

1 **A new strategy on Early diagnosis of cognitive impairment via novel cross-**
2 **lingual language markers: a non-invasive description and AI analysis for the**
3 **cookie theft picture**

4

5 Jintao Wang^{1,2#}, Junhui Gao^{3#}, Jinwen Xiao¹, Jianping Li¹, Haixia Li¹, Xinyi Xie¹, Rundong Tan³,

6 Yuyuan Jia³, Xinjue Zhang³, Chen Zhang³, Dake Yang³, Gang Xu⁴, Rujin Ren^{1*}, Gang Wang^{1,2*}

7 1 Department of Neurology and Institute of Neurology, Ruijin Hospital, Shanghai Jiao Tong University

8 School of Medicine, Shanghai, 200025, People's Republic of China.

9 2 Department of Neurology, Renji Hospital, Shanghai Jiao Tong University, School of Medicine,

10 China.

11 3 Shanghai Nuanhe Brain Technology Co. Ltd., Shanghai, China.

12 4 School of Public Health, Shanghai Jiao Tong University, Shanghai, 200025, China.

13 **Corresponding author:** Gang Wang, MD, PhD, Department of Neurology & Institute of Neurology,

14 Ruijin Hospital affiliated with Shanghai Jiao Tong University School of Medicine, Shanghai, 200025,

15 China (wgneuron@hotmail.com); Department of Neurology, Renji Hospital, Shanghai Jiao Tong

16 University, Shanghai School of Medicine, China.

17 Rujing Ren, MD, PhD, Department of Neurology & Institute of Neurology, Ruijin Hospital affiliated

18 with Shanghai Jiao Tong University School of Medicine, Shanghai, 200025, China

19 (docotorren2001@126.com); Department of Neurology, Renji Hospital, Shanghai Jiao Tong University,

20 Shanghai School of Medicine, China.

21 # Equal contributors

22 * Corresponding author

23 **Figures:** 3

24 **Tables:** 3

25 **References:** 31

26 **Word:** 4635

27

28

29 **Abstract**

30

31 **Background:** Cognitive impairment (CI), including Alzheimer's disease (AD) and
32 mild cognitive impairment (MCI), has been a major research focus for early diagnosis.
33 Both speech assessment and artificial intelligence (AI) have started to be applied in
34 this field, but faces challenges with limited language type assessment and ethical
35 concerns due to the "black box" nature. Here, we explore a new strategy with patient
36 led non-invasive observation for a novel cross-lingual digital language marker with
37 both diagnostic accuracy, scalability and interpretability.

38 **Methods:** Speech data was recorded from the cookie theft task in 3 cohorts. And
39 automatic speech recognition (ASR), Networkx package, jieba library and other tools
40 were used to extract visual, acoustic and language features. The SHAP model was
41 used to screen features. Logistic regression and support vector machine and other
42 methods were used to build the model, and an independent cohort was used for
43 external verification. Finally, we used AIGC technology to further reproduce the
44 entire task process.

45 **Results:** In Chinese environment, we built 3 models of NC/aMCI, NC/AD, and
46 NC/CI (aMCI+AD) through Cohort 1 (NC n=57, aMCI n=62, AD n=66), with
47 accuracy rates of 0.83, 0.79, and 0.79 respectively. The accuracy was 0.75 in the
48 external scalability verification of Cohort 3 (NC n=38, CI n=62). Finally, we built a
49 cross-lingual (Chinese and English) model through Cohort 1 and 2, built a NC/aMCI
50 diagnosis model, and the diagnostic accuracy rate was 0.76. Lastly, we successfully
51 recreate the testing process through Text-to-Image' and Animation Generation.

52 **Discussion:** The visual features created by our research group and combines acoustic
53 and linguistic features were used to build a model for early diagnosis of cognitive
54 impairment, and a cross-lingual model covering English and Chinese, which performs
55 well in external verification of independent cohorts. Finally, we innovatively used AI-
56 generated videos to show the subject's task process to the physician to assist in
57 judging the patient's diagnosis.

58

59 **Keyword:** Alzheimer's disease, Amnestic mild cognitive impairment, speech test,
60 Artificial Intelligence, interpretability

61

62 INTRODUCTION

63 Alzheimer's disease (AD) is the most common neurodegenerative disease,
64 characterized by persistent memory decline. It is reported that 55 million¹ people
65 suffer from AD and other dementia worldwide and 13 million in China². Thus,
66 dementia has been considered one of the global health dilemmas. Detecting
67 individuals who are at the early stage of dementia is essential but challenging,
68 especially in the disease modified therapy (DMT) era of monoclonal antibodies that
69 emphasizes early diagnosis and treatment. In the past few decades, major progress has
70 been made in the development of biofluid or neuroimaging biomarkers for early
71 screening and/or early diagnosis. However, these methods are limited for application
72 by being invasive and/or expensive. On the contrary, previous studies showed that
73 verbal categorical fluency test showed the highest performance in differentiating AD
74 with the heath³, and series of studies⁴⁻⁶ have demonstrated that language deficits
75 precede memory impairment. Therefore, early diagnosis through language features is
76 feasible⁷. Previous studies from our team and others suggested that cognitive
77 impairments have been effectively diagnosed through acoustic and linguistic features
78 such as percentage of silence duration (PSD)⁸⁻¹¹. These features come from the picture
79 description task "cookie-theft", which has become one of the most commonly used
80 tests of language function and was originally part of the Boston Diagnostic Aphasia
81 Examination (BDAE) manual¹². However, these methods only consider language
82 characteristics and ignore other cognitive abilities involved in describing the task.
83 Screening of cognitive function by past traits may be inadequate, and interpretability
84 is insufficient.

85 Meanwhile, advances in artificial intelligence (AI) technology have sparked new
86 research interest for easy and even remote detection, diagnosis and treatment of
87 dementia¹³. Natural language process (NLP) models have achieved relatively better
88 prediction accuracy, even higher than 90%¹⁴. However, clinical medicine relies on the
89 transparency of decision-making, and the logic of black-box models violates medical
90 ethics. Clinicians could not reasonably accept and understand the decision-making
91 process with no explainable AI (black-box models). Therefore, explainable AI has
92 become a hot topic of research in academia, industry, and government. The
93 interpretability of AI in the medical field has received widespread attention due to its
94 high-risk nature. Additionally, few studies have addressed the problem of cross-
95 language screening, and common issues between different languages are difficult to
96 discover. Most studies are limited to the detection of small samples in a single
97 language with difficult reproducibility. In the present, unfortunately, there are no
98 established and widely accepted methods so far¹⁵.

99 Therefore, we here construct a novel language-related digital model with both
100 accuracy and interpretability. Especially, the "cookie theft" task in our study was
101 participant-led, without unnecessary prompts of the physician, which can fully reflect
102 their comprehensive cognition rather than only language ability. So, firstly we create a

103 set of new features named visual features, which can describe the entities and the
104 relationship paths between the entities from the speech to reflect the task process.
105 Together with acoustic or linguistic features mentioned in the previous studies, we
106 construct models to distinguish NC (normal control) from aMCI (mild cognitive
107 impairment), NC from CI (cognitive impairment), and NC from AD. Secondly, these
108 new features can solve the cross-language issues which not handled well before¹⁶
109 regardless of language type (Mandarin or English). Finally, we use Artificial
110 Intelligence Generated Content (AIGC) to reproduce the task process, allowing
111 physicians to participate in the classification to reduce overall errors in clinical setting.

112 **METHOD**

113 **1.Participant**

114 Cohort1: This is a cross-sectional study, with a total of 185 participants recruited
115 from Ruijin Hospital Affiliated to Shanghai Jiao Tong University School of Medicine,
116 Shanghai, in which 57 were NC, 62 were aMCI, and 66 participants were diagnosed
117 with early phase AD. The registration number is ChiCTR2000036718 on the website
118 associated with this study (<https://www.chictr.org.cn>). All participants were recruited
119 between August 2020 and August 2023 from the memory clinic of Ruijin Hospital.
120 The authors asserted that all procedures contributing to this work comply with the
121 ethical standards of the relevant national and institutional committees on human
122 experimentation and with the Helsinki Declaration of 1975, as revised in 2008. All
123 procedures involving human subjects/patients were approved by the Ethics
124 Committee of the Ruijin Hospital (approval number: 2020-261). All included
125 individuals provided written consent.

126 Cohort2: In order to build cross-language models and better generalization
127 capabilities, the DementiaBank corpus was kindly used in the present study¹⁷. This
128 corpus contained recordings of 74 controls and 25 aMCI patients, from July 1983 to
129 April 1988 (last modified in November 2018) involving the participants given a
130 picture description task, which was originally designed for the Boston Diagnostic
131 Aphasia Examination. The task required each participant to describe events depicted
132 in the picture, the same as performed by participants in our center (Cookie Theft
133 picture description task).

134 Cohort3: To further verify the accuracy of the model, we randomly included 100
135 external validation cohorts from the Alzheimer's disease and other dementia clinical
136 cohorts of Ruijin Hospital (Approval number: 2022-097). These include 62 CI
137 patients(MCI due to AD and mild AD) proven by AV45 PET scans and 38 matched
138 cases with normal cognition. The speech task was performed as what has mentioned
139 above.

140 **2.Clinical assessment and diagnosis**

141 To exclude other causes of cognitive impairment, we performed cranial MRI or
142 computed tomography (CT) to exclude confounding factors such as stroke or
143 intracranial space-occupying lesions. Serum folic acid, vitamin B12 levels, and
144 thyroid function were tested to exclude endocrine and metabolic disorders. Clinical
145 and demographic data including age, gender, and level of education were also
146 collected. All subjects underwent neuropsychological tests including the Mini-Mental
147 State Examination (MMSE) , Clinical dementia rating scale(CDR) and the Cookie-
148 theft picture description task from the Boston Diagnostic Aphasia Scales.

149 After clinical assessment, the participants were categorized into three groups: (i) a NC
150 group, who were considered as cognitively healthy after the clinical consultation; (ii)
151 an AD group, whose diagnosis was based on the clinical probable criteria for
152 diagnosis of AD issued by the National Institute on Aging-Alzheimer’s Association
153 workgroups in 2011¹⁸; and (iii) an aMCI group, in which patients had a memory
154 complaint corroborated by at least one informant, and a diagnosis was conducted
155 using the Petersen criteria¹⁹. Participants were excluded if they had any other
156 neurological diseases, any systemic disease which can lead to cognitive dysfunction,
157 psychiatric disorders, or severe hearing or vision impairment. 62 of the subjects
158 accepted dual-phase [18F] AV45 PET scans, with a resolution of $3.76 \times 3.76 \times 4.9$
159 mm^3 (field of view = 157 mm). Forty-seven planes were obtained with a voxel size of
160 $1.95 \times 1.95 \times 3.2 \text{ mm}^3$. A transmission scan was performed for attenuation correction
161 before the PET acquisition. For [18F] AV45 PET, each participant underwent a 10-
162 minutes early acquisition (composed of ten 1-minute dynamical frames) that began
163 immediately after the intravenous injection of $\sim 4 \text{ MBq/kg}$ of [18F] AV45, and a 10-
164 minutes late acquisition (beginning 50-minutes after injection).

165 Subjects with $\text{MMSE} \geq 15$ was the include according to the same standards as before⁸,
166 and all enrolled patients were aMCI or mild AD patients.

167 **3.Recording protocol**

168 Subjects performed a Cookie Theft picture description task, during which they
169 were given a picture and were told to discuss everything they could see happening in
170 the picture in 1 min while being recorded. The mean time duration of the records is
171 $42.87 \pm 18.72 \text{ s}$. The RSF cohort individuals’ speech was recorded under the following
172 configuration parameters of Cool Edit Pro software: a frequency of 160000 Hz,
173 creating a 16-bit mono recording, and environmental noise was limited to under 45 dB.
174 The Pitt cohort and Cohort 3 records (the mean time duration of the records is
175 $51.1 \pm 23.63 \text{ s}$ seconds and $43.52 \pm 18.42 \text{ s}$) were converted to the audio configuration
176 parameters identical to the RSF recording using the Cool Edit Pro software.

177 **4.Information generation and processing**

178 **4.Feature Engineering**

179 **4.1 Sound Feature Extraction**

180 The automatic speech recognition (ASR) software for cognitive impairment v1.3
181 (developed by our team, China Software Copyright number 2016SR164680) for
182 speech analysis was used, according to our previous study^{8,9}. Each sample was
183 analyzed by ASR software for cognitive impairment using v1.3 to extract the
184 speech/silence parameters. The sum of all silent periods divided by the total speech
185 time is the definition of PSD (ratio of total silent pause duration to total speech
186 duration), expressed as a percentage. The definition of basic parameters set in our
187 software was according to Pakhomov et al. Who had developed the measurements of
188 spontaneous speech from the Cookie Theft picture description task for patients with
189 dementia. Silence is defined as the summed duration of all silent segments of the
190 recording, including general short pauses, general long pauses, and hesitation-
191 associated pauses.

192 **4.2 Speech-to-Text Conversion**

193 Upon obtaining the audio, after confirming its integrity, the software was utilized for
194 conversion, resulting in a transcript. The audio was then replayed to proofread the
195 transcript, with discrepancies from the original audio requiring modification. The
196 revised transcript constituted the final version.

197 **4.3 Part-of-Speech Tagging of Text**

198 For both Chinese and English texts, specific parts of speech (adjectives, adverbs,
199 prepositions, etc.) in the description texts of AD patients were statistically analyzed
200 using the Jieba library. Specifically, by utilizing the Jieba library (the best library for
201 Chinese text processing), text segmentation was performed for Chinese texts,
202 followed by part-of-speech tagging to obtain the part of speech for each word.
203 Subsequently, the total word count and the count of specific parts of speech were
204 calculated to determine the proportion of each part of speech in the text. Analyzing
205 the distribution of specific parts of speech in the given text aids in understanding the
206 linguistic characteristics of the text.

207 **4.4 Visual Feature Extraction**

208 **Construction of Spatial Semantic Graphs**

209 Initially, the two-dimensional centroid coordinates (x_i , y_i) of entities were calculated,
210 and spatial semantic graphs were constructed for each participant. For each participant,
211 the text was scanned from start to finish to extract entities, obtaining an entity list. A
212 directed graph was then constructed based on the entity list and their coordinates.

213 **Feature Extraction of Spatial Semantic Graphs**

214 Networkx (v2.8.4 <https://github.com/networkx/networkx>) were utilized to extract
215 features from graphs. Representative strength features include: the number of edges,
216 path length, number of left-right switches, and diameter of the gaze area. Features
217 representing efficiency include entity density. Features representing attention include:
218 the number of loops, proportion of left-right descriptions in the graph, and graph
219 density.

220 **5. Model Construction**

221 **5.1 Feature selection as the digital markers**

222 SHAP (SHapley Additive exPlanations v0.44.0)²⁰ was used to explain the predictions
223 of machine learning (ML) models based on the Shapley value concept from
224 cooperative game theory. Shapley values measure each participant's contribution to
225 the outcome of the game, and in ML, they are used to quantify each feature's impact
226 on the model's output.

227 At the same time, based on the assumption that language impairment will worsen with
228 cognitive impairment, only features with consistent trends among NC, aMCI, and AD
229 groups will be included as digital markers and ranked according to Shapley values.
230 Specifically, Random Forest Classifier was chosen as the ML model and fitted on the
231 training set. Predictions were made on the test set, and the model's accuracy
232 (accuracy_score) was computed. Then, an explainer was created using the SHAP
233 library to calculate and visualize Shapley values, explaining the model's contribution
234 to the predictions. Finally, force_plot was used to display Shapley values for the
235 certain sample.

236 **5.2 Construction of Single-Language Models**

237 We performed the construction of three models for NC (n=57) and aMCI (n=62), NC
238 (n=57) and AD (n=66), and NC (n=57) and CI (AD+aMCI) (n=128), respectively.
239 The data were partitioned into training and testing sets using the train_test_split
240 function in Sklearn package (v1.0.2)²¹, with 80% allocated for training and 20% for
241 testing to evaluate model performance. Feature filtering and selection included
242 accurate evaluation of features' consistency with the disease, balance among three
243 factors (expression intensity, expression efficiency, attention), etc. High-credibility
244 models recognized in clinical research were selected: logistic regression (LR), support
245 vector machine (SVM), random forest (RF), k-nearest neighbors (KNN), etc.

246 **5.3 Construction of Cross-Language Models**

247 This study presents an analysis of classification performance on a dataset comprising
248 218 samples of mixed-language texts from both normal individuals (NC) and those
249 with amnesiac mild cognitive impairment (aMCI). The dataset encompasses 8 lexical
250 diversity features, 9 visual features, and 6 pause-related features. The data were
251 partitioned into training and testing sets using the train_test_split function, with 80%
252 allocated for training and 20% for testing to evaluate model performance. A logistic
253 regression model was initialized and fitted to the training set. The fitted model was
254 then used to predict on the testing set, yielding probability values. These probabilities
255 were utilized to compute parameters of the Receiver Operating Characteristic (ROC)
256 curve, including True Positive Rate (TPR) and False Positive Rate (FPR) for model
257 evaluation. In the graphical representation, the orange curve represents the ROC curve,
258 with the Area Under the Curve (AUC) serving as one of the metrics for evaluating
259 model performance. Higher AUC values signify superior model performance.

260

261 **6. Process Reproduction**

262 **6.1 Image Generation**

263 Our method is based on the open-source project stable-diffusion-webui v4.7²², which
264 implements Stable Diffusion, a deep learning model for image generation. We
265 introduced the Control Net plugin (v11p-sd15)²³ into this method, which extends the

266 functionality of Stable Diffusion to better control the layout and content of images.
267 Specifically, we used Text-Guided and Image-Guided functions to redraw the
268 generated images. The Text-Guided function allows us to use textual descriptions to
269 guide the image generation process, such as specifying the theme or content direction
270 of the image. The Image-Guided function allows us to use reference images to guide
271 image generation, ensuring consistency in content and style between the generated
272 image and the reference image. In summary, our method combines Stable Diffusion
273 algorithm with the ControlNet plugin, as well as Text-Guided and Image-Guided
274 functions, to achieve fine control and redraw of the image generation process, thereby
275 generating high-quality images with specific content and layout. Specific to our
276 language description, after inputting the content of the text conversion, we can obtain
277 high-quality restored pictures.

278 **6.2 Animation Generation**

279 Same as the Image Generation process, iterating through the text describing the
280 images by participants, starting with a blank image, a new image was generated if the
281 number of entities increased. Finally, all images were concatenated in the order they
282 were generated to create the animation. Here, the duration of each image corresponds
283 to the time in the original audio

284 **7 Statistical analysis**

285 The demographic information analysis was performed using SPSS (Version 26.0). T test and one-
286 way ANOVA were used for the group differences, Turkey test was used for post-hoc analysis.
287 Chi-square test and Fisher test were used to detect the frequency differences between groups. The
288 Pearson correlation was used for the association between features and MMSE subscore. $P < 0.05$
289 were considered significant

290

291 **RESULTS**

292 **1. Demographic information**

293 Cohort1:

294 There were 57 NC, 62 aMCI, and 66 AD patients in the cohort1. Gender (female in
295 NC: 59.65%, aMCI: 58.06%, and AD: 48.48 %) and educational level (NC: $14.98 \pm$
296 3.15 years, aMCI: 14.76 ± 3.18 years, and AD: 14.68 ± 3.51 years) showed no
297 significant difference among the NC, aMCI, and AD groups in the cohort, and mean
298 age was 69.60 ± 7.71 , 72.82 ± 7.43 , and 73.21 ± 8.63 years for the NC, aMCI,
299 and AD groups within this cohort ($P= 0.0269$), respectively. However, there were
300 significant differences between groups' mean MMSE scores (NC: 29.11 ± 0.99 ,
301 aMCI: 25.64 ± 4.31 , and AD: 18.85 ± 3.38), and subitems relating to each cognitive

302 domain (Table 1, all $P < 0.0001$), and the post-hoc comparison results are shown in
303 Table 1.

304 Cohort2:

305 There were 74 NC, 25 aMCI patients in the Pitt cohorts. Gender (female makeup in
306 NC: 63.51%, and aMCI : 44.00%) and educational level (NC: 17.19 ± 0.99 years,
307 and aMCI : 17.00 ± 2.16 years), and mean age (63.84 ± 8.29 , and 67.04 ± 8.76
308 years for the NC, and aMCI groups) showed no significant difference among the
309 NC, and aMCI group ($P > 0.05$). There were significant differences between groups'
310 mean MMSE scores (NC: 29.09 ± 1.05 , and aMCI : 26.96 ± 2.39) ($P = 0.0002$).

311 Cohort3:

312 There were 38 NC, 68 CI patients in the Alzheimer's disease and other dementia
313 clinical cohorts. Gender (female makeup in NC: 42.11%, and CI : 53.23%) and
314 educational level (NC: 11.42 ± 3.58 years, and CI : 10.58 ± 4.63 years) , and
315 mean age (70.74 ± 4.92 , and 69.44 ± 10.13 years for the NC, and CI) showed no
316 significant difference among the NC, and CI groups ($P > 0.05$). There were
317 significant differences between groups' mean MMSE scores (NC: 28.55 ± 2.06 , and
318 CI : 21.25 ± 6.30) ($P < 0.0001$).

319 **Table 1**

320

321 **2. Feature selection as the digital markers**

322 We ultimately selected 23 features as the digital markers for cognitive impairment
323 related to the subjects' voice, language, and vision. Their IDs and descriptions are
324 shown in the table below.

325

326 **Table 2**

327 We randomly selected the SHAP values for three samples Figure 1. According to the
328 plots, the contribution value of each feature influences the model's prediction for a
329 specific sample. The plot shows how the model's prediction for a sample is composed
330 of the impacts of individual features. In these samples, pause features (Pause_1-3) and
331 the proportions of specific parts of speech (mainly nouns and prepositions) play a
332 significant role in the model's predictions and how they influence the model's
333 decisions. By calculating the average SHAP values across the three samples, we
334 determine the importance of each feature. Pause_1, Text_n, and Text_p were the most
335 important features in the model. Together with other features, they were included as
336 the digital markers.

337 **Table 3**

338

339 **3. Employing Machine Learning Methods Accepted by Clinicians**

340

341 A ML model was constructed to perform a binary classification between NC/aMCI,
342 NC/AD, NC/CI (aMCI+AD). The ROC curves comparing PSD based classification
343 sensitivity and specificity among NC, aMCI, and AD patients are shown in Fig.1A–D.
344 The AUCs of the curves are 0.83, 0.79, and 0.79 in NC/aMCI, NC/AD, and NC/CI
345 (aMCI+AD). Further, the sensitivity and specificity of NC/aMCI, NC/AD, and NC/CI
346 (aMCI+AD) is 0.71/0.71, 0.84/0.70, and 0.78/0.79 respectively. The results are shown
347 in the following figure. And the weight of each features were listed in the
348 Supplementary Table 2.

349

350

Figure1

351

4. Constructing Cross-Lingual Models for Further Interpretability Support

353

354 The visual features are related to the described order during the examination, and are
355 independent of the language, wording, and sentence structure. We use visual features
356 to solve cross-language diagnostic problems losslessly. In order to clarify the suitable
357 features of the model, all the features need to simultaneously satisfied the
358 'consistency' of change in both Chinese and English, which the mean value of this
359 indicator satisfies a monotonic change across the three groups: $NC > aMCI > AD$, or
360 $NC < aMCI < AD$. Ultimately, the cross-lingual model achieved an accuracy of 0.76
361 and a sensitivity of 0.75.

362

Figure 2

363

5. External validation of machine learning models

365 To evaluate the effectiveness of our model in distinguishing NC from CI, we used an
366 independent external cohort, and all CI patients in this cohort underwent AV45-PET
367 examination and were confirmed to be $A\beta$ positive. After external validation, the
368 model's prediction accuracy reached 75%, with sensitivity and specificity of 68.24%
369 and 73.33%.

370

6. association of visual feature and cognitive domain

372 From the single language model, the visual_0 and visual_4 has the biggest weight for
373 the model, we found they are significantly associated with the MMSE attention and
374 delayed memory subscore(visual_0 with memory: $R^2=0.03891$, $P=0.0492$, visual_4
375 with memory: $R^2=0.1451$, $P<0.0001$, visual_4 with attention: $R^2=0.09499$, $P=0.0018$).
376 Supplementary Table 1 and **Figure 3**.

377

7. Recreating the Testing Process Through Text-to-Image' and Animation Generation

379

380

381 The first method for process recreation is "Text-to-Image" to gain the richer content
382 and more cartoon-like features image compared to the other AIGC for example
383 Image-to-Image method. However, the "Text-to-Image" method cannot show the
384 patient's task description at each time point, and may miss some information. Most
385 importantly, this method is not reproducible and cannot be compared across multiple
386 samples. Therefore, we further used animation generation to completely display the
387 process. The original picture was segmented, and was divided into different entities,
388 which will present according to the task process. The generated animation can inform
389 the doctor which entities were described, which were not described, and the duration
390 of each entity's description.

391

392

Figure 3

393 DISCUSSION

394 In the present study, we successfully constructed a strategy with both accuracy and
395 interpretability, and innovatively created new digital language markers(Figure4). Here,
396 the cookie-theft task was participant-led, and the physician's role was minimized,
397 providing only necessary prompts to fully reflect their cognitive integration ability.
398 Therefore, we draw on the entity path diagram (EPD) of graph theory to create visual
399 features that could reflect information including language but not limited to language.
400 In addition, with visual features we constructed a cross-language cognitive screening
401 model successfully. Finally, this study innovatively used AIGC to more intuitively
402 cooperate with clinicians in clinical applications through generative images or videos.
403 Regarding to good cost performance and easily handling compared to traditional body
404 fluids and neuroimaging biomarkers, we believe new digital markers have better
405 accessibility advantages. It may be particularly suitable for early screening in
406 multilevel referral systems for cognitive impairment, and particularly in less well-
407 resourced or remote regions.

408

Figure 4

409 **1 Models via SHAP method with better accuracy and interpretability reduce** 410 **black box effect confusion**

411 This study adheres to the principle of interpretability in feature selection, and selects
412 interpretable features related to speech, text, and vision through SHAP²⁰ to build a
413 classification model. This approach follows the rigor of medical research and ensures
414 that the features used can be used as digital or biomarkers. For our cognitive
415 screening in the Chinese language environment, we constructed classification
416 diagnosis models with accuracy of 0.83, 0.79, and 0.79 in NC/aMCI, NC/AD, and
417 NC/aMCI+AD through ML models such as SVM. We obtained a good prediction
418 accuracy, compared with other previous studies^{11,14,25-27}. But unlike those technologies
419 based on large language models, including word2vec, Bert, and GPT, the use of

420 SHAP and SVM makes our model more interpretable rather than an unknown black
421 box process. To further verify the reliability of our model, we introduced an
422 independent cohort 3 for external validation in which all patients had abnormal β
423 amyloid deposits through AV45 PET scans. In cohort 3, our model also achieved a
424 prediction accuracy of 75%, indicating that our model has high external scalability
425 and high value in the diagnosis of cognitive impairment. We further analyzed the
426 weight of each feature in the model of NC/aMCI+AD. Among them, the acoustic
427 features headed by ratio of hesitation/phrasal counts, and PSD discovered by our
428 previous research. The Scale among the visual features also have a greater
429 contribution in this model. We believe the pauses in the speech can well reflect the
430 patient's cognitive ability^{8,9}. A larger PSD indicates worse cognitive function. Scale
431 can refer to the hierarchy or complexity of the graph. In our model, the smaller the
432 complexity, the more likely the patient to be cognitively impaired. People with
433 cognitive impairment describe the task more simply. In addition, we hold that visual
434 features can reflect more cognitive domains than speech ability. We found that visual
435 features with high weights in the model such as visual_0 and visual_4 are closely
436 related to MMSE subscores, especially to cognitive domains such as memory and
437 attention, which further proves the reliability of visual features in cognitive screening.

438

439 **2 Visual features solve cross-language problems well**

440 At present, the biggest barriers in the methods for the speech detection of cognitive
441 impairment is the fact that: most models are language dependent. For example, Yan
442 and colleagues relied on semantic features and existence of transcription tools from
443 any language to English and/or powerful NLP models²⁸. Most of these methods
444 convert samples of different languages into the same language, or only use acoustic
445 features for model construction¹⁴. Translation may lose real information, and the black
446 box properties of the NLP model are also limited. The visual features are less
447 dependent on language type and more related to task completion. Therefore, with the
448 same speech collection conditions, we merge the English speech samples from the Pitt
449 database¹⁷ and the Chinese speech data in our cohort to truly classify the cognitive
450 levels of different language environments in the same model. We constructed a cross-
451 language model with an higher accuracy of 0.76 compared with 0.70 in Fraser and
452 colleagues' model¹⁶ in classifying a mixed sample of Swedish and English. Compared
453 with in single language environment, the model accuracy was similar which means
454 the important role of visual features in these model.

455 **3 Image and Animation Generation**

456 In picture description tasks, many studies have implemented relatively automated
457 screening processes, but often a high degree of automation may lead to potential
458 errors that are difficult to detect by clinicians. Therefore, we innovatively made use of

459 AIGC technology combined with visual features to generate images/videos that can
460 reproduce the completion process of language tasks. The images/videos will facilitate
461 more intuitive inspection by clinicians. Information such as the completeness of the
462 content or the intensity of the description, the order of the description, and the spacing
463 will be fully displayed in the picture. In our three patient examples of NC, aMCI, and
464 AD(Figure 3), different patients presented different completion results and different
465 processes, which very intuitively shows the completion status of patients with
466 different cognitive or language abilities.

467 **Figure3**

468 **4 Limitation**

469 First, although, we have considered confounding factors such as education level, age,
470 and gender in the model construction, however, language as a high-level brain
471 function is still affected by factors such as religion, emotion, personality traits, and
472 language habits^{6,15}. In the future, these factors should be taken into consideration.
473 Secondly, How to avoid the learning effect after multiple tests and reduce the
474 deviation is the direction we are working for. Like other cross-sectional studies^{10,29,30},
475 the conclusions of this study do not have the highest evidence value, and we need
476 more efforts to construct longitudinal cohorts to obtain more reliable conclusions.
477 Thirdly, the etiology of MCI is a heterogeneous³¹. Thus, it is envisioned as future
478 work the implementation of multilingual or language independent systems, supported
479 by extensive and diverse databases (that still must be gathered, with genders, ages,
480 disease severity), as well as the automation of the features selection and extraction.
481 Better decision models, task oriented, are also required.

482 Generally, the present study successfully creates new digital tool from a new
483 perspective, and uses digital markers and AI to further improve the ability to diagnose
484 early cognitive impairment across languages. Meanwhile, as a vision-related
485 parameter, it can also reflect advanced cognitive functions such as attention and
486 observation. Therefore, our results suggest that visual features can be used not only to
487 screen for cognitive disorders such as AD, but also for diseases related to cognitive
488 changes such as attention deficit and hyperactivity disorder, depression. With the
489 rapid rise of AI, its application has immeasurable prospects. Making full use of new
490 digital markers to diagnose patients with early cognitive impairment will be simpler
491 and more efficient than traditional methods in future.

492 **ACKNOWLEDGEMENT**

493 Sincere thanks to each subject participating in this clinical trial.

494 **FOOTNOTES**

495 Jintao Wang, Junhui Gao contributed equally.

496 **CONTRIBUTORS:** Conceptualisation—Gang Wang, Jintao Wang, Junhui Gao.
497 Methodology—Junhui Gao, Jintao Wang, Rundong Tan, Yuyuan Jia, Xinjue Zhang,
498 Chen Zhang, Dake Yang, Gang Xu, Rujin Ren, and Gang Wang. Investigation—
499 Jintao Wang, Jinwen Xiao, Jianping Li, Haixia Li, Xinyi Xie, and Gang Wang.
500 Visualisation—Jintao Wang, Junhui Gao, Rundong Tan, Yuyuan Jia, Xinjue Zhang,
501 Chen Zhang, Dake Yang, and Gang Wang. Funding acquisition—Gang Wang. Project
502 administration—Gang Wang. Supervision—Gang Wang. Writing (original draft)—
503 Gang Wang, Jintao Wang and Junhui Gao. Writing (review and editing)—Gang
504 Wang, Gang Xu and Rujin Ren. Guarantor—Gang Wang.

505 **FUNDING:** This work was supported by the Ministry of Science and Technology of
506 the People's Republic of China (2021ZD0201804, GW).

507 **COMPETING INTERESTS:** Junhui Gao, Rundong Tan, Yuyuan Jia, Xinjue Zhang,
508 Chen Zhang, Dake Yang are current employees of Shanghai Nuanhe Brain
509 Technology Co. Ltd., Shanghai, China. Other authors declare that they have no
510 competing interests.

511 **SUPPLEMENTAL MATERIAL:**

512 **REFERENCE**

- 513 1. Estimation of the global prevalence of dementia in 2019 and forecasted prevalence in 2050: an analysis for
514 the Global Burden of Disease Study 2019. *The Lancet Public Health* 2022; **7**(2): e105-e25.
- 515 2. Ren R, Qi J, Lin S, et al. The China Alzheimer Report 2022. *Gen Psychiatr* 2022; **35**(1): e100751.
- 516 3. Battista P, Salvatore C, Berlinger M, Cerasa A, Castiglioni I. Artificial intelligence and neuropsychological
517 measures: The case of Alzheimer's disease. *Neuroscience and Biobehavioral Reviews* 2020; **114**: 211-28.
- 518 4. Bäckman L, Jones S, Berger A-K, Laukka EJ, Small BJ. Cognitive impairment in preclinical Alzheimer's
519 disease: a meta-analysis. *Neuropsychology* 2005; **19**(4): 520-31.
- 520 5. Meilan JJG, Martinez-Sanchez F, Carro J, Carcavilla N, Ivanova O. Voice Markers of Lexical Access in
521 Mild Cognitive Impairment and Alzheimer's Disease. *Curr Alzheimer Res* 2018; **15**(2): 111-9.
- 522 6. Folia V, Liampas I, Siokas V, et al. Language performance as a prognostic factor for developing Alzheimer's
523 clinical syndrome and mild cognitive impairment: Results from the population-based HELIAD cohort. *J Int*
524 *Neuropsychol Soc* 2023; **29**(5): 450-8.
- 525 7. Olmos-Villaseñor R, Sepulveda-Silva C, Julio-Ramos T, et al. Phonological and Semantic Fluency in
526 Alzheimer's Disease: A Systematic Review and Meta-Analysis. *Journal of Alzheimer's Disease : JAD* 2023; **95**(1).
- 527 8. Wang H-L, Tang R, Ren R-J, et al. Speech silence character as a diagnostic biomarker of early cognitive
528 decline and its functional mechanism: a multicenter cross-sectional cohort study. *BMC Med* 2022; **20**(1): 380.
- 529 9. Qiao Y, Xie XY, Lin GZ, et al. Computer-Assisted Speech Analysis in Mild Cognitive Impairment and
530 Alzheimer's Disease: A Pilot Study from Shanghai, China. *J Alzheimers Dis* 2020; **75**(1): 211-21.
- 531 10. Mueller KD, Kosciak RL, Hermann BP, Johnson SC, Turkstra LS. Declines in Connected Language Are
532 Associated with Very Early Mild Cognitive Impairment: Results from the Wisconsin Registry for Alzheimer's

- 533 Prevention. *Front Aging Neurosci* 2017; **9**: 437.
- 534 11. König A, Satt A, Sorin A, et al. Automatic speech analysis for the assessment of patients with predementia
535 and Alzheimer's disease. *Alzheimer's & Dementia (Amsterdam, Netherlands)* 2015; **1**(1): 112-24.
- 536 12. Roth C. Boston Diagnostic Aphasia Examination. In: Kreutzer JS, DeLuca J, Caplan B, eds. *Encyclopedia of*
537 *Clinical Neuropsychology*. New York, NY: Springer New York; 2011: 428-30.
- 538 13. Parsapoor M. AI-based assessments of speech and language impairments in dementia. *Alzheimers Dement*
539 2023.
- 540 14. Amini S, Hao B, Zhang L, et al. Automated detection of mild cognitive impairment and dementia from voice
541 recordings: A natural language processing approach. *Alzheimers Dement* 2022.
- 542 15. Vigo I, Coelho L, Reis S. Speech- and Language-Based Classification of Alzheimer's Disease: A Systematic
543 Review. *Bioengineering (Basel)* 2022; **9**(1).
- 544 16. Fraser KC, Lundholm Fors K, Eckerström M, Öhman F, Kokkinakis D. Predicting aMCI Status From
545 Multimodal Language Data Using Cascaded Classifiers. *Front Aging Neurosci* 2019; **11**: 205.
- 546 17. Becker JT, Boller F, Lopez OL, Saxton J, McGonigle KL. The natural history of Alzheimer's disease.
547 Description of study cohort and accuracy of diagnosis. *Archives of Neurology* 1994; **51**(6): 585-94.
- 548 18. McKhann GM, Knopman DS, Chertkow H, et al. The diagnosis of dementia due to Alzheimer's disease:
549 recommendations from the National Institute on Aging-Alzheimer's Association workgroups on diagnostic
550 guidelines for Alzheimer's disease. *Alzheimer's & Dementia : the Journal of the Alzheimer's Association* 2011;
551 **7**(3): 263-9.
- 552 19. Petersen RC, Smith GE, Waring SC, Ivnik RJ, Tangalos EG, Kokmen E. Mild cognitive impairment: clinical
553 characterization and outcome. *Archives of Neurology* 1999; **56**(3): 303-8.
- 554 20. Lundberg SM, Lee S-I. A unified approach to interpreting model predictions. Proceedings of the 31st
555 International Conference on Neural Information Processing Systems. Long Beach, California, USA: Curran
556 Associates Inc.; 2017. p. 4768–77.
- 557 21. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: Machine Learning in Python. *J Mach Learn Res*
558 2011; **12**(null): 2825–30.
- 559 22. Rombach R, Blattmann A, Lorenz D, Esser P, Ommer B. High-Resolution Image Synthesis with Latent
560 Diffusion Models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) 2021*:
561 10674-85.
- 562 23. Zhang L, Rao A, Agrawala M. Adding Conditional Control to Text-to-Image Diffusion Models. *2023*
563 *IEEE/CVF International Conference on Computer Vision (ICCV) 2023*: 3813-24.
- 564 24. Reuben D, Levin J, Frank J, et al. Closing the dementia care gap: Can referral to Alzheimer's Association
565 chapters help? *Alzheimer's & Dementia : the Journal of the Alzheimer's Association* 2009; **5**(6): 498-502.
- 566 25. He R, Chapin K, Al-Tamimi J, et al. Automated Classification of Cognitive Decline and Probable
567 Alzheimer's Dementia Across Multiple Speech and Language Domains. *Am J Speech Lang Pathol* 2023: 1-12.
- 568 26. Wang R, Kuang C, Guo C, et al. Automatic Detection of Putative Mild Cognitive Impairment from Speech
569 Acoustic Features in Mandarin-Speaking Elders. *J Alzheimers Dis* 2023.
- 570 27. Agbavor F, Liang H. Predicting dementia from spontaneous speech using large language models. *PLOS*
571 *Digit Health* 2022; **1**(12): e0000168.
- 572 28. Conneau A, Khandelwal K, Goyal N, et al. Unsupervised cross-lingual representation learning at scale. *arXiv*
573 *preprint arXiv:191102116* 2019.
- 574 29. Mahajan P, Baths V. Acoustic and Language Based Deep Learning Approaches for Alzheimer's Dementia
575 Detection From Spontaneous Speech. *Front Aging Neurosci* 2021; **13**: 623607.
- 576 30. Martínez-Nicolás I, Llorente TE, Martínez-Sánchez F, Meilán JGG. Ten Years of Research on Automatic

577 Voice and Speech Analysis of People With Alzheimer's Disease and Mild Cognitive Impairment: A Systematic
578 Review Article. *Front Psychol* 2021; **12**: 620251.
579 31. Gauthier S, Reisberg B, Zaudig M, et al. Mild cognitive impairment. *Lancet (London, England)* 2006;
580 **367**(9518): 1262-70.
581

582 **Table 1 The demographic feature of subjects**

	cohort 1 (n=185)			statis tic value	P	coh ort2 (n=99)		stati stic valu e	P	cohort 3 (n=100)		stati stic valu e	P
	NC (n=57)	aMCI (n=62)	AD (n=66)			NC (n=74)	aMCI (n=25)			NC (n=38)	CI (n=62)		
num													
Age (years)	69.60 ± 7.71	72.82 ± 7.43	73.21 ± 8.63	F=3.690	0.0269	63.84 ± 8.29	67.04 ± 8.76	t=-1.646	0.1029	70.74 ± 4.92	69.44 ± 10.13	t' = 0.859	0.3925
Gender, female, n, %	34 (59.65%)	36 (58.06%)	32 (48.48%)	$\chi^2=1.865$	0.3936	47 (63.51%)	11 (44.00%)	$\chi^2=2.933$	0.088	16 (42.11%)	33 (53.23%)	$\chi^2=1.166$	0.2802
edu catio n (years)	14.98 ± 3.15	14.76 ± 3.18	14.68 ± 3.51	F=0.135	0.8738	17.19 ± 0.99	17.00 ± 2.16	t' = 0.423	0.675	11.42 ± 3.58	10.58 ± 4.63	t=0.957	0.3409
Apo E e4 type , n, %	11/54, 20.37%	20/57, 35.09%	38/66, 57.58%	$\chi^2=1.782$	0.0001	-	-	-	-	-	-	-	-
MM SE	29.11 ± 0.99	25.23 ± 4.31	18.85 ± 3.38	F=15.694	<0.0001	29.09 ± 1.05	26.96 ± 2.39	t' = 4.329	0.002	28.55 ± 2.06	21.25 ± 6.30	t' = 7.253	<0.0001

583

584 Table 2 The IDs and Descriptions of the Three Categories of Features

ID (Sound)	ID	Meaning
Sound (aroustic feature)	Pause_1	Percentage of silence duration, %
	Pause_2	Ratio of silence/speech counts, %
	Pause_3	Ratio of hesitation/speech counts, %
	Pause_4	Ratio of long pause/speech counts, %
	Pause_5	Ratio of short pause/speech counts, %
	Pause_6	Ratio of hesitation/phrasal counts, %
language	Text_p	Preposition
	Text_d	Adverb
	Text_r	Pronoun
	Text_c	Conjunction
	Text_v	Verb
	Text_a	Adjective
visual	Visual_0 (scale)	The number of edges in the entity sequence The diameter of the graph. The length of the longest shortest path between all pairs of entities (node pairs)
	Visual_1 (diameter)	
	Visual_2 (Total path length)	The sum of the lengths of the edges that form the entity sequence
	Visual_3 (The center of the graph)	The number of central nodes. The center refers to the set of nodes whose eccentricity is equal to the radius. Defined as: $2m/n(n-1)$, where m is the number of edges and n is the number of nodes.
	Visual_4 (Density 1)	
	Visual_5 (Density 2)	Scale/diameter. (Visual_0/Visual_5)
	Visual_6 (ring)	The number of rings, for example, if the entity sequence described by the subject includes ' , girl, mother, plate, girl, ', a ring appears.
	Visual_7 (Left and right switching)	The number of times the entity sequence described by the subject switches between the left and right sides of the graph The ratio of the number of times the subject focuses on the number of entities on the left to the number of times the subject focuses on the number of entities on the right in the entity sequence described. If an entity is described twice, the count is 2.
Visual_8 (Left and right gaze ratio)		

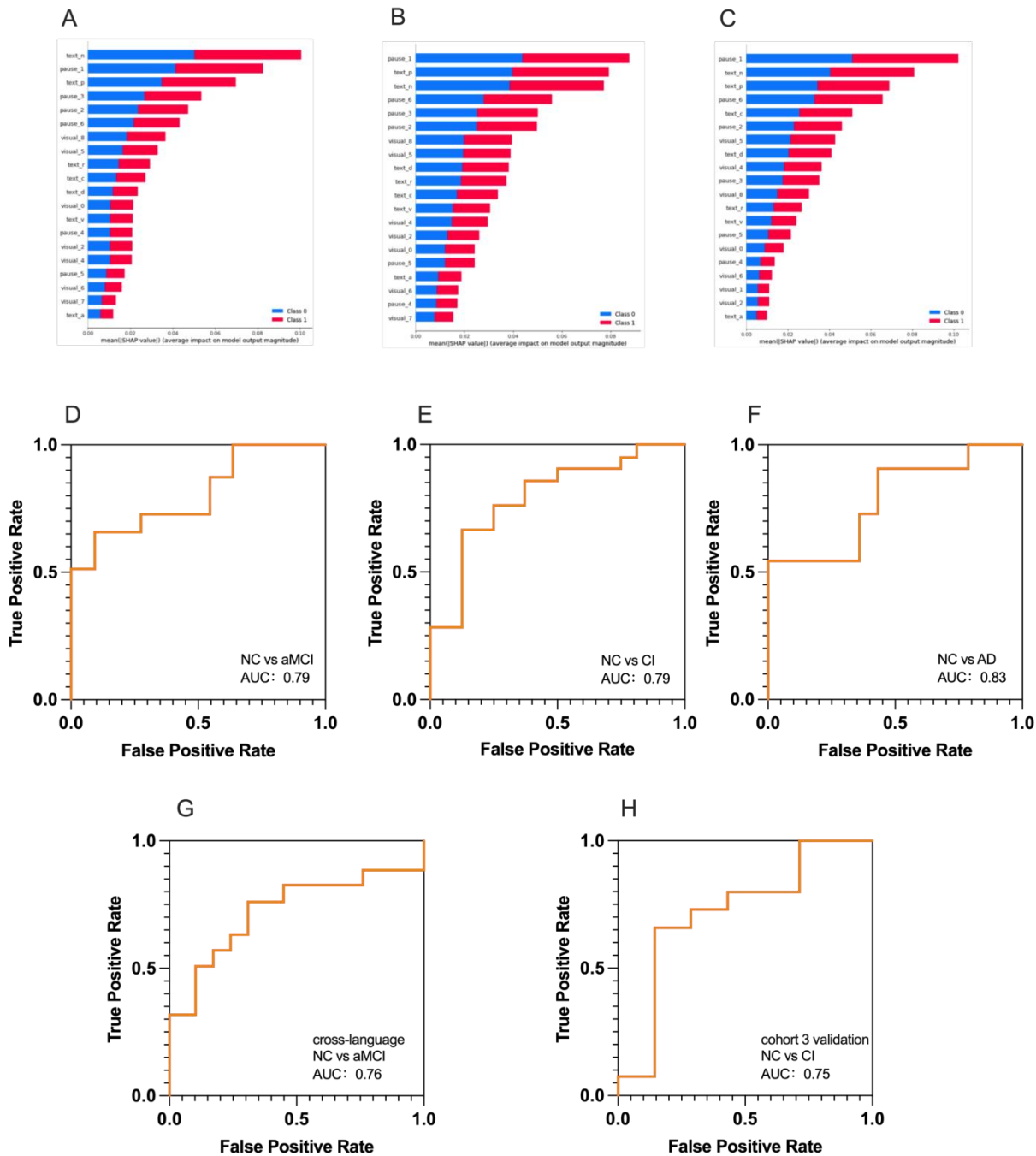
585

586 Table 3 Feature Importance Ranking

	Sample1	Sample2	Sample3	Average
Pause_1 (PSD)	2	1	1	1.33
Text_n (Noun)	1	3	2	2.00
Text_p (Preposition)	3	2	3	2.67
Pause_6 (hesitation counts)	6	4	4	4.67
Pause_2 (RSD)	5	6	6	5.67
Pause_3 (hesitation)	4	5	10	6.33
Visual_5 (Density 2)	8	8	7	7.67
Visual_8 (Left and right gaze ratio)	7	7	11	8.33
Text_d (Adverb)	9	9	8	8.67

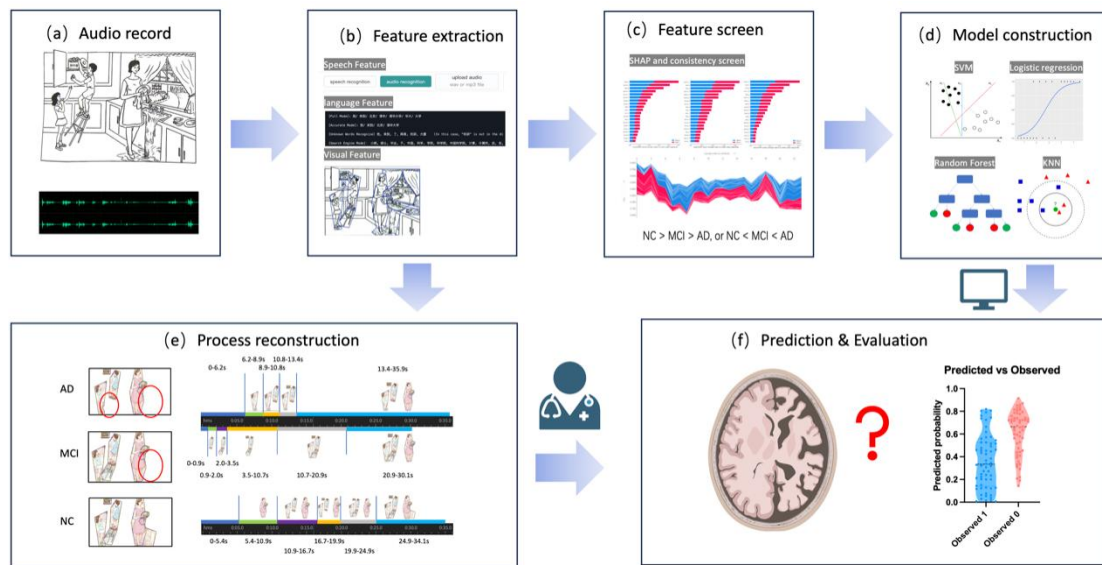
587

588 Figure1
 589 $R^2=0.1451$, $P<0.0001$



590
 591 Figure1 (A-C) The contribution value of different features to prediction of randomly
 592 selected three samples, (D-F) Diagnostic model for NC/aMCI, NC/AD, NC/CI. And
 593 the area under the curve. (G) In a cross-language environment, the diagnostic ROC
 594 curve of diagnostic model for NC/aMCI, and the area under the curve. (G) Validation
 595 of the NC/CI diagnostic model in an external cohort 3 and the area under the curve.
 596

597 Figure 2



598

599

600 Figure2 Language Cognition Screening and Reproduction Flowchart

601

602

603

604

605

606

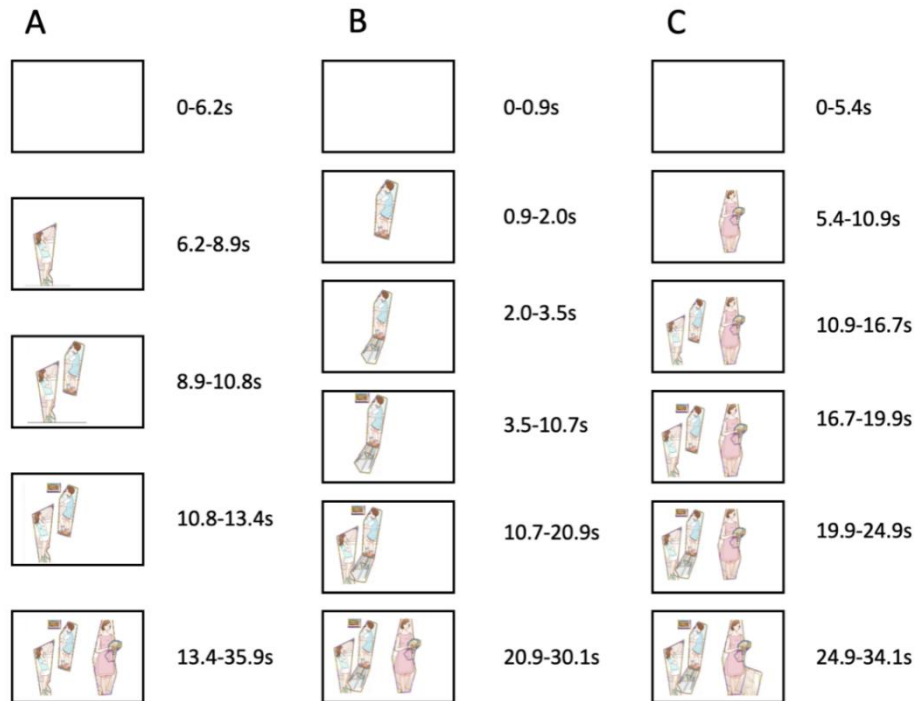
607

608

609

610

611

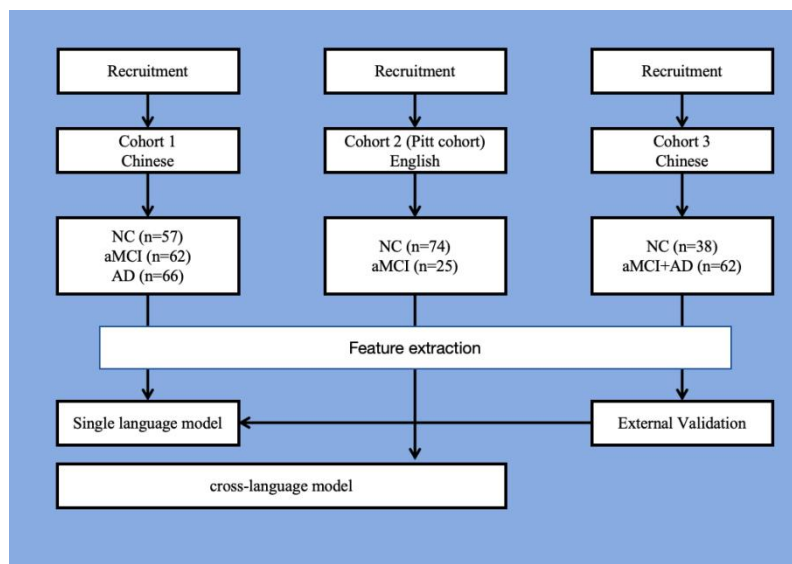


612

613 Figure3 Generated animation screenshot of different timelines. (A-C) the animation of

614 AD, aMCI, NC respectively

615



616 Figure 4 Study design of the clinical trial

617

618

619