN and SPE, the imbalance of SEN and SPE indicates that the missed diagnosis or misdiagnosis rate of these methods was high. A previous work²⁹ achieved SEN of 42.11%, and the SPE of 82.43% in the task of pMCI *vs.* sMCI, which means only 42.11% pMCI patients were correctly diagnosed and 17.57% sMCI patients were misdiagnosed. The proposed method achieved balanced and satisfactory SEN and SPE for both tasks, which demonstrates that

IV. DISCUSSIONS

588 our method can conduct a reliable diagnosis. Furthermore, the five-fold cross-validation ap-
589 proach has been performed in this work. The mean values and standard deviation of ACC
590 and AUC are as demonstrated in Table 5. The proposed method achieved the best results
591 on both tasks, which had the ACC of 0.9861 ± 0.0182 and the AUC of 0.9908 ± 0.0143 on
592 AD *vs.* CN task, the ACC of 0.8449 ± 0.0332 and the AUC of 0.8569 ± 0.0214 on pMCI
593 *vs.* sMCI task. The quantification biases of these metrics were effectively reduced by the
594 use of ensemble learning, which was lower than that of each base classifier. The results with
595 low quantification bias generated by our method indicate that the proposed method is able
596 to generate a robust diagnosis, which is also in good agreement with the effect of ensemble
597 learning.

598 In this work, MHSA voting is proposed to aggregate the outputs of base classifiers as
599 previous studies^{34,36} typically adopted the common voting approaches which ignore the inter-
600 actions among base classifiers. The majority voting, weighted voting, and SVM-based voting
601 are commonly used for the aggregation in the ensemble. However, these common voting ap-
602 proaches sometimes may cause a decrease or stay flat on results after fusion. The reason for
603 this is that majority voting and weighted voting are not data-adaptive, they assign the fixed
604 weight to each base classifier. Though SVM-based voting is a learnable ensemble approach,
605 it leaves the interactions among the ensemble members out of consideration. MHSA voting
606 has been shown to have an improvement on results after fusion. This implies that the in-
607 teractions among the base classifiers can play a role during their fusion, and MHSA voting
608 can exploit the interactions to generate better classification results during the fusion of base
609 classifiers.

610 While MIL strategy has been applied in the diagnosis of different diseases, to our knowl-
611 edge, rare studies have explored it in the diagnosis of AD based on slice-level. We incorpo-
612 rated entropy-based MIL strategy into base CNN classifiers to use more effective information
613 contained in sMRI. As shown in Fig. 5, MIL strategy can indeed further improve the perfor-
614 mance of both tasks in contrast to non-MIL methods, in which the proposed entropy-based
615 MIL strategy has been shown to achieve the best classification results. Due to AD-related
616 pathological areas having the uneven distribution in sMRI, non-MIL methods are easily af-
617 fected, thereby resulting in sub-optimal performance in two tasks. Compared with non-MIL
618 methods, MIL methods consider the relationships between slices, which is beneficial to im-
619 proving the utilization of information contained in sMRI. The normal MIL methods (i.e.,
620 mean MIL method, and maximum MIL method) consider that the relationships between
621 slices have no difference, and the slices have similar feature expression abilities. Neverthe-
622 less, the slices with abundant tissue information are generally getting more attention in

623 clinical diagnosis, and radiologists also focus on these slices. Similar to the habit of ra-
 624 diologists' review of MRI, the proposed entropy-based MIL method measures the feature
 625 expression abilities of different slices according to their information entropy, which can gen-
 626 erate more reasonable embedding-level features for further classification. And therefore, the
 627 entropy-based MIL method has better performance than normal MIL methods.

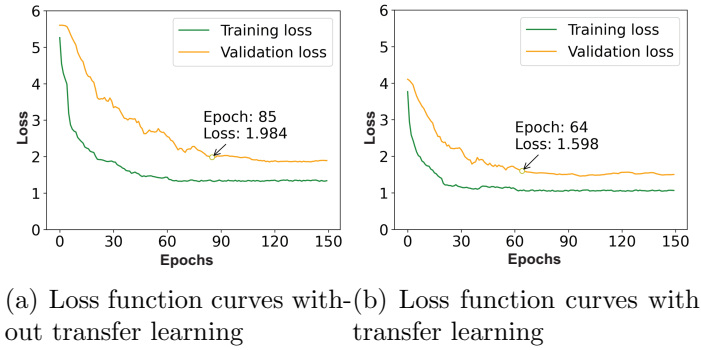


Figure 8: Comparison of the loss function curves achieved by training without and with transfer learning for task 2.

628 Transfer learning improved the classification results in terms of ACC and AUC by
 629 $\sim 4\%$ across two tasks. This situation is consistent with existing studies^{28,29}. The results
 630 demonstrate that the two tasks are correlated, and the supplementary information from
 631 AD and CN subjects implicitly enriches the features in the task of pMCI *vs.* sMCI during
 632 training. In addition, we also analyzed the influence of transfer learning on training duration.
 633 Here, we trained the proposed method for task 2 without early stopping and set the epochs to
 634 150. Fig. 8 shows the loss function curves with and without transfer learning during training.
 635 As observed in Fig. 8, the training loss has a faster downward trend than the validation loss,
 636 and after the convergence of training, the validation loss is slightly more than the training
 637 loss. With transfer learning, the initial values of training and validation losses (epoch 1)
 638 were lower than that without transfer learning, and the validation loss converged to about
 639 1.6 after epoch 64. The validation loss converged to about 2.0 after epoch 85 when transfer
 640 learning strategy was not adopted. These results show that the model can fit the data better
 641 and faster when using transfer learning. In this work, early stopping was adopted with the
 642 patience of 10 epochs on the validation loss, and the training time was five minutes per
 643 epoch. For the task of pMCI *vs.* sMCI, the training lasted about 6 hours, which can save
 644 about 1.7 hours in contrast to that training without transfer learning. For the task of AD
 645 *vs.* CN, the training time was about 7.5 hours.

646 As different imaging modalities and clinical data can provide various information about
 647 AD patients, we adopted multimodal data (sMRI and clinical information) to develop an
 648 ensemble framework. In this work, the use of multimodal data led to an overall improvement

IV. DISCUSSIONS

649 in both tasks, which improved the diagnosis performance and reduced the quantification
650 bias. From Table 2, it can be observed that most studies using multimodal data have
651 better performance than the studies using single modal data. Moreover, we analyzed the
652 sensitivity of clinical information to two tasks. For the task of AD *vs.* CN, the use of clinical
653 information only can also obtain satisfactory performance, while for the task of pMCI *vs.*
654 sMCI, the use of clinical information only cannot achieve good results. These results show
655 that the clinical information is more sensitive to the task of AD *vs.* CN than that to the task
656 of pMCI *vs.* sMCI. It can be also learned that cognitive and neuropsychological measures in
657 clinical information change greatly from normal cognition to dementia, and these measures
658 have no significant change in CN or MCI stages. This inference is consistent with previous
659 research^{50,51}.

660 AD is an irreversible neurodegenerative disease with concealed, slow, and non-lethal
661 characteristics, which is also a serious social problem. The dementia symptoms caused
662 by AD gradually worsen over several years. In general, a person with AD lives 4 to 8
663 years after diagnosis but can live as long as 20 years, depending on other factors (e.g.,
664 earlier diagnosis or intervention). At present, AD has no cure, some treatments can only
665 temporarily slow the worsening of dementia symptoms and improve the quality of life for AD
666 patients and their caregivers. Earlier diagnosis of AD is crucial for prolonging the lifespan
667 and improving the quality of life for those with AD. Our proposed ensemble framework is
668 able to generate reliable and robust results for the diagnosis of AD and the prediction of
669 MCI conversion, which has great practical significance for the earlier diagnosis of AD. The
670 detailed analyses of the results give an important indication that the proposed ensemble
671 framework can potentially be employed in the reliable diagnosis of AD and prediction of
672 MCI conversion. Furthermore, due to the characteristics of AD, the collection of AD samples
673 is difficult in clinical settings. With ensemble learning, the dilemma caused by the limited
674 number of samples can be solved to some extent³². The proposed method is based on
675 ensemble learning, which makes our method potential to perform reliable diagnoses under
676 limited data, thereby reducing the burden of physicians collecting data. In many clinical
677 settings, because it is difficult to identify the exact cause of dementia, multiple diagnostic
678 tests are typically adopted to determine if a person has AD, including brain imaging, mental
679 cognitive status tests, etc. To closer to practical clinical application and obtain a more
680 reliable diagnosis, we also adopted the multimodal data in this work. It is worth noting that
681 our method can also achieve satisfactory results only using sMRI.

682 This current work has some limitations despite its successful performance in AD diag-
683 nosis and MCI conversion prediction. The black-box nature is a common limitation of deep

684 learning methods, which is also the main reason that limits the widespread application of
685 medical artificial intelligence (AI). In clinical settings, to determine whether a person suffers
686 from a certain disease, it needs to undergo a detailed clinical examination, and the physicians
687 confirm the condition of this person according to the clinical test results. In this process,
688 the basis for the diagnosis is detailed and clear. For medical AI, the details of algorithmic
689 decision-making should also be exposed like clinical diagnosis, which is currently difficult.
690 Note that conceptual understanding and experiences owned by physicians are impossible for
691 AI to fully learn. To deploy an explainable AI in medical practices, it still requires the nec-
692 essary human oversight⁵². The interactive deep learning with the “human in the loop” can
693 be potentially considered as a robust way to handle explainability. This human-in-the-loop
694 deep learning combines the conceptual understanding and experiences owned by physicians
695 with the effectiveness of deep learning, which can ensure that decision-making is controllable
696 and clinically justified. As a high level of accountability is required in the medical field, ma-
697 chined decisions and predictions need to be explained clearly, our future work will include
698 exploring human-in-the-loop deep learning.

699 V. CONCLUSION

700 In this paper, a robust ensemble framework is proposed for reliable diagnosis of AD and
701 prediction of MCI conversion. Specifically, three base CNN classifiers with different scales
702 of network width, depth, and resolution are designed to capture detailed features in sMRI.
703 To better use effective information contained in sMRI, we incorporate entropy-based MIL
704 strategy into base CNN classifiers, which can take the information densities of slices into
705 account to generate more reasonable features for classification. Additionally, one shallow
706 classifier (i.e., MLP) is employed to analyze the clinical information. The final diagnosis is
707 achieved by MHSA voting approach that aggregates the predictions of base classifiers while
708 considering the interactions among them. Extensive experimental results on ADNI database
709 show that the proposed ensemble framework has reliable and competitive performance in
710 both tasks.

711 VI. ACKNOWLEDGMENTS

712 This work was supported by the Young Scientist Studio of Harbin Institute of Technology.
713 Data used in preparation of this article were obtained from the Alzheimer’s Disease Neu-

VI. ACKNOWLEDGMENTS

714 roimaging Initiative (ADNI) database (adni.loni.usc.edu). As such, the investigators within
715 the ADNI contributed to the design and implementation of ADNI and/or provided data
716 but did not participate in analysis or writing of this report. A complete listing of ADNI
717 investigators can be found at: [http://adni.loni.usc.edu/wp-content/uploads/how_to_](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf)
718 [apply/ADNI_Acknowledgement_List.pdf](http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf).

719 VII. CONFLICT OF INTEREST DISCLOSURE

720 The authors have no conflicts of interest to declare.

722 References

- 723 ¹ D. E. Barnes and K. Yaffe, The projected effect of risk factor reduction on Alzheimer’s
724 disease prevalence, *The Lancet Neurology* **10**, 819–828 (2011).
- 725 ² R. Wang et al., Generation of synthetic PET images of synaptic density and amyloid
726 from 18F-FDG images using deep learning, *Medical physics* **48**, 5115–5129 (2021).
- 727 ³ M. Grundman et al., Mild cognitive impairment can be distinguished from Alzheimer
728 disease and normal aging for clinical trials, *Archives of neurology* **61**, 59–66 (2004).
- 729 ⁴ R. C. Petersen, R. O. Roberts, D. S. Knopman, B. F. Boeve, Y. E. Geda, R. J. Ivnik,
730 G. E. Smith, and C. R. Jack, Mild cognitive impairment: ten years later, *Archives of*
731 *neurology* **66**, 1447–1455 (2009).
- 732 ⁵ N. I. Bradfield and D. Ames, Mild cognitive impairment: narrative review of tax-
733 onomies and systematic review of their prediction of incident Alzheimer’s disease de-
734 mentia, *BJPsych bulletin* **44**, 67–74 (2020).
- 735 ⁶ A. Nordberg, PET imaging of amyloid in Alzheimer’s disease, *The lancet neurology* **3**,
736 519–527 (2004).
- 737 ⁷ S. Lehericy, M. Marjanska, L. Mesrob, M. Sarazin, and S. Kinkingnehun, Magnetic
738 resonance imaging of Alzheimer’s disease, *European radiology* **17**, 347–362 (2007).
- 739 ⁸ K. Blennow and H. Hampel, CSF markers for incipient Alzheimer’s disease, *The Lancet*
740 *Neurology* **2**, 605–613 (2003).

Last edited *Date* :

- 741 ⁹ C. R. Jack Jr, D. S. Knopman, W. J. Jagust, L. M. Shaw, P. S. Aisen, M. W. Weiner,
742 R. C. Petersen, and J. Q. Trojanowski, Hypothetical model of dynamic biomarkers of
743 the Alzheimer’s pathological cascade, *The Lancet Neurology* **9**, 119–128 (2010).
- 744 ¹⁰ G. B. Frisoni, N. C. Fox, C. R. Jack, P. Scheltens, and P. M. Thompson, The clinical
745 use of structural MRI in Alzheimer disease, *Nature Reviews Neurology* **6**, 67–77 (2010).
- 746 ¹¹ I. Arevalo-Rodriguez, N. Smailagic, M. R. i Figuls, A. Ciapponi, E. Sanchez-Perez, A. Gi-
747 annakou, O. L. Pedraza, X. B. Cosp, and S. Cullum, Mini-Mental State Examination
748 (MMSE) for the detection of Alzheimer’s disease and other dementias in people with
749 mild cognitive impairment (MCI), *Cochrane Database of Systematic Reviews* (2015).
- 750 ¹² S. E. O’Bryant et al., Staging dementia using Clinical Dementia Rating Scale Sum of
751 Boxes scores: a Texas Alzheimer’s research consortium study, *Archives of neurology* **65**,
752 1091–1095 (2008).
- 753 ¹³ P. Doraiswamy, F. Bieber, L. Kaiser, K. Krishnan, J. Reuning-Scherer, and B. Gulanski,
754 The Alzheimer’s Disease Assessment Scale: patterns and predictors of baseline cognitive
755 performance in multicenter Alzheimer’s disease trials, *Neurology* **48**, 1511–1517 (1997).
- 756 ¹⁴ A. Estévez-González, J. Kulisevsky, A. Boltes, P. Otermín, and C. García-Sánchez, Rey
757 verbal learning test is a useful tool for differential diagnosis in the preclinical phase
758 of Alzheimer’s disease: comparison with mild cognitive impairment and normal aging,
759 *International journal of geriatric psychiatry* **18**, 1021–1028 (2003).
- 760 ¹⁵ M. Lin, S. Momin, Y. Lei, H. Wang, W. J. Curran, T. Liu, and X. Yang, Fully auto-
761 mated segmentation of brain tumor from multiparametric MRI using 3D context deep
762 supervised U-Net, *Medical Physics* **48**, 4365–4374 (2021).
- 763 ¹⁶ R. Haweel, A. Shalaby, A. Mahmoud, N. Seada, S. Ghoniemy, M. Ghazal, M. F.
764 Casanova, G. N. Barnes, and A. El-Baz, A robust DWT–CNN-based CAD system
765 for early diagnosis of autism using task-based fMRI, *Medical physics* **48**, 2315–2326
766 (2021).
- 767 ¹⁷ M. Jo and S.-H. Oh, A preliminary attempt to visualize nigrosome 1 in the substantia
768 nigra for Parkinson’s disease at 3T: An efficient susceptibility map-weighted imaging
769 (SMWI) with quantitative susceptibility mapping using deep neural network (QSMnet),
770 *Medical Physics* **47**, 1151–1160 (2020).

- 771 ¹⁸ H.-I. Suk et al., Hierarchical feature representation and multimodal fusion with deep
772 learning for AD/MCI diagnosis, *NeuroImage* **101**, 569–582 (2014).
- 773 ¹⁹ E.-J. Hwang, H.-G. Kim, D. Kim, H. Y. Rhee, C.-W. Ryu, T. Liu, Y. Wang, and G.-H.
774 Jahng, Texture analyses of quantitative susceptibility maps to differentiate Alzheimer’s
775 disease from cognitive normal and mild cognitive impairment, *Medical physics* **43**, 4718–
776 4728 (2016).
- 777 ²⁰ E. Moradi et al., Machine learning framework for early MRI-based Alzheimer’s conver-
778 sion prediction in MCI subjects, *Neuroimage* **104**, 398–412 (2015).
- 779 ²¹ I. Beheshti et al., Classification of Alzheimer’s disease and prediction of mild cognitive
780 impairment-to-Alzheimer’s conversion from structural magnetic resource imaging using
781 feature ranking and a genetic algorithm, *Computers in biology and medicine* **83**, 109–119
782 (2017).
- 783 ²² S. Basaia et al., Automated classification of Alzheimer’s disease and mild cognitive
784 impairment using a single MRI and deep neural networks, *NeuroImage: Clinical* **21**,
785 101645 (2019).
- 786 ²³ P. Calvini et al., Automatic analysis of medial temporal lobe atrophy from structural
787 MRIs for the early assessment of Alzheimer disease, *Medical physics* **36**, 3737–3747
788 (2009).
- 789 ²⁴ J. Koikkalainen et al., Multi-template tensor-based morphometry: application to anal-
790 ysis of Alzheimer’s disease, *NeuroImage* **56**, 1134–1144 (2011).
- 791 ²⁵ M. Liu, D. Zhang, and D. Shen, Relationship induced multi-template learning for di-
792 agnosis of Alzheimer’s disease and mild cognitive impairment, *IEEE transactions on*
793 *medical imaging* **35**, 1463–1474 (2016).
- 794 ²⁶ Y. Shi, H.-I. Suk, Y. Gao, S.-W. Lee, and D. Shen, Leveraging coupled interaction for
795 multimodal Alzheimer’s disease diagnosis, *IEEE transactions on neural networks and*
796 *learning systems* **31**, 186–200 (2019).
- 797 ²⁷ T. Tong et al., Multiple instance learning for classification of dementia in brain MRI,
798 *Medical image analysis* **18**, 808–818 (2014).
- 799 ²⁸ P. Coupé et al., Scoring by nonlocal image patch estimator for early detection of
800 Alzheimer’s disease, *NeuroImage: clinical* **1**, 141–152 (2012).

- 801 29 M. Liu, J. Zhang, E. Adeli, and D. Shen, Landmark-based deep multi-instance learning
802 for brain disease diagnosis, *Medical image analysis* **43**, 157–168 (2018).
- 803 30 S.-H. Wang, P. Phillips, Y. Sui, B. Liu, M. Yang, and H. Cheng, Classification of
804 Alzheimer’s disease based on eight-layer convolutional neural network with leaky rectified
805 linear unit and max pooling, *Journal of medical systems* **42**, 1–11 (2018).
- 806 31 C. Wu et al., Discrimination and conversion prediction of mild cognitive impairment
807 using convolutional neural networks, *Quantitative imaging in medicine and surgery* **8**,
808 992 (2018).
- 809 32 Y. Cao, T. A. Geddes, J. Y. H. Yang, and P. Yang, Ensemble deep learning in bioinfor-
810 matics, *Nature Machine Intelligence* **2**, 500–508 (2020).
- 811 33 A. Loddo, S. Buttau, and C. Di Ruberto, Deep learning based pipelines for Alzheimer’s
812 disease diagnosis: a comparative study and a novel deep-ensemble method, *Computers
813 in biology and medicine* **141**, 105032 (2022).
- 814 34 W. Kang et al., Multi-model and multi-slice ensemble learning architecture based on 2D
815 convolutional neural networks for Alzheimer’s disease diagnosis, *Computers in Biology
816 and Medicine* **136**, 104678 (2021).
- 817 35 J. Y. Choi and B. Lee, Combining of Multiple Deep Networks via Ensemble General-
818 ization Loss, Based on MRI Images, for Alzheimer’s Disease Classification, *IEEE Signal
819 processing letters* **27**, 206–210 (2020).
- 820 36 H. Wang, Y. Shen, S. Wang, T. Xiao, L. Deng, X. Wang, and X. Zhao, Ensemble of 3D
821 densely connected convolutional network for diagnosis of mild cognitive impairment and
822 Alzheimer’s disease, *Neurocomputing* **333**, 145–156 (2019).
- 823 37 T. G. Dietterich, Ensemble methods in machine learning, in *International workshop on
824 multiple classifier systems*, pages 1–15, Springer, 2000.
- 825 38 C. R. Jack Jr et al., The Alzheimer’s disease neuroimaging initiative (ADNI): MRI meth-
826 ods, *Journal of Magnetic Resonance Imaging: An Official Journal of the International
827 Society for Magnetic Resonance in Medicine* **27**, 685–691 (2008).
- 828 39 J. G. Sled, A. P. Zijdenbos, and A. C. Evans, A nonparametric method for automatic
829 correction of intensity nonuniformity in MRI data, *IEEE transactions on medical imaging*
830 **17**, 87–97 (1998).

- 831 40 V. Fonov, A. Evans, R. McKinstry, C. Alml, and D. Collins, Unbiased nonlinear average
832 age-appropriate brain templates from birth to adulthood, *NeuroImage* **47**, S102 (2009),
833 Organization for Human Brain Mapping 2009 Annual Meeting.
- 834 41 D. Holland et al., Subregional neuroanatomical change as a biomarker for Alzheimer's
835 disease, *Proceedings of the National Academy of Sciences* **106**, 20954–20959 (2009).
- 836 42 K. He, X. Zhang, S. Ren, and J. Sun, Deep residual learning for image recognition, in
837 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages
838 770–778, 2016.
- 839 43 J. Shi, X. Zheng, Y. Li, Q. Zhang, and S. Ying, Multimodal neuroimaging feature
840 learning with multimodal stacked deep polynomial networks for diagnosis of Alzheimer's
841 disease, *IEEE journal of biomedical and health informatics* **22**, 173–183 (2017).
- 842 44 S. Liu et al., Multimodal neuroimaging feature learning for multiclass diagnosis of
843 Alzheimer's disease, *IEEE Transactions on Biomedical Engineering* **62**, 1132–1140
844 (2014).
- 845 45 R. Cui and M. Liu, Hippocampus Analysis by Combination of 3-D DenseNet and Shapes
846 for Alzheimer's Disease Diagnosis, *IEEE journal of biomedical and health informatics*
847 **23**, 2099–2107 (2018).
- 848 46 C. Lian, M. Liu, J. Zhang, and D. Shen, Hierarchical Fully Convolutional Network for
849 Joint Atrophy Localization and Alzheimer's Disease Diagnosis Using Structural MRI,
850 *IEEE Transactions on Pattern Analysis and Machine Intelligence* **42**, 880–893 (2020).
- 851 47 Y. Chen and Y. Xia, Iterative sparse and deep learning for accurate diagnosis of
852 Alzheimer's disease, *Pattern Recognition* **116**, 107944 (2021).
- 853 48 J. Zhang, B. Zheng, A. Gao, X. Feng, D. Liang, and X. Long, A 3D densely connected
854 convolution neural network with connection-wise attention mechanism for Alzheimer's
855 disease classification, *Magnetic Resonance Imaging* **78**, 119–126 (2021).
- 856 49 C. Ju, A. Bibaut, and M. van der Laan, The relative performance of ensemble methods
857 with deep convolutional neural networks for image classification, *Journal of Applied*
858 *Statistics* **45**, 2800–2818 (2018).
- 859 50 P. K. Crane et al., Development and assessment of a composite score for memory in
860 the Alzheimer's Disease Neuroimaging Initiative (ADNI), *Brain imaging and behavior*
861 **6**, 502–516 (2012).

- 862 ⁵¹ R. C. Petersen et al., Alzheimer's disease neuroimaging initiative (ADNI): clinical char-
863 acterization, *Neurology* **74**, 201–209 (2010).
- 864 ⁵² E. Tjoa and C. Guan, A survey on explainable artificial intelligence (xai): Toward
865 medical xai, *IEEE transactions on neural networks and learning systems* **32**, 4793–4813
866 (2020).