

# **Developing an Accessible Dementia Assessment Tool: Leveraging a Residual Network, the Trail Making Test, and Demographic Data**

**Authors:** Jingmei Yang<sup>1</sup>, Samad Amini<sup>1</sup>, Boran Hao<sup>1</sup>, Seho Park<sup>2</sup>, Cody Karjadi<sup>3</sup>, Lance San Souci<sup>2</sup>, Vijaya B. Kolachalama<sup>2,4,5</sup>, Stephanie Cosentino<sup>6</sup>, Stacy L. Andersen<sup>2</sup>, Rhoda Au<sup>2,3,7</sup>, and Ioannis Ch. Paschalidis<sup>1,4,\*</sup>

<sup>1</sup> Department of Electrical & Computer Engineering, Division of Systems Engineering, and Department of Biomedical Engineering, Boston University

<sup>2</sup> Department of Medicine, Chobanian & Avedisian School of Medicine, Boston University

<sup>3</sup> Framingham Heart Study, Boston University

<sup>4</sup> Faculty of Computing & Data Sciences, Boston University

<sup>5</sup> Department of Computer Science, Boston University

<sup>6</sup> Department of Neurology, Columbia University Irving Medical Center

<sup>7</sup> Departments of Anatomy & Neurobiology, Neurology, and Epidemiology, Chobanian & Avedisian School of Medicine and School of Public Health, Boston University

## **Corresponding Author:**

Ioannis Ch. Paschalidis

Rafik B. Hariri Institute for Computing and Computational Science & Engineering,

Boston University, 665 Commonwealth Ave., Boston, MA 02215, USA

e-mail: [yannisp@bu.edu](mailto:yannisp@bu.edu), Tel: [617-353-0434](tel:617-353-0434), <http://sites.bu.edu/paschalidis>

1 **Abstract**

2 **Background:** The global burden of Alzheimer’s disease and related dementias is rapidly  
3 increasing, particularly in low- and middle-income countries where access to specialized  
4 healthcare is limited. Neuropsychological tests are essential diagnostic tools, but their  
5 administration requires trained professionals, creating screening barriers. Automated  
6 computational assessment presents a cost-effective solution for global dementia screening.

7 **Objective:** To develop and validate an artificial intelligence-based screening tool using the Trail  
8 Making Test (TMT), demographic information, completion times, and drawing analysis for  
9 enhanced dementia detection.

10 **Methods:** We developed: (1) non-image models using demographics and TMT completion times,  
11 (2) image-only models, and (3) fusion models. Models were trained and validated on data from  
12 the Framingham Heart Study (FHS) ( $N = 1,252$ ), the Long Life Family Study (LLFS) ( $N = 1,613$ ),  
13 and the combined cohort ( $N = 2,865$ ).

14 **Results:** Our models, integrating TMT drawings, demographics, and completion times, excelled  
15 in distinguishing dementia from normal cognition. In the LLFS cohort, we achieved an Area Under  
16 the Receiver Operating Characteristic Curve (AUC) of 98.62%, with sensitivity/specificity of  
17 87.69%/98.26%. In the FHS cohort, we obtained an AUC of 96.51%, with sensitivity/specificity of  
18 85.00%/96.75%.

19 **Conclusions:** Our method demonstrated superior performance compared to traditional approaches  
20 using age and TMT completion time. Adding images captures subtler nuances from the TMT  
21 drawing that traditional methods miss. Integrating the TMT drawing into cognitive assessments

- 22 enables effective dementia screening. Future studies could aim to expand data collection to include
- 23 more diverse cohorts, particularly from less-resourced regions.
- 24 **Keywords:** Alzheimer's disease, Artificial Intelligence, Dementia, Trail Making Test

## 25 **Introduction**

26 Dementia, of which *Alzheimer's Disease (AD)* is the most common form, profoundly impacts  
27 memory, thinking, and daily functioning. The global prevalence of dementia is increasing rapidly  
28 and is projected to reach 139 million cases by 2050.<sup>1</sup> Currently, about 60% of dementia cases are  
29 in low- and middle-income countries, a figure projected to rise to 71% by 2050.<sup>2</sup> Dementia not  
30 only impacts individuals but also imposes substantial financial and emotional strains on families  
31 and societies, with costs exceeding 1.3 trillion dollars annually and anticipated to rise to 2.8 trillion  
32 dollars by 2030.<sup>3</sup> Dementia, with its substantial global economic impact, often remains  
33 undiagnosed, with only 20-50% of cases identified in high-income countries and fewer in low-  
34 income regions.<sup>4</sup> The growing global population of older adults intensifies the need for large-scale  
35 screening to effectively manage, monitor, and even predict this age-related disease.<sup>5,6</sup>

36 To address the need for large-scale dementia screening, pen-and-paper drawing tasks within  
37 neuropsychological test batteries have become essential. These include the *Clock Drawing Test*  
38 (*CDT*),<sup>7</sup> *Trail Making Test (TMT)*, *Pentagon Drawing Test (PDT)*,<sup>8</sup> and *Rey-Osterrieth Complex*  
39 *Figure Test (RCFT)*, which are extensively used in detecting neurocognitive disorders such as  
40 *Alzheimer's Disease and Related Dementias (ADRD)*. Although these drawing tests can be easily  
41 implemented in a pen-and-paper format, they still require administration and interpretation by  
42 trained professionals. Recent research has utilized digital pen technology to collect drawing tests  
43 electronically. Leveraging digital data, studies have used deep learning techniques to develop  
44 automated scoring systems, demonstrating promising performance in the PDT,<sup>9-12</sup> RCFT,<sup>13,14</sup> and  
45 CDT.<sup>15,16</sup>

46 However, research on developing an automated diagnostic tool using the TMT has been relatively

47 limited. The TMT consists of two parts: (i) in Part A (TMT-A), numbered circles are displayed  
48 on a piece of paper, and participants are instructed to use a pen to draw lines to connect the numbers  
49 in sequential order as quickly as possible; (ii) in TMT Part B (TMT-B), a series of numbers and  
50 letters are displayed and participants are instructed to connect the numbers and letters in alternate  
51 sequences.<sup>17</sup> Over the past decades, research has shown that the TMT is a useful tool for detecting  
52 cognitive impairment, with completion time being a strong indicator. This has led to the  
53 development and validation of adapted TMT versions across many countries.<sup>18-20</sup> A major area of  
54 research on TMT involves establishing normative data for older adult populations, investigating  
55 differences based on education level, gender, and age across culturally and geographically diverse  
56 countries.<sup>21-23</sup>

57 Additionally, digital adaptations of the TMT, known as dTMT, have introduced more nuanced  
58 measurements. Studies comparing these digital iterations with the traditional pen-and-paper format  
59 affirm their efficacy in cognitive assessment.<sup>24</sup> The dTMT's key advantage is its capability to  
60 gather additional data, offering a richer, more comprehensive analysis than its original  
61 counterpart.<sup>25-27</sup> While there are studies demonstrating the utility of TMT in discriminating  
62 between cognitive impairment and normal cognition utilizing geographically diverse datasets, a  
63 significant portion of prior studies have relied on data from a single center or limited geographical  
64 area with small sample sizes. Moreover, the majority of these studies utilized only the completion  
65 time of the TMT as a predictor variable, with few incorporating the actual drawings into their  
66 models. As a result, the predictive potential of TMT drawings remains largely unexplored. To fully  
67 explore the diagnostic capabilities of the TMT for dementia, large-scale, multi-center studies are  
68 needed that analyze both the TMT completion time and the drawings, in conjunction with  
69 demographic information.

70 To accelerate the screening and diagnosis of cognitive impairment globally, our study aimed to  
71 develop an accessible detection tool that can be readily implemented worldwide, including in  
72 resource-limited regions. Accordingly, our model utilized easily collectible information such as  
73 age, gender, education level, the completion time of TMT-A, and the TMT-A drawing.  
74 Specifically, TMT-A was selected over TMT-B given its use of universally recognized Arabic  
75 numerals, as opposed to alphabet letters used in TMT-B which may not be as widely recognized,  
76 especially in non-English speaking countries. Using data from the *Framingham Heart Study (FHS)*  
77 and *Long Life Family Study (LLFS)*, our multifaceted study first explored the ability of  
78 demographics and completion time, individually and in combination, to distinguish individuals at  
79 risk for developing dementia. We then fine-tuned two vision networks to evaluate the predictive  
80 potential of TMT-A drawings. Ultimately, we developed a fusion model that integrated three key  
81 components: (1) the probability score output by the fine-tuned vision model when using a TMT-  
82 A drawing as input, (2) demographic characteristics, and (3) the TMT-A completion time. By  
83 combining these signals in our fusion model, we significantly enhanced the model's overall  
84 performance both within and across studies.

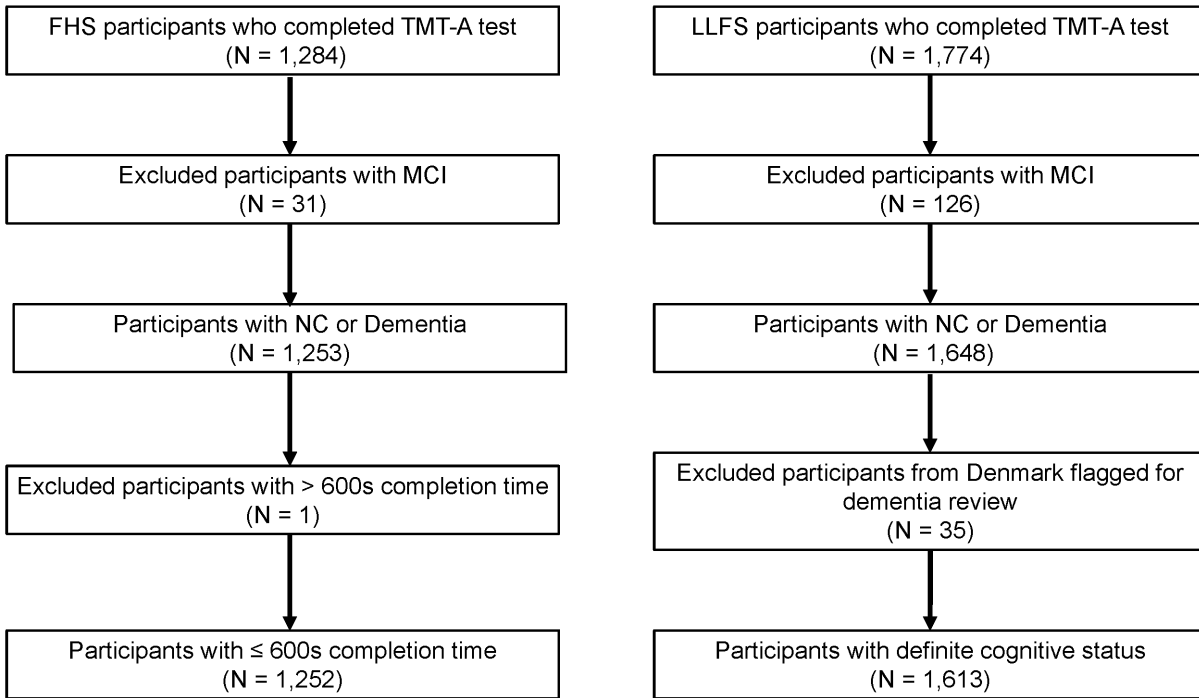
## 85 **Methods**

### 86 *Study cohorts*

87 We used data collected from 1,252 participants in the FHS cohort and 1,613 participants in the  
88 LLFS cohort. FHS is the longest ongoing longitudinal transgenerational cohort study of chronic  
89 disease.<sup>28</sup> LLFS is a multicenter (Boston, New York, Pittsburgh, and Denmark) longitudinal study  
90 of human longevity and healthy aging.<sup>29</sup> Data used in our study were collected across several  
91 geographic regions which are different in environment, culture, and demographics. All  
92 participants provided written informed consent. Study protocols and consent forms were approved  
93 by the Boston University Medical Campus Institutional Review Board and the Institutional Review  
94 Boards of the LLFS field sites as well as the LLFS coordinating center at Washington University  
95 St. Louis.

96 The cognitive status was provided by each study cohort. In the FHS cohort, a participant's  
97 cognitive status was determined by the dementia diagnostic review panel.<sup>30</sup> A dementia diagnosis  
98 for those showing signs of cognitive decline was reached by consensus between at least one  
99 neurologist and one neuropsychologist, based on neurology exams, medical records, and brain  
100 imaging.<sup>31</sup> In the FHS dataset, there were 1,233 participants with normal cognition (NC) and 19  
101 with dementia. Similar to the FHS cohort, an impaired cognitive status in the LLFS cohort was  
102 determined by a dementia diagnostic review panel based on cognitive testing and informant  
103 interviews. Specifically, participants from Denmark who were flagged for a dementia review were  
104 omitted due to the absence of case adjudication in Denmark, with only those classified as normal  
105 being included. Given the gradual progression of cognitive decline and the limited available  
106 samples of Mild Cognitive Impairment (MCI) cases, we focused our analysis on distinguishing

107 between normal cognition and dementia cases. Consequently, 1,548 participants were identified  
 108 with normal cognition, and 65 participants were diagnosed with dementia. The participant selection  
 109 process for both the FHS and LLFS datasets is illustrated in Figure 1.



110 **Figure 1:** Flow diagram for the participant selection process in the Framingham Heart Study (FHS) and the Long Life  
 111 Family Study (LLFS). Participants with mild cognitive impairment (MCI), those with greater than 600s completion  
 112 time on the Trail Making Test Part A, and those without a definite cognitive status were excluded.

113 *Data preparation*

114 During the in-person TMT, information such as gender, age, education, and type of Apolipoprotein  
 115 E (ApoE) alleles was documented. A digital pen recorded  $x$  and  $y$  coordinates approximately  
 116 every 13 milliseconds. The data corresponding to each pen stroke, which is defined as continuous  
 117 drawing without lifting the pen, was saved in a separate section of a *.txt* file.

118 We deliberately chose not to use the raw data collected by the digital pen, as digital pens are not

119 readily available, and older adults are often unfamiliar with such technology. Our goal was to  
120 develop an automated tool that does not heavily depend on resources or expertise predominantly  
121 available in developed countries and regions. While digital pens capture rich information including  
122 pressure sensitivity, precise temporal dynamics, and stroke-level metrics, these advantages come  
123 at the cost of accessibility and widespread applicability. There were additional reasons for not  
124 relying on the digital pen detailed trajectory data. Specifically, LLFS test administrators recorded  
125 participant names during the TMT administration. By converting to a static image format, we  
126 could easily anonymize the data through targeted image cropping, while preserving the essential  
127 elements of the TMT drawing itself. Further, an image-based approach allows our method to  
128 be applied to TMT drawings collected through various means, including those administered with  
129 traditional paper-and-pencil methods and subsequently digitized. This universality enhances the  
130 potential for retrospective analyses of existing datasets and enables wider implementation across  
131 diverse clinical environments without requiring specialized equipment.

132 To diminish our reliance on digital pens, we developed a preprocessing pipeline. This pipeline  
133 was designed to extract  $x$  and  $y$  coordinates, compute the duration of each stroke, and derive the  
134 overall completion time. We then plotted these extracted coordinates from a participant on a blank  
135 canvas to create an image and stored it as a *.png* file. This approach preserves the essential spatial  
136 and visuomotor patterns evident in TMT performance while allowing for standardized processing  
137 across diverse collection methods.

138 Our datasets are highly imbalanced, with a significantly smaller proportion of dementia cases  
139 compared to normal cognition cases. To address this imbalance and mitigate potential overfitting,  
140 we implemented a comprehensive set of strategies. Following our previous work, we generated  
141 additional training data through image augmentation, applying transformations more frequently to

142 the minority class.<sup>15</sup> These transformations include rotations of  $\pm 10$  degrees, zooming in/out by  
143  $\pm 15\%$ , width and height shifts of  $\pm 10\%$ , shearing of  $\pm 10\%$ , and image resizing to  $224 \times 224$  pixels.  
144 We applied these transformations disproportionately, based on the ratio of positive to negative  
145 cases. This approach enabled us to expand our training dataset with more varied positive cases.  
146 For non-image and fusion models, we adopted oversampling in the training sets, where samples  
147 from the minority class were randomly duplicated to match the majority class distribution.  
148 Additionally, we implemented stratified 5-fold cross-validation to ensure that each fold maintained  
149 the same proportion of dementia and normal cognition cases as the overall dataset, preventing  
150 potential sampling biases during model evaluation. To further optimize model performance on  
151 imbalanced data, we refined the classification threshold based on F1 scores for each fold, typically  
152 resulting in thresholds greater than the standard 0.5. This approach helped achieve a more  
153 appropriate balance between sensitivity and specificity in our predictions. We also applied  
154 regularization techniques to prevent overfitting. Specifically, we used weight decay ( $\ell_2$ -norm  
155 penalty with weight  $\lambda = 0.001$ ) during model training, which penalizes large weights and  
156 encourages the model to learn more generalizable patterns. Detailed data preprocessing steps and  
157 imputation procedures are provided in Appendix A.

### 158 *Statistical analysis and performance metrics*

159 We performed the Kolmogorov-Smirnov test for continuous variables and the  $\chi^2$  test for categorical  
160 variables to determine whether the distributions across the *Normal Cognition (NC)* and *Dementia*  
161 groups are significantly different.<sup>32,33</sup> A significance level of 0.05 was used. A  $p$ -value less than  
162 0.05 indicates that the distribution of a given feature is significantly different across cognitive  
163 statuses.

164 The data were randomly divided using stratified five-fold cross-validation. A model was trained  
165 on four folds and tested on the remaining fold. This training process was repeated five times.  
166 Performance metrics for all models were reported as the mean across the five runs, along with the  
167 standard deviation. The performance metrics included Area Under the Receiver Operating  
168 Characteristic Curve (AUC), sensitivity, specificity, weighted F1 score, and accuracy.

## 169 *Models*

170 As demonstrated in Figure 2, we developed three types of binary classification models aimed at  
171 identifying participants with dementia from those with normal cognition: (i) image-only models,  
172 (ii) non-image models, and (iii) fusion models.

173 The TMT-A drawings were used as inputs for our image-only models and two backbones were  
174 employed: *Residual Network (ResNet)* and *Vision Transformer (ViT)*.<sup>34,35</sup> Specifically, we selected  
175 the *ResNet-50* variant and the *ViT* base variant. ResNet-50 is a 50-layer convolutional neural  
176 network, which includes 48 convolutional layers, one *MaxPool* layer, and one average pooling  
177 layer. Conversely, the ViT base variant has 12 transformer layers, a hidden size of 768 dimensions,  
178 and 12 attention heads; it is designed to process 196 patches of  $16 \times 16$  pixels each. In our study,  
179 the ResNet-50 and ViT base variants are referred to as ResNet and ViT, respectively. Given the  
180 relatively small sample size, we applied transfer learning and initialized the selected backbones  
181 with weights pre-trained on ImageNet.<sup>36</sup> We fine-tuned these two backbones by appending a fully  
182 connected layer to each and training only that layer while freezing all other layers. Notably, this  
183 approach significantly reduces the complexity of training; only 1,538 trainable parameters in ViT  
184 and 4,098 in ResNet need to be trained. In the model training process, a batch size of 32 and a  
185 total of 50 epochs were used to train each vision model. Images were resized to  $224 \times 224$  pixels

186 and were normalized by dividing the value of each pixel by 255, thus rescaling the pixel values to  
187 a range between 0 and 1. The Adam optimizer was utilized for updating model parameters, with  
188 a learning rate of  $3 \times 10^{-4}$  and a weight decay of  $1 \times 10^{-3}$ . A dynamic adjustment mechanism  
189 for the learning rate was implemented, reducing the learning rate by a factor of 0.2 if no  
190 improvement in validation loss was observed for 5 consecutive epochs. The best model state was  
191 saved based on the lowest validation loss observed over the training course.

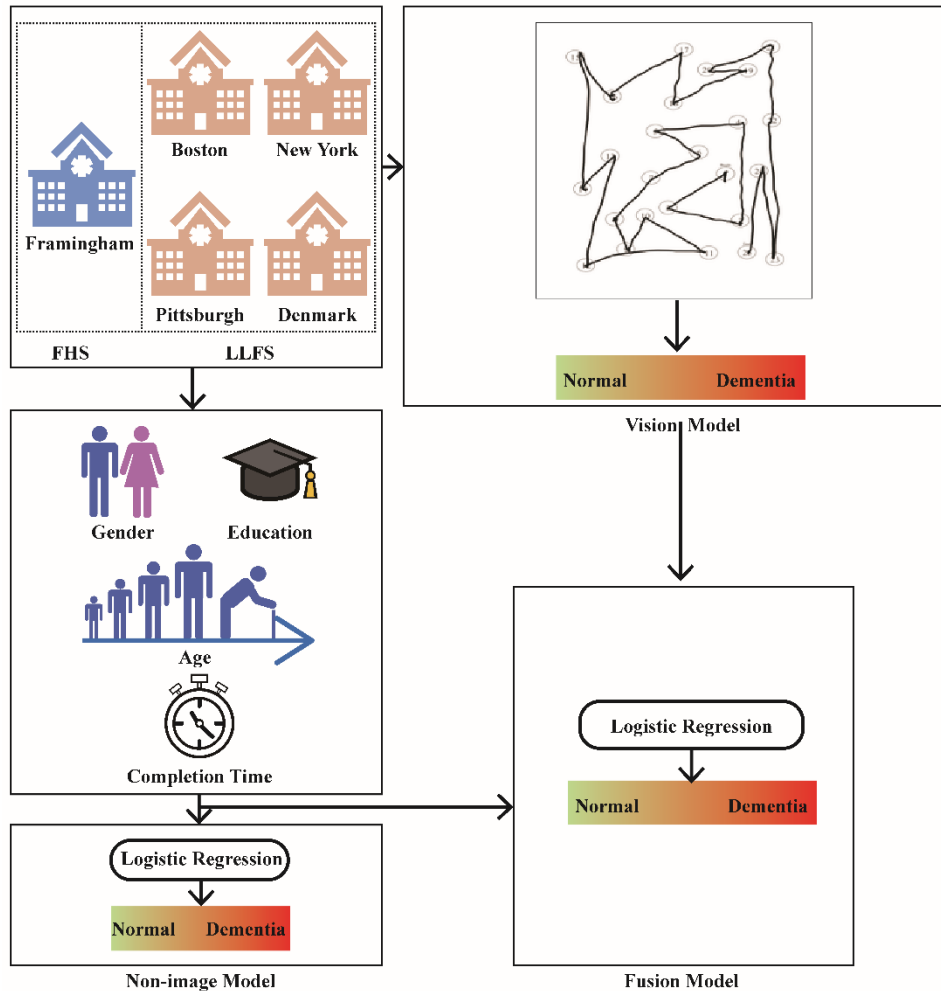
192 We used factors such as age, gender, education, and completion time to develop our non-image  
193 models. Although ApoE is known for its significant predictive value in dementia, we intentionally  
194 excluded it from our model development since it isn't routinely assessed in dementia evaluations.  
195 We developed a wide range of traditional machine learning classifiers, including Logistic  
196 Regression (LR), Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), and  
197 Extreme Gradient Boosting (XGBoost).<sup>37</sup> These classifiers were trained solely on demographic  
198 features (referred to as baseline), specifically age, gender, and education, to determine the best  
199 model for subsequent analysis. As the LR classifier showed the highest overall performance, we  
200 selected the LR algorithm for our fusion model (cf. Appendix B). To understand how non-imaging  
201 features affected our model's performance, we used the following sets of features as inputs to  
202 logistic regression models for comparison: (1) age alone, (2) education alone, (3) gender alone, (4)  
203 completion time alone, (5) baseline (age, education, and gender), (6) age and completion time,  
204 and (7) baseline and completion time.

205 The fusion model integrates multiple features through a logistic regression model. Specifically, we  
206 combine: (1) the probability score derived from the fine-tuned vision model's analysis of TMT-A  
207 drawings, (2) demographic information, and (3) completion time. All these features serve as inputs  
208 to the logistic regression model, which then predicts the final probability of dementia. Figure 3

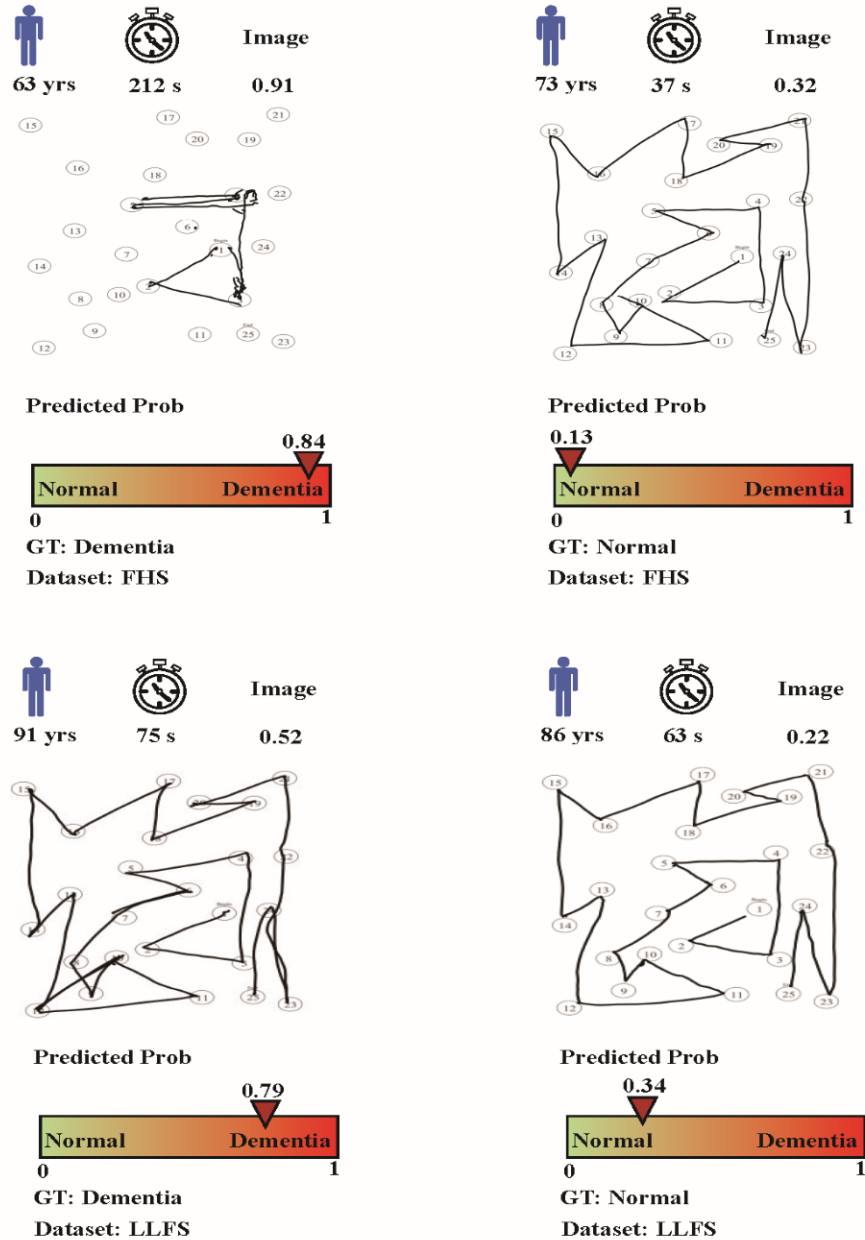
209 illustrates how the best-performing fusion model combines age, completion time, and the ResNet-  
210 50 derived image probability to generate predictions, alongside participants' ground truth  
211 cognitive status.

212

213



214 **Figure 2:** Dementia detection framework. This study leveraged data from two cohorts: the Framingham Heart Study  
 215 (FHS) cohort ( $N = 1,252$ , single-center in Framingham, MA) and the Long Life Family Study (LLFS) cohort ( $N =$   
 216  $1,613$ , multi-center in Boston, New York, Pittsburgh, and Denmark). Our study used two forms of data collected  
 217 during an in-person neuropsychological test: demographic characteristics (age, gender, and education) and digital pen  
 218 data from the Trail Making Test Part A (TMT-A), including the TMT-A completion time and  $x$  and  $y$  coordinates.  
 219 These coordinates were plotted and saved as *.png* images. We developed three types of dementia detection models:  
 220 (i) image-only models, where the *Vision Transformer (ViT)* and *Residual Network (ResNet)* were fine-tuned using  
 221 TMT-A drawings; (ii) non-image models, utilizing demographic data and the TMT-A completion time; and (iii) fusion  
 222 models, which combine non-image features with the image scores derived from the fine-tuned vision models to  
 223 differentiate individuals with dementia from their normal counterparts.



224 **Figure 3:** Visualization of the proposed fusion model. The figure includes four subfigures, each representing a  
 225 participant from either the FHS or LLFS cohort. Two of the participants are diagnosed with dementia, while the other  
 226 two are cognitively normal. For each participant, their age, the TMT-A completion time, and the probability of  
 227 dementia derived from the ResNet-50 model based on the TMT-A image are displayed. The bar below each subfigure  
 228 shows the predicted probability of dementia when combining age, the TMT-A completion time, and the image-derived  
 229 probability as predictors in the proposed fusion model. The Ground Truth (GT) label indicates the actual cognitive  
 230 status of each participant.

231 **Results**

232 *Characteristics of the participants*

233 Descriptive statistics based on cognitive status (NC and dementia) are presented in Table 1 and  
234 Table 2. The feature distributions for the FHS and LLFS datasets are provided in Appendix C. In  
235 both datasets, several differences emerged when comparing the characteristics of participants in  
236 the dementia group with those in the normal cognition group. Participants flagged as demented  
237 were generally older. They took longer to complete the TMT-A. The dementia group also had a  
238 higher percentage of females than males. Additionally, a greater proportion of participants in the  
239 dementia group had an education of high school level or lower compared to those in the normal  
240 group.

241 When comparing the overall characteristics across datasets, distinct differences become evident.  
242 Specifically, participants from LLFS are older than those from FHS on average. Additionally,  
243 LLFS has a higher proportion of individuals with an education level of high school or lower.  
244 Typically, individuals from LLFS take a longer time to complete the test.

245 **Table 1:** Descriptive characteristics and statistical analysis in the FHS dataset. This table compares  
 246 both numerical and categorical features between the Normal Cognitive (NC) and dementia groups  
 247 within the FHS cohort. For Age and Completion time, this table shows mean values for the NC,  
 248 dementia, and overall groups. Kolmogorov-Smirnov tests were conducted to assess the  
 249 significance of differences in these features between NC and dementia groups. For Gender,  
 250 Education, and Apolipoprotein E (ApoE), column percentages are presented for each group. Chi-  
 251 square ( $\chi^2$ ) tests were conducted to determine the significance of differences in these features  
 252 between NC and dementia groups. A significance level of  $\alpha = 0.05$  was used, with a  $p$ -value less  
 253 than 0.05 indicating a significant difference in feature distribution by cognitive status.

Feature	NC	Dementia	Overall	$p$ -value
	N = 1233 (98.48%)	N = 19 (1.52%)	N = 1252 (100.00%)	
Age	69.74 $\pm$ 12.31	84.79 $\pm$ 8.71	70.25 $\pm$ 12.42	< 0.01
Gender				0.37
Female	713 (57.83%)	13 (68.42%)	741 (57.76%)	
Male	520 (42.17%)	6 (31.58%)	542 (42.24%)	
Education				0.24
High school or lower	228 (18.49%)	5 (26.32%)	242 (18.86%)	
College or above	1005 (81.51%)	14 (73.68%)	1041 (81.14%)	
Completion time	37.56 $\pm$ 21.89	105.44 $\pm$ 72.9	39.14 $\pm$ 25.20	< 0.01
ApoE				0.09
$\epsilon 2/\epsilon 2$	7 (0.57%)	0 (0.00%)	7 (0.55%)	
$\epsilon 2/\epsilon 3$	148 (12.00%)	5 (26.32%)	157 (12.24%)	
$\epsilon 2/\epsilon 4$	23 (1.87%)	0 (0.00%)	23 (1.79%)	

---

$\varepsilon 3 / \varepsilon 3$	812 (65.86%)	7 (36.84%)	840 (65.47%)
$\varepsilon 3 / \varepsilon 4$	226 (18.33%)	7 (36.84%)	237 (18.47%)
$\varepsilon 4 / \varepsilon 4$	17 (1.38%)	0 (0.00%)	19 (1.48%)

---

254

255

256 **Table 2:** Descriptive characteristics and statistical analysis in the LLFS dataset. This table  
 257 compares both numerical and categorical features between the Normal Cognitive (NC) and  
 258 dementia groups within the LLFS cohort. For Age and Completion time, this table shows mean  
 259 values for the NC, dementia, and overall groups. Kolmogorov-Smirnov tests were conducted to  
 260 assess the significance of differences in these features between NC and dementia groups. For  
 261 Gender, Education, and Apolipoprotein E (ApoE), column percentages are presented for each  
 262 group. Chi-square ( $\chi^2$ ) tests were conducted to determine the significance of differences in these  
 263 features between NC and dementia groups. A significance level of  $\alpha = 0.05$  was used, with a p-  
 264 value less than 0.05 indicating a significant difference in feature distribution by cognitive status.

<b>Feature</b>	<b>NC</b>	<b>Dementia</b>	<b>Overall</b>	<b>p-value</b>
	N = 1548 (95.97%)	N = 65 (4.03%)	N = 1613 (100.00%)	
Age	69.67 ± 8.68	94.97 ± 7.38	71.74 ± 10.74	< 0.01
Gender				0.64
Female	841 (54.33%)	35 (53.85%)	939 (54.00%)	
Male	707 (45.67%)	30 (46.15%)	800 (46.00%)	
Education				0.02
High school or lower	479 (30.94%)	29 (44.62%)	556 (31.97%)	
College or above	1069 (69.06%)	36 (55.38%)	1183 (68.03%)	
Completion time	38.68 ± 17.57	135.05 ± 91.24	44.36 ± 33.05	< 0.01
ApoE				0.21
$\epsilon_2/\epsilon_2$	9 (0.58%)	0 (0.00%)	9 (0.52%)	
$\epsilon_2/\epsilon_3$	214 (13.82%)	5 (7.69%)	241 (13.86%)	
$\epsilon_2/\epsilon_4$	30 (1.94%)	1 (1.54%)	37 (2.13%)	

---

$\varepsilon_3/\varepsilon_3$	999 (64.53%)	50 (76.92%)	1129 (64.92%)
$\varepsilon_3/\varepsilon_4$	274 (17.7%)	9 (13.85%)	299 (17.19%)
$\varepsilon_4/\varepsilon_4$	22 (1.42%)	0 (0.00%)	24 (1.38%)

---

265

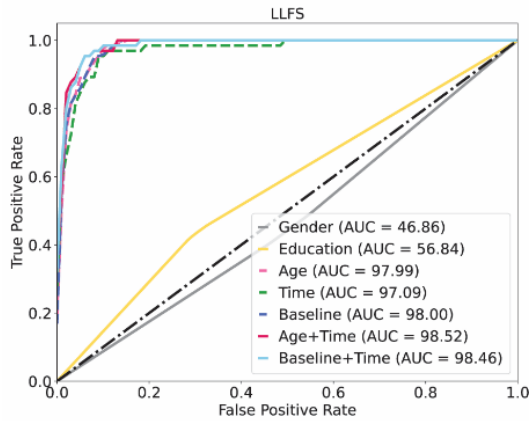
266 *Non-image models*

267 Table 3 presents the mean performance metrics, across five runs, of the non-image models  
268 specifically designed for distinguishing participants with dementia from those with normal  
269 cognition. Both gender and education, when evaluated individually, showed poor predictive  
270 power. Their AUC values were just around 50%, indicating that their performance was nearly  
271 equivalent to random guessing. In contrast, age and completion time were strong individual  
272 predictors. The higher AUC scores for age and completion time suggest that these features  
273 considerably contributed to the model's ability to discern individuals with dementia. When we  
274 integrated the completion time with demographics, we observed increases in both AUC and  
275 sensitivity across all datasets. Since completion time is the time taken to complete the TMT-A,  
276 this finding suggests that effective screening for potential dementia requires the inclusion of some  
277 form of cognitive testing. The risk of developing dementia cannot be inferred from demographics  
278 alone.

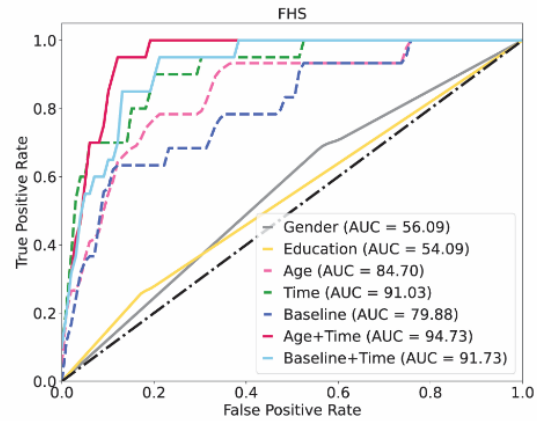
279 The best-performing non-image models are the models that used both age and completion time as  
280 features. Specifically, compared to the metrics of the baseline model, this combination enhanced  
281 the AUC by 14.96% and sensitivity by 18.33% in the FHS dataset. For the LLFS dataset, this  
282 combination achieved an AUC of 98.50%, a sensitivity of 81.54%, and a specificity of 98.64%.  
283 Like the results for the individual datasets, the top-performing non-image model for the combined  
284 dataset utilized both age and completion time as predictors. As shown in Figure 4, mixing the two  
285 datasets led to decreases in AUC values compared to using the LLFS dataset alone, but resulted in  
286 increases in AUC values for the FHS dataset.

287 **Table 3:** Performance metrics (in %) of non-image models using demographics and completion  
288 time. The models were trained and evaluated on the LLFS, FHS, and combined datasets,  
289 respectively, to differentiate participants with dementia from those with normal cognition. This  
290 table summarizes the mean values  $\pm$  standard deviations in percentages for a 5-fold cross-  
291 validation process. Evaluation metrics included AUC, sensitivity, specificity, F1 score, and  
292 accuracy.

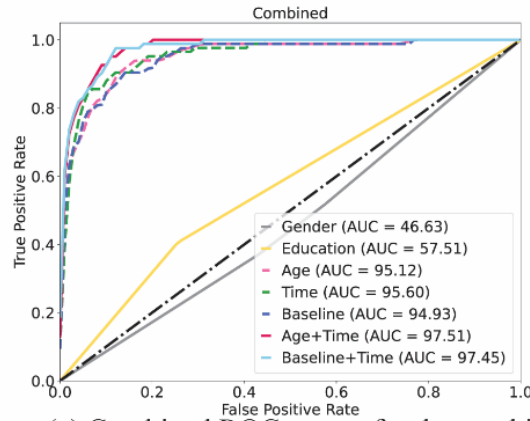
Dataset	Feature	AUC	Sensitivity	Specificity	F1 score	Accuracy
LLFS	Gender	46.86 $\pm$ 2.07	43.07 $\pm$ 4.21	50.65 $\pm$ 4.86	63.68 $\pm$ 4.10	50.34 $\pm$ 4.56
	Education	56.84 $\pm$ 6.37	44.61 $\pm$ 11.41	69.06 $\pm$ 1.75	77.74 $\pm$ 1.43	68.07 $\pm$ 2.04
	Age	97.98 $\pm$ 0.78	78.46 $\pm$ 10.03	97.81 $\pm$ 0.66	97.21 $\pm$ 0.59	97.02 $\pm$ 0.67
	Time	97.09 $\pm$ 2.01	73.85 $\pm$ 11.67	98.19 $\pm$ 1.49	97.35 $\pm$ 0.93	97.21 $\pm$ 1.18
	Baseline	97.97 $\pm$ 0.69	78.46 $\pm$ 13.76	98.06 $\pm$ 0.68	97.40 $\pm$ 0.23	97.27 $\pm$ 0.26
	Age+Time	<b>98.50 <math>\pm</math> 0.59</b>	<b>81.54 <math>\pm</math> 15.95</b>	98.64 $\pm$ 0.48	98.00 $\pm$ 0.42	97.95 $\pm$ 0.35
	Baseline+Time	98.46 $\pm$ 0.68	81.54 $\pm$ 15.95	98.77 $\pm$ 0.84	98.11 $\pm$ 0.42	98.08 $\pm$ 0.46
FHS	Gender	56.09 $\pm$ 21.09	70.00 $\pm$ 41.08	42.17 $\pm$ 1.50	58.27 $\pm$ 1.73	42.57 $\pm$ 1.98
	Education	54.09 $\pm$ 13.07	26.67 $\pm$ 25.28	89.22 $\pm$ 9.87	92.15 $\pm$ 4.99	88.26 $\pm$ 9.29
	Age	84.72 $\pm$ 10.68	46.67 $\pm$ 31.51	92.55 $\pm$ 14.64	94.20 $\pm$ 8.70	91.87 $\pm$ 13.96
	Time	91.00 $\pm$ 6.11	46.67 $\pm$ 19.19	98.70 $\pm$ 1.12	98.08 $\pm$ 0.76	97.92 $\pm$ 1.03
	Baseline	79.87 $\pm$ 13.16	46.67 $\pm$ 31.51	95.46 $\pm$ 3.60	96.14 $\pm$ 2.17	94.73 $\pm$ 3.64
	Age+Time	<b>94.83 <math>\pm</math> 2.30</b>	<b>65.00 <math>\pm</math> 33.54</b>	96.59 $\pm$ 3.99	97.07 $\pm$ 2.01	96.08 $\pm$ 3.47
	Baseline+Time	91.75 $\pm$ 7.49	55.00 $\pm$ 27.39	96.18 $\pm$ 5.19	96.72 $\pm$ 2.99	95.52 $\pm$ 5.06
Combined	Gender	46.62 $\pm$ 4.73	46.40 $\pm$ 13.87	46.85 $\pm$ 5.97	61.23 $\pm$ 4.95	46.84 $\pm$ 5.45
	Education	57.51 $\pm$ 2.88	40.44 $\pm$ 6.11	74.58 $\pm$ 0.49	82.33 $\pm$ 0.26	73.58 $\pm$ 0.36
	Age	95.19 $\pm$ 2.82	57.13 $\pm$ 8.71	98.99 $\pm$ 0.37	97.71 $\pm$ 0.30	97.76 $\pm$ 0.31
	Time	95.63 $\pm$ 2.65	70.29 $\pm$ 9.09	97.91 $\pm$ 1.26	97.36 $\pm$ 0.74	97.1 $\pm$ 1.04



(a) Combined ROC curves for the LLFS dataset.



(b) Combined ROC curves for the FHS dataset.



(c) Combined ROC curves for the combined dataset.

Baseline	$94.96 \pm 2.72$	$60.59 \pm 12.75$	$98.60 \pm 0.39$	$97.52 \pm 0.37$	$97.49 \pm 0.36$
Age+Time	<b><math>97.56 \pm 1.41</math></b>	<b><math>71.32 \pm 15.51</math></b>	$98.78 \pm 0.62$	$98.02 \pm 0.39$	$97.97 \pm 0.45$
Baseline+Time	$97.44 \pm 1.40$	$70.15 \pm 15.25$	$98.92 \pm 0.36$	$98.09 \pm 0.26$	$98.08 \pm 0.21$

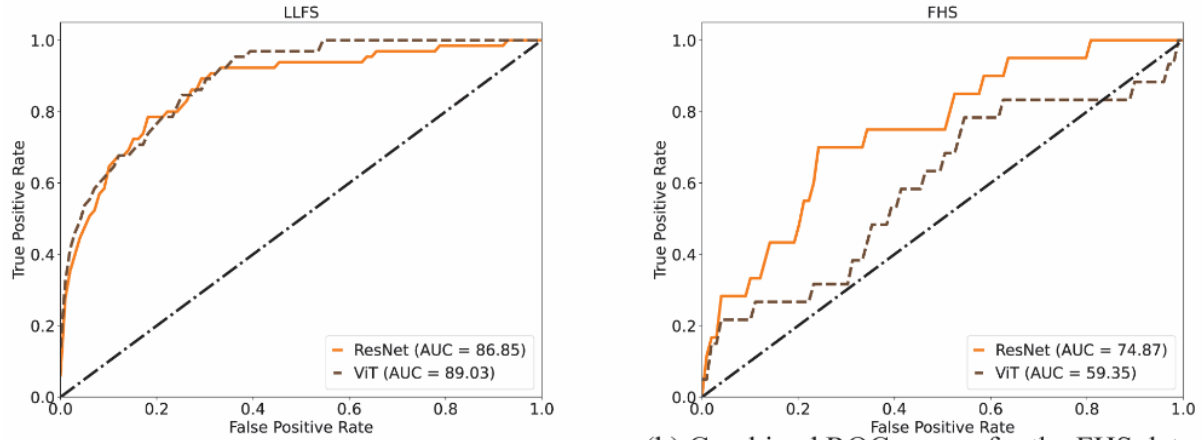
Time: completion time; Baseline: age, gender, and education.

293 **Figure 4:** Mean ROC curves of non-image models using demographics and completion time. The models were trained  
 294 and evaluated on the LLFS, FHS, and combined datasets, respectively, to differentiate participants with dementia from  
 295 those with normal cognition.

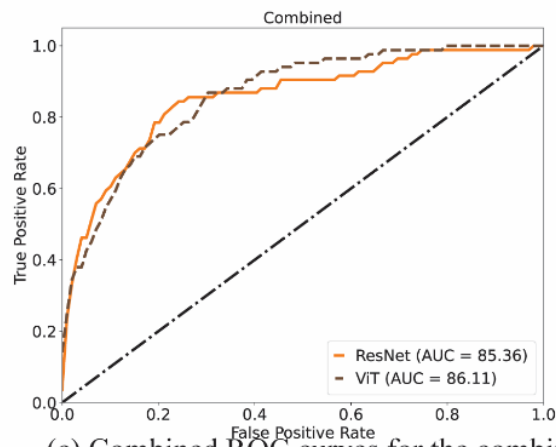
296 *Image models*

297 Our comparison of two image models, ViT and ResNet, on the LLFS and FHS datasets revealed  
298 performance disparity (cf. Figure 5). In the LLFS dataset, ViT outperformed ResNet, achieving a  
299 higher AUC of 89.03% compared to 86.85% from ResNet. In contrast, in the FHS dataset, ResNet  
300 demonstrated significantly better performance with an AUC of 74.87% versus 59.35% from ViT.  
301 When evaluating the combined dataset of LLFS and FHS, the models achieved more comparable  
302 performance. ViT achieved an AUC of 86.11%, while ResNet achieved 85.36%. The variability  
303 between the datasets appears to balance out when examining the aggregated dataset. Additionally,  
304 the image-only models achieved AUCs between 59.35% to 89.03%, indicating room for  
305 improvement in predictive performance by incorporating additional features like demographics and  
306 completion time.

307



(a) Combined ROC curves for the LLFS dataset. (b) Combined ROC curves for the FHS dataset.



(c) Combined ROC curves for the combined dataset.

308 **Figure 5:** Mean ROC curves of ViT and ResNet models using the TMT-A drawings. The models were trained and  
 309 evaluated on the LLFS, FHS, and combined datasets, respectively, to differentiate participants with dementia from  
 310 those with normal cognition.

## 311 *Fusion models*

312 The fusion models, detailed in Table 4, incorporated non-imaging features along with probability  
313 scores derived from fine-tuned vision models to enhance dementia identification. Across all  
314 datasets, the AUC exceeded 96% for each fusion model, demonstrating strong discriminative  
315 performance in differentiating between demented and cognitively healthy individuals. The best  
316 AUC (98.74%) was achieved with the LLFS dataset using age, completion time, and the  
317 probability score predicted by ViT. The highest sensitivity (92.32%) was obtained from the LLFS  
318 dataset using age, completion time, and the probability score derived from ResNet. Overall, ResNet  
319 outperformed ViT in terms of sensitivity, given that there was no significant difference in AUCs  
320 and other metrics. Models utilizing age, completion time, and the image score from ResNet  
321 achieved better AUCs and sensitivities compared to models that additionally incorporate gender  
322 and education level, alongside age, completion time, and image scores from ResNet, in both the  
323 FHS and the combined datasets. Therefore, the optimal feature combination is age, completion  
324 time, and the probability score from ResNet. For the LLFS dataset, this combination yielded an  
325 AUC of 98.62%, sensitivity of 87.62%, and specificity of 98.26%. Similarly, for the FHS dataset,  
326 it achieved an AUC of 96.51%, sensitivity of 85.00%, and specificity of 96.75%.

327 Figure 6 and Table 5 show the performance improvements by adding completion time alone and  
328 both completion time and ResNet-derived image score to the models using only age as a single  
329 predictor. Adding completion time alone demonstrated moderate gains while incorporating both  
330 temporal data and visual signals led to substantial improvements across all datasets. Notably, the  
331 addition of temporal information and visual signals to the demographic data resulted in marked  
332 increases in sensitivity across all datasets. Specifically, for the LLFS dataset, we observed a 3.08%

333 improvement by adding temporal information to age, and a 9.23% improvement by incorporating  
334 both temporal and visual signals in terms of sensitivity. For the FHS dataset, the improvements  
335 were 18.33% and 38.33%, respectively; for the combined dataset, the increases were 14.19% and  
336 16.55%, respectively. Moreover, AUC also exhibited a pronounced boost with the full fusion model  
337 for the FHS dataset. These gains indicated that fusing the visual signal from a participant's drawing  
338 and the time taken to complete the TMT-A could considerably enhance predictive performance.

339 The ablation study in Appendix D compared different model configurations to assess the impact  
340 of various features on model performance. The study revealed that while the Age+Image model  
341 improved sensitivity compared to the traditional Age+Time model, it slightly compromised other  
342 metrics. The proposed Age+Time+Image model consistently outperformed others across most  
343 metrics, demonstrating the value of including completion time. Further analysis of error-related  
344 features collected by trained examiners during the test showed no significant improvement when  
345 added to the Age+Time+Image model. This suggests that the ResNet-derived image probability  
346 likely captures this information. In conclusion, the Age+Time+Image model provides the best  
347 balance of predictive power and ease of administration, aligning with the goal of developing an  
348 accessible dementia screening tool.

349 Since the risk of developing Alzheimer's disease and related dementias increases with aging, we  
350 reapplied our pipeline to a subset of the combined dataset that included only participants aged 65  
351 years and older. Table 6 reports the results for this subsample.

352 As depicted in Figure 7, age was the dominant factor with a significantly higher coefficient than  
353 all other features. After excluding younger participants, age remained the top contributing feature,  
354 but a decrease in the coefficient was observed. Although the probability score derived from ResNet

355 ranks as the least important factor in distinguishing individuals with cognitive dysfunction across  
356 most datasets, its coefficients are notably above zero. This suggests that the visual subtleties  
357 captured by ResNet offer a complementary signal.

358

359 Table 4: Performance metrics of fusion models using combinations of demographics, completion  
360 time, and probability score derived from ResNet and ViT. The models were trained and evaluated  
361 on the LLFS, FHS, and combined datasets, respectively, to differentiate participants with dementia  
362 from those with normal cognition. This table summarized the mean values  $\pm$  standard deviations  
363 in percentages for a 5-fold cross-validation process. Evaluation metrics included AUC, sensitivity,  
364 specificity, F1 score, and accuracy.

	Dataset	AUC	Sensitivity	Specificity	F1 score	Accuracy
LLFS	Age+Time+Image (ResNet)	<b>98.62 <math>\pm</math> 0.91</b>	87.69 $\pm$ 8.77	98.26 $\pm$ 0.67	97.97 $\pm$ 0.58	97.83 $\pm$ 0.66
	Baseline+Time+Image (ResNet)	98.59 $\pm$ 0.97	<b>92.31 <math>\pm</math> 7.69</b>	98.25 $\pm$ 0.93	98.17 $\pm$ 0.57	98.01 $\pm$ 0.72
FHS	Age+Time+Image (ResNet)	<b>96.51 <math>\pm</math> 2.37</b>	<b>85.00 <math>\pm</math> 22.36</b>	96.75 $\pm$ 3.40	97.47 $\pm$ 1.93	96.56 $\pm$ 3.12
	Baseline+Time+Image (ResNet)	94.26 $\pm$ 6.09	80.00 $\pm$ 20.92	96.10 $\pm$ 4.11	97.03 $\pm$ 2.44	95.84 $\pm$ 3.97
Combined	Age+Time+Image (ResNet)	<b>97.95 <math>\pm</math> 1.02</b>	<b>73.68 <math>\pm</math> 7.50</b>	98.81 $\pm$ 0.52	98.14 $\pm$ 0.31	98.08 $\pm$ 0.39
	Baseline+Time+Image (ResNet)	97.80 $\pm$ 1.26	70.15 $\pm$ 9.69	98.88 $\pm$ 0.68	98.08 $\pm$ 0.37	98.04 $\pm$ 0.48
LLFS	Age+Time+Image (ViT)	<b>98.74 <math>\pm</math> 0.75</b>	80.00 $\pm$ 16.85	98.90 $\pm$ 0.54	98.14 $\pm$ 0.48	98.14 $\pm$ 0.38
	Baseline+Time+Image (ViT)	98.66 $\pm$ 0.89	<b>81.54 <math>\pm</math> 13.97</b>	98.97 $\pm$ 0.62	98.28 $\pm$ 0.66	98.26 $\pm$ 0.65
FHS	Age+Time+Image (ViT) FHS	<b>97.71 <math>\pm</math> 1.21</b>	61.67 $\pm$ 19.19	98.62 $\pm$ 1.67	98.32 $\pm$ 0.95	98.08 $\pm$ 1.48
	Baseline+Time+Image (ViT)	97.29 $\pm$ 1.95	<b>71.67 <math>\pm</math> 24.01</b>	97.81 $\pm$ 2.50	97.98 $\pm$ 1.59	97.44 $\pm$ 2.44
Combined	Age+Time+Image (ViT)	<b>98.15 <math>\pm</math> 0.68</b>	66.55 $\pm$ 19.64	99.10 $\pm$ 0.72	98.11 $\pm$ 0.25	98.15 $\pm$ 0.26
	Baseline+Time+Image (ViT)	98.10 $\pm$ 0.84	<b>70.15 <math>\pm</math> 16.86</b>	98.63 $\pm$ 0.86	97.87 $\pm$ 0.33	97.80 $\pm$ 0.47

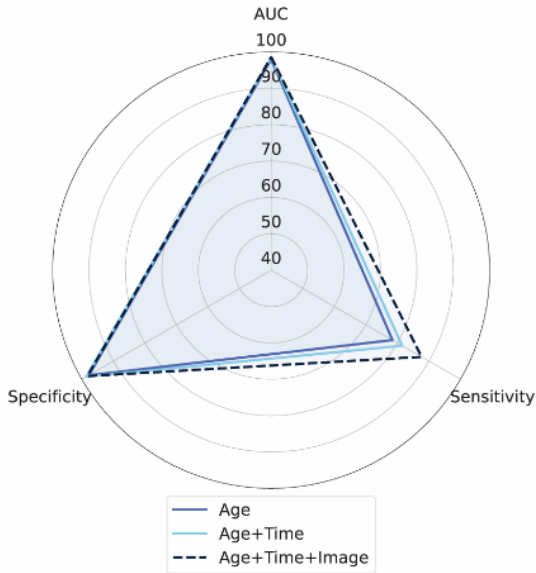
Time: completion time; Baseline: age, gender, and education; Image (ResNet): probability score derived from ResNet; Image (ViT): probability score derived from ViT.

365

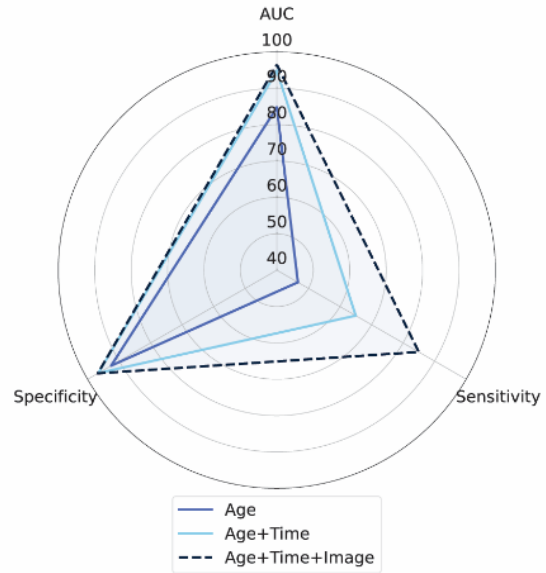
366 **Table 5:** Performance metrics for models using age only, age with completion time, and age with  
367 both completion time and the probability score generated by ResNet. This table compares the  
368 performance of models that use a single demographic feature without any cognitive assessment,  
369 models that incorporate the demographic feature and completion time (commonly used in  
370 traditional cognitive assessments), and models that use demographic, temporal information, and  
371 TMT-A drawing from digital cognitive assessments.

Dataset	Feature	AUC	Sensitivity	Specificity	F1 score	Accuracy
LLFS	Age	97.98 ± 0.78	78.46 ± 10.03	97.81 ± 0.66	97.21 ± 0.59	97.02 ± 0.67
	Age+Time	98.50 ± 0.59	81.54 ± 15.95	98.64 ± 0.48	98.00 ± 0.42	97.95 ± 0.35
	Age+Time+Image (ResNet)	<b>98.62 ± 0.91</b>	<b>87.69 ± 8.77</b>	98.26 ± 0.67	97.97 ± 0.58	97.83 ± 0.66
FHS	Age	84.72 ± 10.68	46.67 ± 31.51	92.55 ± 14.64	94.20 ± 8.70	91.87 ± 13.96
	Age+Time	94.83 ± 2.30	65.00 ± 33.54	96.59 ± 3.99	97.07 ± 2.01	96.08 ± 3.47
	Age+Time+Image (ResNet)	<b>96.51 ± 2.37</b>	<b>85.00 ± 22.36</b>	96.75 ± 3.40	97.47 ± 1.93	96.56 ± 3.12
Combined	Age	95.19 ± 2.82	57.13 ± 8.71	98.99 ± 0.37	97.71 ± 0.30	97.76 ± 0.31
	Age+Time	97.56 ± 1.41	71.32 ± 15.51	98.78 ± 0.62	98.02 ± 0.39	97.97 ± 0.45
	Age+Time+Image (ResNet)	<b>97.95 ± 1.02</b>	<b>73.68 ± 7.50</b>	98.81 ± 0.52	98.14 ± 0.31	98.08 ± 0.39

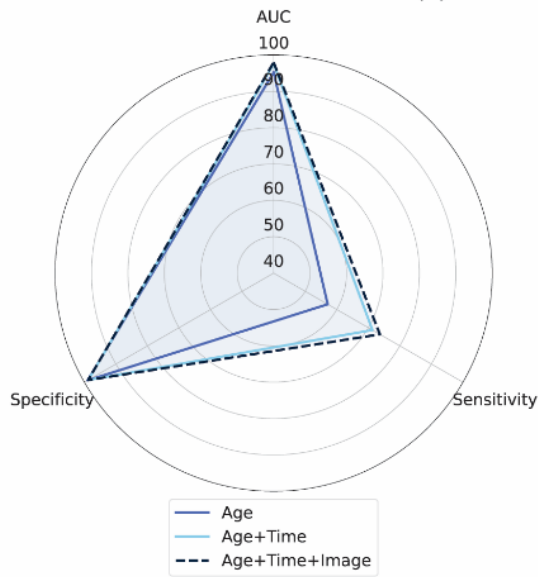
Time: completion time; Image (ResNet): probability score derived from ResNet.



(a) The LLFS dataset.



(b) The FHS dataset.



(c) The combined dataset.

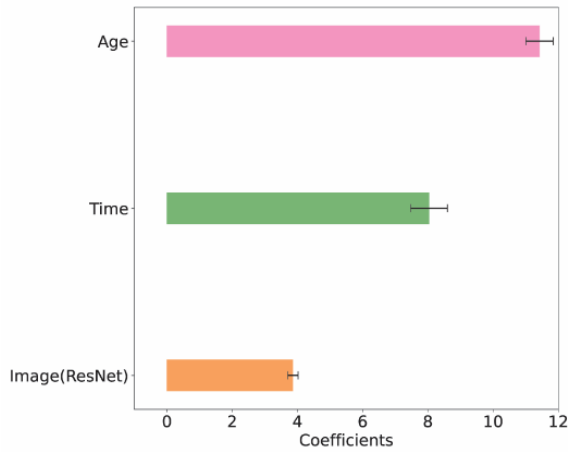
372 **Figure 6:** Comparison of key performance metrics: from age to integrating temporal and visual information. This  
 373 figure presents radar plots for the LLFS, FHS, and combined datasets, respectively. Each plot highlights the  
 374 improvements in three key metrics — AUC, sensitivity, and specificity — when augmenting the baseline model first  
 375 with the completion time alone, and then with both completion time and ResNet-derived image score. Other metrics,  
 376 not showing significant improvements, were excluded from this visualization.

377 **Table 6:** Performance metrics of models trained and evaluated on the subset of older adults aged  
 378 65 and over from both cohorts. This table summarizes mean values  $\pm$  standard deviations in  
 379 percentages for a five-fold cross-validation process.

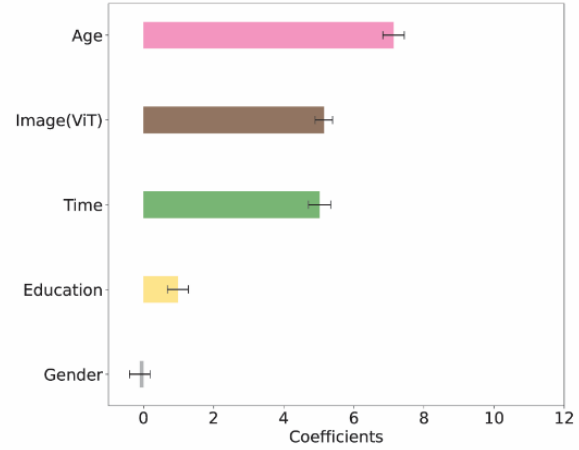
Feature	AUC	Sensitivity	Specificity	F1 score	Accuracy
Gender	47.42 $\pm$ 4.34	45.95 $\pm$ 11.78	48.89 $\pm$ 6.44	62.14 $\pm$ 5.36	48.76 $\pm$ 5.87
Education	56.28 $\pm$ 4.44	40.96 $\pm$ 9.80	71.61 $\pm$ 2.05	79.34 $\pm$ 1.12	70.37 $\pm$ 1.75
Age	94.28 $\pm$ 2.95	56.62 $\pm$ 9.17	98.38 $\pm$ 0.92	96.68 $\pm$ 1.04	96.7 $\pm$ 1.16
Time	94.38 $\pm$ 1.61	65.00 $\pm$ 11.12	97.67 $\pm$ 0.72	96.51 $\pm$ 0.62	96.36 $\pm$ 0.70
Baseline	94.15 $\pm$ 2.69	60.15 $\pm$ 12.87	98.28 $\pm$ 0.77	96.76 $\pm$ 0.97	96.74 $\pm$ 1.03
Age+Time	96.57 $\pm$ 0.97	71.18 $\pm$ 7.40	98.38 $\pm$ 1.05	97.35 $\pm$ 0.50	97.28 $\pm$ 0.72
Baseline+Time	96.41 $\pm$ 1.31	72.28 $\pm$ 3.28	98.43 $\pm$ 0.61	97.44 $\pm$ 0.43	97.38 $\pm$ 0.53
Age+Time+Image (ResNet)	97.04 $\pm$ 1.49	69.85 $\pm$ 7.28	98.63 $\pm$ 0.63	97.49 $\pm$ 0.31	97.47 $\pm$ 0.41
Baseline+Time+Image (ResNet)	97.09 $\pm$ 1.90	68.60 $\pm$ 8.01	98.63 $\pm$ 0.73	97.44 $\pm$ 0.32	97.43 $\pm$ 0.44

Time: completion time; Baseline: age, gender, and education; Image (ResNet): probability score derived from ResNet.

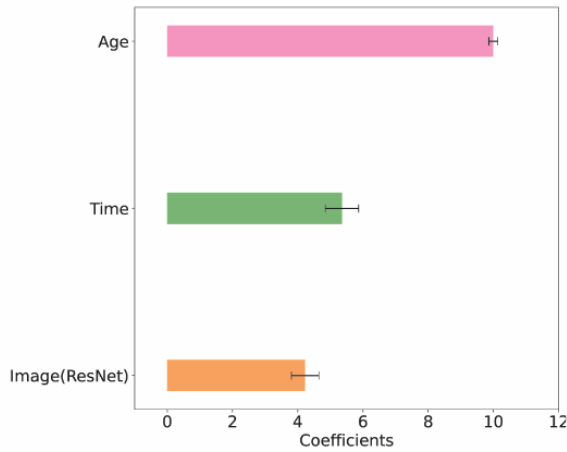
380



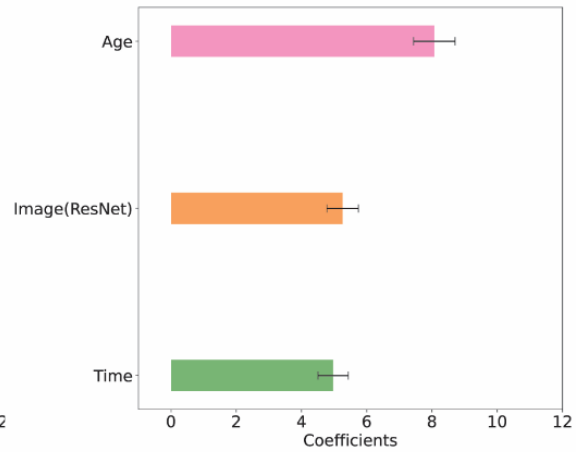
(a) The combined dataset.



(b) The combined dataset including only participants aged 65 years and older.



(c) The LLFS dataset.



(d) The FHS dataset.

381 **Figure 7:** Feature importance for the best-performing fusion models. This figure shows the mean coefficients of the  
 382 fusion models using the combination of age, completion time, and probability score derived from ResNet. These  
 383 coefficients were obtained from the fusion models trained on the combined dataset, the subset of older adults aged 65  
 384 and over, the LLFS dataset, and the FHS dataset, respectively  
 385 .

## 386 **Discussion**

387 This work demonstrated the ability of multimodal AI models to distinguish individuals with  
388 dementia from cognitively healthy controls. We found that integrating visual, temporal, and  
389 demographic data significantly enhances performance, surpassing single-modality models. That  
390 is, models using age alone are less predictive. Incorporating completion time with age adds value,  
391 showing substantially improved predictive performance. Furthermore, the inclusion of the visual  
392 predictor, TMT-A image, enhances discriminative capability beyond age and completion time. As  
393 shown in Table 5, incorporating the predictor of completion time, used in traditional cognitive  
394 assessments, along with age, results in sensitivity improvements of 3.08% and 17.33% for the  
395 LLFS and FHS datasets, respectively, compared to age alone. Moreover, the inclusion of the visual  
396 predictor, TMT-A image, significantly increases sensitivity by 6.15% and 20% when integrated  
397 with the traditional cognitive predictors, age, and completion time, for the LLFS and FHS datasets,  
398 respectively. In addition, we showed that the discriminative capability of these fusion models  
399 remains consistent across various cohorts. Our findings highlight the potential of utilizing  
400 geographically diverse data sources to automate large-scale dementia screening. Importantly, our  
401 evaluation of the clinical significance of these improvements reveals meaningful real-world impact  
402 (see Appendix E). Using established metrics such as clinical significance ratio and Number  
403 Needed to Diagnose (NND), we demonstrate that the fusion model captures a substantial  
404 proportion of previously missed cases, with particularly strong performance for the FHS dataset  
405 where an NND of only 3.33 indicates a good diagnostic efficiency in clinical practice.

406 One strength of our framework is its accessibility. While neuroimaging and biofluid analyses are  
407 sensitive and accurate for diagnosing Alzheimer’s disease, their limited accessibility and reliance

408 on specialized personnel make them less practical, particularly in resource-limited regions.<sup>38,39</sup> In  
409 contrast, our proposed framework addresses these challenges through deliberate design choices.  
410 We opted not to use ApoE. Individuals in resource-limited regions face challenges in obtaining  
411 genetic information, such as limited access to genetic testing facilities, the costs associated with  
412 testing, and the availability of trained professionals to administer a blood test. Similarly, we made  
413 a strategic decision not to use raw digital pen recordings directly. Instead, we extracted the overall  
414 completion time and converted the raw data into images, then fine-tuned deep learning models to  
415 learn spatial information from TMT-A drawings. This approach represents a deliberate  
416 prioritization of accessibility over the capture of fine-grained data. Digital pen technology, while  
417 capable of recording nuanced information about drawing dynamics, timing, and pressure, remains  
418 inaccessible in many settings due to cost, technological barriers, and unfamiliarity among older  
419 adults – our primary target population. Our best-performing model utilizes age, TMT-A  
420 completion time, and TMT-A drawings – all readily obtainable variables that eliminate the need  
421 for specialized equipment. This approach enables diverse implementation pathways: individuals  
422 can complete the test using standard pen and paper, time it with a basic stopwatch, and then digitize  
423 the drawing through various means. For example, drawings can be collected through traditional  
424 paper-and-pencil methods in clinical settings, completed at home and photographed using a  
425 smartphone camera, or captured via telehealth platforms during remote assessments. This  
426 versatility extends the potential reach of automated cognitive assessment by supporting both  
427 prospective implementation across diverse clinical environments and retrospective analysis of  
428 existing datasets. This flexibility addresses a critical gap in cognitive assessment infrastructure  
429 globally, particularly in settings where specialized neuropsychological expertise is limited. The  
430 performance of our image-only models suggests that essential diagnostic information is preserved

431 in the spatial patterns of TMT-A drawings. This finding is consistent with a growing body of  
432 evidence supporting image-based approaches for cognitive assessments. In the clock drawing tests,  
433 several works have demonstrated that deep learning models can effectively screen for dementia and  
434 identify cognitive decline using the visual features of drawings.<sup>15,40,41</sup> This pattern extends to other  
435 neuropsychological instruments, including the Rey-Osterrieth Complex Figure Test, where image-  
436 based deep learning approaches have shown comparable efficacy in predicting cognitive  
437 impairment.<sup>13</sup> These findings across multiple drawing-based assessments suggest that spatial  
438 patterns and visual features captured in static images contain sufficient information for effective  
439 cognitive evaluation.

440 Another advantage of our framework is its scalability. We selected TMT-A over the more  
441 challenging TMT-B to facilitate large-scale screening. TMT-A, which simply involves connecting  
442 Arabic numerals, is easier and more universally recognized than TMT-B, which requires  
443 alternating between numerals and alphabetic letters. This difference is particularly relevant for  
444 older adults, non-English speaking populations in regions with lower education levels, where the  
445 complexity of TMT-B may lead to a higher failure rate due to misunderstood instructions or the  
446 difficulty of the task, as demonstrated in previous studies.<sup>42,43</sup> The simplicity and broader global  
447 recognition of TMT-A make it a more practical choice for implementing our framework widely.

448 In addition, our framework can adapt to other cohorts, as we use multi-center, multi-cohort data  
449 from geographically diverse populations across various states and nations. Our framework was  
450 applied to and validated on both individual and combined cohorts, demonstrating its  
451 generalizability across different cohorts. This versatility indicates the potential for global adaptation  
452 of our method for large-scale cognitive impairment screening, particularly in settings where  
453 dementia expertise is limited or unavailable.

454 Our framework demonstrates broad applicability to general populations, as it was developed using  
455 community-based cohorts comprising participants with diverse cognitive profiles. Unlike  
456 dementia-specific cohorts, which primarily include individuals with diagnosed dementia or  
457 elevated risk factors, our community-based approach captures a more representative spectrum of  
458 cognitive function. Consequently, our model provides a more ecologically valid approach for  
459 identifying dementia cases within heterogeneous populations. The performance characteristics  
460 observed in our models, particularly the trade-off between sensitivity and specificity, should be  
461 interpreted within this community-based research context. While specialized clinical settings such  
462 as Alzheimer’s Disease Research Centers (ADRCs) or memory clinics offer comprehensive  
463 assessments with higher diagnostic accuracy and potentially more balanced case distributions, they  
464 introduce selection bias. Our approach reflects the natural prevalence of cognitive impairment in  
465 general populations, where most individuals are cognitively normal. This better represents the  
466 intended application environment for accessible screening tools and minimizes potential  
467 distributional shift during real-world implementation. A better dementia care pathway would likely  
468 combine approaches like ours with comprehensive assessment, using community-based digital  
469 screening for initial identification and referral of high-risk individuals to specialized centers,  
470 allowing for more efficient allocation of resources while maximizing population coverage.

471 The performance metrics of models trained on the FHS dataset demonstrate greater variability  
472 compared to those trained on the LLFS dataset, as indicated by higher standard deviations across  
473 evaluation metrics. This performance discrepancy can be attributed to fundamental differences in  
474 cohort characteristics and study designs (detailed in Appendix F). Age distribution represents a  
475 primary factor influencing model performance, with FHS participants being younger on average  
476 ( $70.25 \pm 12.42$  years) compared to those in LLFS ( $71.74 \pm 10.74$  years). This younger demographic

477 profile corresponds with a substantially lower prevalence of dementia in FHS (1.52% vs. 4.03% in  
478 LLFS), resulting in fewer dementia cases available for model training (19 cases vs. 65 cases), and  
479 thereby creating a significant class imbalance challenge for deep learning algorithms. Variations  
480 in performance may also stem from differences in cohort selection criteria. While the FHS primarily  
481 examines cardiovascular health in a community-based population, the LLFS cohort consists of  
482 individuals selectively recruited from families with exceptional longevity, emphasizing factors  
483 associated with healthy aging.<sup>44</sup> These divergent selection criteria result in populations with distinct  
484 cognitive aging profiles. Data collection settings also differ meaningfully between the two studies.  
485 LLFS conducted in-home assessments, potentially including participants who may be unable to  
486 travel to clinical settings due to health or mobility limitations. This approach may capture a  
487 broader spectrum of cognitive impairment compared to FHS, where all assessments required clinic  
488 visits. The cumulative effect of these differences manifests in the increased variability observed in  
489 FHS model performance, reflecting the inherent challenges of detecting subtle cognitive changes  
490 in a relatively younger population with fewer overt cases of dementia.

491 Although this study offers valuable insights, it also has several limitations. Our study is limited  
492 by the small sample size and highly imbalanced dataset. The area under the precision-recall curve  
493 (AUPRC) revealed low values, likely due to the scarcity of dementia cases. While our models  
494 achieved decent sensitivity, this came at the cost of precision, which lowered the AUPRC. These  
495 limitations highlight the need for larger datasets with more dementia cases. However, addressing  
496 this challenge is inherently difficult as our datasets are derived from community-based studies  
497 where the prevalence of dementia is naturally low.

498 We compromised on some clinically valuable information to enhance the feasibility of  
499 implementation in different settings by choosing not to use the time-series data collected by digital

500 pens, despite its potential to provide additional clinically meaningful insights. Features such as  
501 stroke-by-stroke timing, pauses between strokes, and writing pressure variations contain  
502 information about cognitive processing speed, motor control, and executive function that is not fully  
503 captured in static images.<sup>45-47</sup> As digital pen technology becomes more widely available and  
504 affordable, future research could compare the performance of models trained on TMT images  
505 versus those utilizing digital pen data. Such studies would help quantify the practical impact of  
506 this tradeoff between data richness and accessibility. Additionally, future work could explore  
507 hybrid approaches that maintain accessibility while incorporating selective temporal features that  
508 can be captured without specialized equipment, such as developing smartphone applications that  
509 capture basic temporal dynamics during test administration.

510 A key limitation of our study was the exclusion of TMT-B data. While TMT-A primarily measures  
511 processing speed through a simple number-sequence task, TMT-B offers a more comprehensive  
512 evaluation of executive functions through its complex alternation between numbers and letters,  
513 making it particularly sensitive for dementia detection. TMT-A's simpler design means its  
514 performance might be more influenced by motor impairments related to factors like frailty, rather  
515 than cognitive decline alone. Due to limited TMT-B data availability in our dataset for  
516 comprehensive experimentation, we focused our analysis on TMT-A. Future studies should  
517 incorporate both tests to develop more comprehensive screening models that can better  
518 differentiate between motor and cognitive impairments.

519 Another limitation is that both FHS and LLFS do not conduct a comprehensive review to determine  
520 the cognitive status of all participants. Instead, the FHS focuses on those at higher risk, prioritizing  
521 dementia reviews for them. Participants who are at lower risk, show no signs of cognitive concerns,  
522 or are relatively young, are presumed to have normal cognition without undergoing a detailed

523 consensus review. It is important to note that the FHS confirms cognitive status, including normal  
524 cognition, at the time of death. Similarly, in the LLFS cohort, participants were selected for review  
525 if they had a clinical dementia rating score greater than 0 or if they were labeled as having cognitive  
526 impairment consistent with dementia by a diagnostic algorithm that considered sex and specific  
527 cognitive scores.<sup>48</sup> Consequently, both datasets include both confirmed cases of normal cognition,  
528 determined through a dementia review, and presumed cases of normal cognition.

529 While our study focuses on distinguishing normal cognition from dementia, we recognize the  
530 clinical importance of detecting earlier stages of cognitive decline. Our preliminary investigations  
531 explored multiple classification scenarios across the cognitive spectrum, including MCI detection  
532 (see Appendix G for comprehensive analysis). These preliminary analyses revealed that our fusion  
533 model performed significantly better when distinguishing between the more distinct cognitive states  
534 of normal cognition and dementia, compared to classifying the more subtle differences between  
535 normal cognition and MCI. This guided our decision to focus the current study on normal vs.  
536 dementia classification, establishing a strong methodological foundation. Future research should  
537 build upon this foundation to develop specialized approaches for the more challenging task of early  
538 MCI detection, potentially incorporating longitudinal data and additional neuropsychological  
539 measures to capture subtle cognitive changes.

540 Future research directions could explore several key areas to enhance model performance. More  
541 extensive fine-tuning of both ResNet and ViT models, allowing additional layers to be updated  
542 during training, may improve the models' ability to learn from limited datasets. For ViT  
543 specifically, investigating the impact of different patch sizes, particularly smaller patches, might  
544 capture more fine-grained information. Given the sparse nature of TMT drawings, which consist  
545 primarily of line segments, future work should also explore architectures specifically designed for

546 sketch analysis, such as graph neural networks. These specialized architectures might better  
547 capture the sequential and structural aspects of TMT drawings.

548 Although our data collection spans geographically diverse locations, there remains a need to  
549 include more demographically and ethnically diverse populations to enhance the generalizability  
550 of our framework. We have already begun addressing this limitation by expanding beyond the  
551 initial FHS cohort through integration with the LLFS dataset, which introduced greater  
552 demographic variability. We also explored the National Alzheimer’s Coordinating Center’s  
553 Uniform Data Set (NACC UDS) for its potential to enhance ethnic diversity through multi-center  
554 cohorts, but encountered methodological challenges as UDS contains TMT completion times  
555 without corresponding drawings. Looking ahead, we are implementing several practical validation  
556 strategies. First, we have developed a web-based platform (developed in our previous CDT work)  
557 that targets global participation across varied demographic and ethnic backgrounds, allowing us to  
558 gather TMT drawings from diverse populations worldwide.<sup>15</sup> As this dataset reaches sufficient  
559 sample size, we will implement external validation and periodic model retraining to improve  
560 generalizability. Additionally, we plan to establish collaborative partnerships with international  
561 research centers to validate our models, particularly focusing on underrepresented populations and  
562 resource-limited environments. Finally, integrating factors such as race, ethnicity, and other  
563 sociodemographic characteristics into our models will further enhance their predictive power across  
564 diverse populations.

565 **Acknowledgements**

566 None

567 **Ethical Considerations**

568 All participants provided written informed consent. Study protocols and consent forms were  
569 approved by the Boston University Medical Campus Institutional Review Board and the  
570 Institutional Review Boards of the Long Life Family Study field sites as well as the Long Life  
571 Family Study coordinating center at Washington University St. Louis.

572 **Consent to Participate**

573 Written informed consent was obtained from all participants.

574 **Consent for Publication**

575 All authors reviewed and approved the final manuscript for publication.

576 **Funding**

577 This research was partially supported by the NSF under grants CCF-2200052, DMS-1664644,  
578 ECCS-2317079, and IIS-1914792; by the ONR under grant N00014-19-1-2571; by the DOE under  
579 grant DE-AC02-05CH11231; by the NIH under grant UL54 TR004130; by Boston University; by  
580 the Karen Toffler Charitable Trust; by National Institute on Aging's Artificial Intelligence and  
581 Technology Collaboratories (AITC) for Aging Research program under grant P30-AG073104; by  
582 the American Heart Association under grant 20SFRN35460031; by Gates Ventures and National  
583 Institutes of Health under grants RF1-AG062109, U19-AG068753; by the Framingham Heart

584 Study's National Heart, Lung, and Blood Institute contract under grants N01-HC-25195,  
585 HHSN268201500001I, and by the NIH National Institute on Aging under grants AG008122,  
586 AG016495, AG062109, and AG068753. The LLFS was supported by National Institute on Aging  
587 under grants U01AG023746, U01AG023712, U01AG023749, U01AG023755, U01AG023744,  
588 and U19 AG063893; SLA was supported by National Institute on Aging under grant K01AG057798.

## 589 **Declaration of Conflicting Interests**

590 The authors declare no conflicts of interest.

## 591 **Data Availability Statement**

592 The data for the study is available after a reasonable request to the authors and approval from the  
593 FHS and LLFS.

## 594 **Author Contributions**

595 Jingmei Yang (Formal analysis; Software; Investigation; Visualization; Writing - Original draft;  
596 Writing - Review & Editing); Samad Amini (Formal analysis; Investigation; Writing - Review &  
597 Editing); Boran Hao (Formal analysis; Investigation; Writing - Review & Editing); Seho Park  
598 (Data curation; Writing - Review & Editing); Cody Karjadi (Data curation; Writing - Review &  
599 Editing); Lance San Souci (Data curation; Writing - Review & Editing); Vijaya B. Kolachalama  
600 (Conceptualization; Formal analysis; Writing - Review & Editing); Stephanie Cosentino  
601 (Conceptualization; Formal analysis; Writing - Review & Editing); Stacy Andersen  
602 (Conceptualization; Data curation; Formal analysis; Investigation; Writing - Review & Editing);  
603 Rhoda Au (Conceptualization; Data curation; Formal analysis; Investigation; Methodology;

- 604 Writing - Review & Editing; Project administration; Supervision; Funding acquisition); Ioannis
- 605 Ch. Paschalidis (Conceptualization; Formal analysis; Investigation; Methodology; Writing -
- 606 Review & Editing; Project administration; Supervision; Funding acquisition).

## References

1. International AD. Dementia Statistics, <https://www.alzint.org/about/dementia-facts-figures/dementia-statistics/> (2023).
2. International AD. World Alzheimer Report 2022, <https://www.alzint.org/resource/world-alzheimer-report-2022/> (2022).
3. Shin J-H. Dementia epidemiology fact sheet 2022. *Ann Rehabil Med* 2022; 46: 53.
4. Organization WH. Ageing and health, <https://www.who.int/news-room/fact-sheets/detail/ageing-and-health> (2022).
5. Amini S, Hao B, Zhang L, et al. Automated detection of mild cognitive impairment and dementia from voice recordings: A natural language processing approach. *Alzheimers Dement* 2023; 19: 946–955.
6. Amini S, Hao B, Yang J, et al. Prediction of Alzheimer’s disease progression within 6 years using speech: a novel approach leveraging language models. *Alzheimers Dement J Alzheimers Assoc* 2024; 20: 5262–5270.
7. Sunderland T, Hill JL, Mellow AM, et al. Clock drawing in Alzheimer’s disease: a novel measure of dementia severity. *J Am Geriatr Soc* 1989; 37: 725–729.
8. Rabin LA, Barr WB, Burton LA. Assessment practices of clinical neuropsychologists in the United States and Canada: A survey of INS, NAN, and APA Division 40 members. *Arch Clin Neuropsychol* 2005; 20: 33–65.
9. Li Y, Guo J, Yang P. Developing an image-based deep learning framework for automatic

- scoring of the pentagon drawing test. *J Alzheimers Dis* 2022; 85: 129–139.
10. Tasaki S, Kim N, Truty T, et al. Explainable deep learning approach for extracting cognitive features from hand-drawn images of intersecting pentagons. *NPJ Digit Med* 2023; 6: 157.
  11. Maruta J, Uchida K, Kurozumi H, et al. Deep convolutional neural networks for automated scoring of pentagon copying test results. *Sci Rep* 2022; 12: 9881.
  12. Ruengchaijatuporn N, Chatnuntawech I, Teerapittayanon S, et al. An explainable self-attention deep neural network for detecting mild cognitive impairment using multi-input digital drawing tasks. *Alzheimers Res Ther* 2022; 14: 1–11.
  13. Park JY, Seo EH, Yoon H-J, et al. Automating Rey Complex Figure Test scoring using a deep learning-based approach: a potential large-scale screening tool for cognitive decline. *Alzheimers Res Ther* 2023; 15: 1–11.
  14. Cheah W-T, Hwang J-J, Hong S-Y, et al. A digital screening system for alzheimer disease based on a neuropsychological test and a convolutional neural network: System development and validation. *MIR Med Inf* 2022; 10: e31106.
  15. Amini S, Zhang L, Hao B, et al. An artificial intelligence-assisted method for dementia detection using images from the clock drawing test. *J Alzheimers Dis* 2021; 83: 581–589.
  16. Handzlik D, Richmond LL, Skiena S, et al. Explainable automated evaluation of the clock drawing task for memory impairment screening. *Alzheimers Dement* 2023; 15: e12441.
  17. Tombaugh TN. Trail Making Test A and B: normative data stratified by age and education. *Arch Clin Neuropsychol* 2004; 19: 203–214.

18. Wei M, Shi J, Li T, et al. Diagnostic accuracy of the Chinese version of the trail-making test for screening cognitive impairment. *J Am Geriatr Soc* 2018; 66: 92–99.
19. Inomoto A, Deguchi J, Fukuda R, et al. Gender-specific factors associated with the Japanese version of the trail making test among Japanese workers. *J Phys Ther Sci* 2023; 35: 547–552.
20. Hashimoto R, Meguro K, Lee E, et al. Effect of age and education on the Trail Making Test and determination of normative data for Japanese elderly people: the Tajiri Project. *Psychiatry Clin Neurosci* 2006; 60: 422–428.
21. Specka M, Weimar C, Stang A, et al. Trail Making Test Normative Data for the German Older Population. *Arch Clin Neuropsychol* 2022; 37: 186–198.
22. Suzuki H, Sakuma N, Kobayashi M, et al. Normative Data of the Trail Making Test Among Urban Community-Dwelling Older Adults in Japan. *Front Aging Neurosci* 2022; 14: 832158.
23. Cangoz B, Karakoc E, Selekler K. Trail Making Test: normative data for Turkish elderly population by age, sex and education. *J Neurol Sci* 2009; 283: 73–78.
24. Fellows RP, Dahmen J, Cook D, et al. Multicomponent analysis of a digital Trail Making Test. *Clin Neuropsychol* 2017; 31: 154–167.
25. Du M, Andersen SL, Cosentino S, et al. Digitally generated trail making test data: analysis using hidden Markov modeling. *Alzheimers Dement* 2022; 14: e12292.
26. Dahmen J, Cook D, Fellows R, et al. An analysis of a digital variant of the Trail Making Test using machine learning techniques. *Technol Health Care* 2017; 25: 251–264.
27. Zhang W, Zheng X, Tang Z, et al. Combination of Paper and Electronic Trail Making Tests

- for Automatic Analysis of Cognitive Impairment: Development and Validation Study. *Med Internet Res* 2023; 25: e42637.
28. Mahmood SS, Levy D, Vasan RS, et al. The Framingham Heart Study and the epidemiology of cardiovascular disease: a historical perspective. *The lancet* 2014; 383: 999–1008.
  29. Wojczynski MK, Jiuan Lin S, Sebastiani P, et al. NIA Long Life family study: Objectives, design, and heritability of cross-sectional and longitudinal phenotypes. *J Gerontol Biol Sci Med Sci* 2022; 77: 717–727.
  30. Au R, Piers RJ, Devine S. How technology is reshaping cognitive assessment: Lessons from the Framingham Heart Study. *Neuropsychology* 2017; 31: 846.
  31. Satizabal CL, Beiser AS, Chouraki V, et al. Incidence of dementia over three decades in the Framingham Heart Study. *N Engl J Med* 2016; 374: 523–532.
  32. Massey Jr FJ. The Kolmogorov-Smirnov test for goodness of fit. *J Am Stat Assoc* 1951; 46: 68–78.
  33. Greenwood PE, Nikulin MS. *A guide to chi-squared testing*. John Wiley & Sons, 1996.
  34. He K, Zhang X, Ren S, et al. Deep residual learning for image recognition. In: *Proc IEEE Comput Soc Conf Comput Vis Pattern Recognit*. 2016, pp. 770–778.
  35. Dosovitskiy A, Beyer L, Kolesnikov A, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *ArXiv Prepr ArXiv201011929*.
  36. Deng J, Dong W, Socher R, et al. Imagenet: A large-scale hierarchical image database. In: *Proc IEEE Comput Soc Conf Comput Vis Pattern Recognit*. IEEE, 2009, pp. 248–255.

37. Hastie T, Tibshirani R, Friedman J. *The elements of statistical learning: Data mining, inference, and prediction*. 2nd ed. Springer, 2009.
38. Olsson B, Lautner R, Andreasson U, et al. CSF and blood biomarkers for the diagnosis of Alzheimer's disease: a systematic review and meta-analysis. *Lancet Neurol* 2016; 15: 673–684.
39. Tartaglia MC, Rosen HJ, Miller BL. Neuroimaging in dementia. *Neurotherapeutics* 2011; 8: 82–92.
40. Chen S, Stromer D, Alabdallah HA, et al. Automatic dementia screening and scoring by applying deep learning on clock-drawing tests. *Sci Rep* 2020; 10: 20854.
41. Sato K, Niimi Y, Mano T, et al. Automated evaluation of conventional clock-drawing test using deep neural network: Potential as a mass screening tool to detect individuals with cognitive decline. *Front Neurol* 2022; 13: 896403.
42. Smith Watts AK, Ahern DC, Jones JD, et al. Trail-making test part B: evaluation of the efficiency score for assessing floor-level change in veterans. *Arch Clin Neuropsychol* 2019; 34: 243–253.
43. Papandonatos GD, Ott BR, Davis JD, et al. Clinical utility of the Trail-Making test as a predictor of driving performance in older adults. *J Am Geriatr Soc* 2015; 63: 2358–2364.
44. Newman AB, Glynn NW, Taylor CA, et al. Health and function of participants in the Long Life Family Study: a comparison with other cohorts. *Aging* 2011; 3: 63.
45. Prange A, Barz M, Heimann-Steinert A, et al. Explainable automatic evaluation of the trail

- making test for dementia screening. In: *Proc SIGCHI Conf Hum Factor Comput Syst.* 2021, pp. 1–9.
46. Prange A, Sonntag D. Modeling users' cognitive performance using digital pen features. *Front Artif Intell* 2022; 5: 787179.
  47. Kobayashi M, Yamada Y, Shinkawa K, et al. Automated early detection of Alzheimer's disease by capturing impairments in multiple cognitive domains with multiple drawing tasks. *J Alzheimers Dis* 2022; 88: 1075–1089.
  48. Cosentino S, Schupf N, Christensen K, et al. Reduced prevalence of cognitive impairment in families with exceptional longevity. *JAMA Neurol* 2013; 70: 867–874.

## 1 **Appendix A: Data preprocessing**

2 In the FHS dataset, one participant, who took over 600 seconds, was identified as an outlier and  
3 excluded from our analysis due to the completion time from that participant significantly deviating  
4 from the distribution, while no outlier was identified in the LLFS dataset. Regarding education  
5 data, nine participants from FHS lacked information. We filled in these missing values using the  
6 most common education level, i.e., college or above. The LLFS dataset had seven such cases, and  
7 the missing values were addressed in a similar manner. In both datasets, all participants reported  
8 their gender. As for the ApoE genotype, the FHS dataset had 45 missing entries, while the LLFS  
9 dataset had 143. These missing entries were replaced with the most common genotype, i.e.,  $\epsilon 3/\epsilon 3$ .  
10 In the FHS cohort, all participants reported their ages. In the LLFS cohort, two participants did not  
11 provide their ages. For those with missing ages, we imputed them using the mean age.  
12 Furthermore, we excluded participants identified with mild cognitive impairment from both the  
13 LLFS and FHS datasets. The missing values in our datasets are relatively low. Specifically, for  
14 education, the missing values are only 0.72% in the FHS dataset and 0.43% in the LLFS dataset. In  
15 terms of age, there are no missing values in the FHS dataset and only 0.12% in the LLFS dataset.  
16 Additionally, there are no missing values for gender in both datasets. For completion time, only a  
17 single case was excluded from the LLFS dataset. Given the relatively low percentages of missing  
18 values for these variables, it is suggested that the missing data and imputation would likely not  
19 have a significant impact on model performance.

20 **Appendix B: Model selection**

21 In this study, a wide range of traditional machine learning models were developed, including Logistic  
22 Regression (LR), Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), and  
23 Extreme Gradient Boosting (XGBoost). These models were built utilizing baseline demographic  
24 features, namely age, education, and gender. As shown in Table A1, the LR model was the best-  
25 performing model. Therefore, it was adopted as the base model for both non-image models and  
26 fusion models.

27 **Table A1:** Model selection.

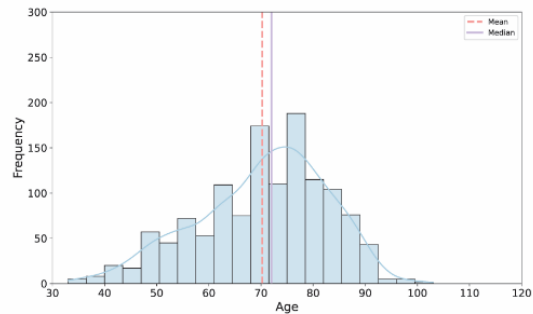
<b>Model</b>	<b>AUC (%)</b>	<b>Sensitivity (%)</b>	<b>Specificity (%)</b>	<b>F1 score (%)</b>	<b>Accuracy (%)</b>
DT	78.9 ± 4.35	56.92 ± 8.77	98.04 ± 1.29	96.47 ± 0.61	96.42 ± 0.94
LR	97.71 ± 0.94	75.39 ± 10.03	98.29 ± 1.52	97.53 ± 1.23	97.39 ± 1.47
RF	89.81 ± 2.89	63.08 ± 12.64	97.41 ± 1.63	96.28 ± 0.73	96.05 ± 1.14
SVM	97.45 ± 0.91	75.39 ± 10.03	97.85 ± 1.5	97.17 ± 0.98	96.96 ± 1.27
XGBoost	95.74 ± 1.03	61.54 ± 7.69	98s.11 ± 1.05	96.73 ± 0.66	96.66 ± 0.86

LR: Logistic Regression; SVM: Support Vector Machine; DT: Decision Tree; RF: Random Forest; XGBoost: Extreme Gradient Boosting.

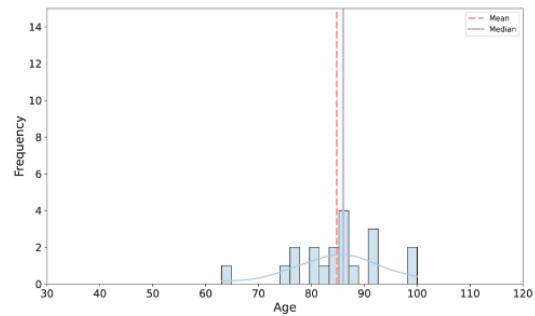
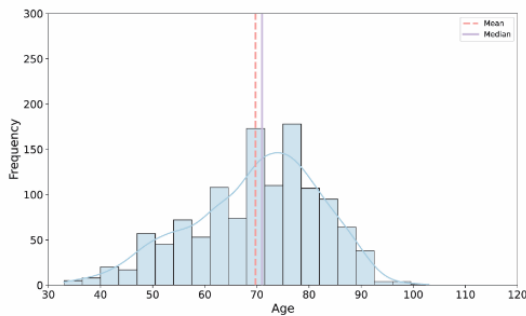
28

## 29 Appendix C: Feature distribution

30 This section presents the feature distributions for the FHS and LLFS datasets. Figure A1 shows  
31 the age distribution for the FHS dataset, while Figure A2 displays the age distribution for the  
32 LLFS dataset. The completion time distributions for the FHS and LLFS datasets are presented in  
33 Figure A3 and Figure A4.

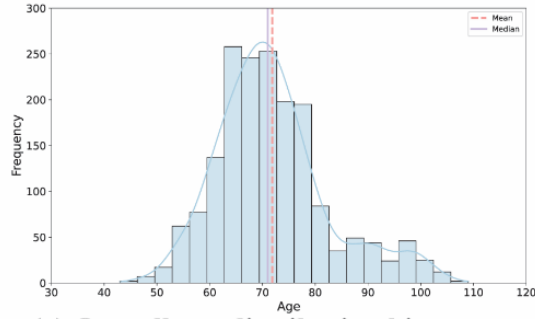


(a) Overall age distribution histogram.

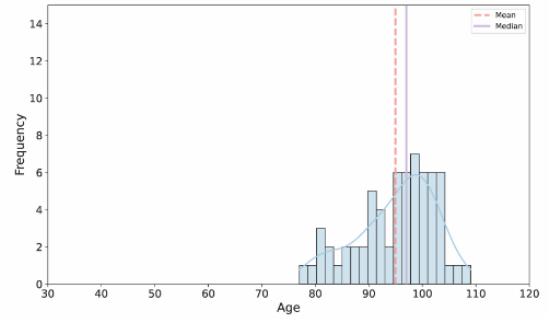
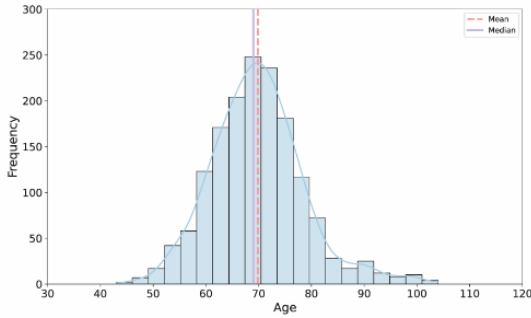


34 (b) Age distribution for the normal cognition group. (c) Age distribution for the dementia group.

35 **Figure A1:** Age distribution for the FHS dataset. (a) Histogram displaying the overall age  
36 distribution. (b-c) Histograms showing the age distributions for the normal cognition and dementia  
37 groups, respectively.

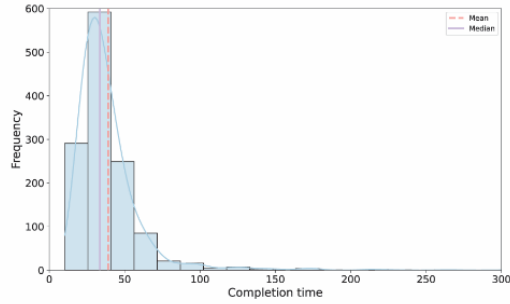


(a) Overall age distribution histogram.

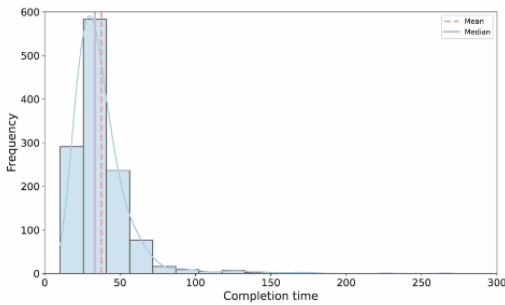


38 (b) Age distribution for the normal cognition group. (c) Age distribution for the dementia group.

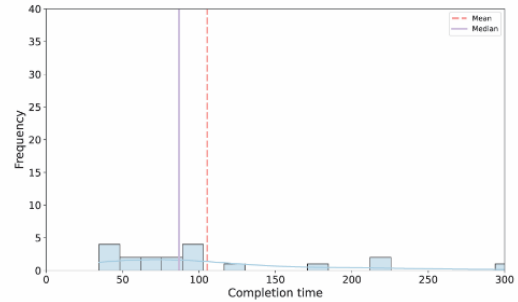
39 **Figure A2:** Age distribution for the LLFS dataset. (a) Histogram displaying the overall age  
 40 distribution. (b-c) Histograms showing the age distributions for the normal cognition and dementia  
 41 groups, respectively.



**(a) Overall completion time distribution histogram.**

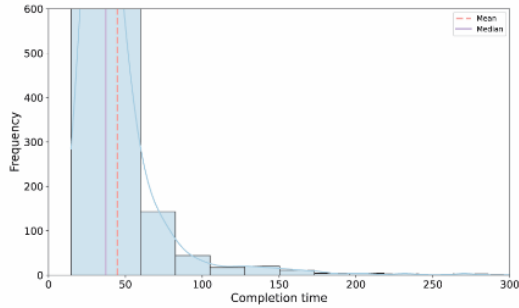


**(b) Completion time distribution for the normal cognition group.**

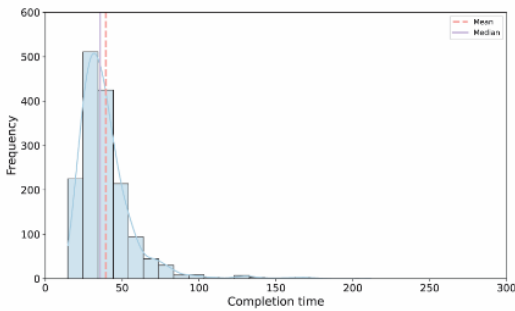


**(c) Completion time distribution for the dementia cognition group.**

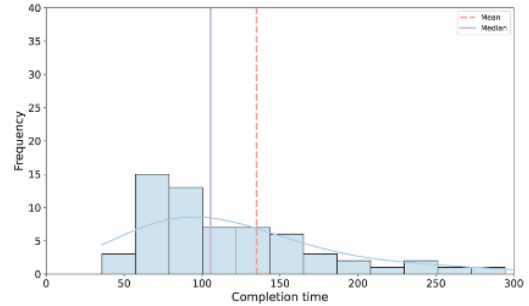
42 **Figure A3:** Completion time distribution for the FHS dataset. (a) Histogram displaying the overall  
 43 distribution of completion times. (b-c) Histograms showing the completion time distributions for  
 44 the normal cognition and dementia groups, respectively. The completion time represents the  
 45 duration taken by participants to complete the tasks.



(a) Overall completion time distribution histogram.



(b) Completion time distribution for the normal cognition group.



(c) Completion time distribution for the dementia cognition group.

46 **Figure A4:** Completion time distribution for the LLFS dataset. (a) Histogram displaying the  
 47 overall distribution of completion times. (b-c) Histograms showing the completion time  
 48 distributions for the normal cognition and dementia groups, respectively. The completion time  
 49 represents the duration taken by participants to complete the task.

## 50 **Appendix D: Comparative analysis of feature combination**

51 We conducted an experiment comparing three models: Age+Time (traditional), Age+Image  
52 (excluding Time), and Age+Time+Image (our proposed model). This experiment aimed to explore  
53 the impact of including completion time. As shown in Table A2, the Age+Image model  
54 demonstrated substantially better sensitivity compared to the Age+Time model across both LLFS  
55 and FHS datasets. However, this improvement came at the cost of slightly lower AUC,  
56 specificity, F1 score, and accuracy. Across all datasets, the Age+Time+Image model consistently  
57 outperformed the Age+Image model in terms of AUC, specificity, F1 score, and accuracy.  
58 Additionally, the Age+Time+Image model achieved higher sensitivity across the FHS and  
59 combined datasets. While the Age+Image model offers decent discriminative power and reduces  
60 the participant effort for timing, the inclusion of completion time in the Age+Time+Image model  
61 offers superior overall performance across most metrics. We propose that the image can be  
62 captured as a snapshot and the time can be tracked using a stopwatch. In settings where  
63 approximate timing is feasible, we recommend using the Age+Time+Image model. Although this  
64 approximation of completion time using a simple stopwatch by the participant may not be as  
65 accurate as that tracked by a digital pen or trained examiner, it still contributes valuable information.  
66 In extremely resource-limited settings or where timing is impractical, the Age+Image model  
67 provides an alternative for initial dementia screenings.

68 We also examined error-related features documented during the TMT administration. We focused  
69 on two key error features: LiftPen, which records the number of times participants lifted their  
70 pen from the paper, and Tremor, indicating the presence of obvious tremor. While examiners also  
71 noted self-corrected and examiner-corrected perceptual errors, these were excluded from our

72 analysis due to lack of statistical significance. We compared our proposed model  
73 (Age+Time+Image) against models incorporating these error-related features to assess their  
74 contribution to predictive performance. As shown in **Table 3**, adding error-related features to our  
75 proposed model (Age+Time+Image) did not improve the AUC and sensitivity. We hypothesize  
76 this may be because the error information is already captured by ResNet. Furthermore, tracking  
77 these features requires documentation by trained professionals, which deviates from our goal of  
78 developing an accessible dementia screening tool and limits the applicability of our approach. Our  
79 findings indicate that the most effective model, combining age, completion time, and image  
80 probability derived from ResNet, captures the most relevant information for discriminating  
81 between individuals with and without dementia. This approach balances predictive power and ease  
82 of administration.

83 **Table A2:** Performance metrics for models using age only, age with completion time, age with  
 84 image probability score derived from ResNet, and age with both completion time and the  
 85 probability score generated by ResNet.

Dataset	Feature	AUC	Sensitivity	Specificity	F1 score	Accuracy
	Age	97.98 ± 0.78	78.46 ± 10.03	97.81 ± 0.66	97.21 ± 0.59	97.02 ± 0.67
	Age+Time	98.50 ± 0.59	81.54 ± 15.95	98.64 ± 0.48	98.00 ± 0.42	97.95 ± 0.35
	Age+Image (ResNet)	98.44 ± 0.82	92.31 ± 7.69	97.87 ± 1.09	97.86 ± 0.71	97.64 ± 0.90
LLFS	Age+Time+Image (ResNet)	98.62 ± 0.91	87.69 ± 8.77	98.26 ± 0.67	97.97 ± 0.58	97.83 ± 0.66
	Age	84.72 ± 10.68	46.67 ± 31.51	92.55 ± 14.64	94.20 ± 8.70	91.87 ± 13.96
	Age+Time	94.83 ± 2.30	65.00 ± 33.54	96.59 ± 3.99	97.07 ± 2.01	96.08 ± 3.47
	Age+Image (ResNet)	94.82 ± 3.13	73.33 ± 18.07	95.94 ± 5.02	96.87 ± 2.86	95.60 ± 4.73
FHS	Age+Time+Image (ResNet)	96.51 ± 2.37	85.00 ± 22.36	96.75 ± 3.40	97.47 ± 1.93	96.56 ± 3.12
	Age	95.19 ± 2.82	57.13 ± 8.71	98.99 ± 0.37	97.71 ± 0.30	97.76 ± 0.31
	Age+Time	97.56 ± 1.41	71.32 ± 15.51	98.78 ± 0.62	98.02 ± 0.39	97.97 ± 0.45
	Age+Image (ResNet)	97.26 ± 1.03	67.65 ± 13.15	98.74 ± 0.85	97.89 ± 0.50	97.83 ± 0.64
Combined	Age+Time+Image (ResNet)	97.95 ± 1.02	73.68 ± 7.50	98.81 ± 0.52	98.14 ± 0.31	98.08 ± 0.39

Time: completion time; Image(ResNet): probability score derived from ResNet.

86

87 **Table A3:** Performance metrics for models using age, the completion time, and the probability  
 88 score generated by ResNet and models using additional error-related features. This table compares  
 89 the performance of our proposed models with the models that use extra error-related features.

Dataset	Feature	AUC	Sensitivity	Specificity	F1 score	Accuracy
	Age+Time+Image (ResNet)	98.62 ± 0.91	87.69 ± 8.77	98.26 ± 0.67	97.97 ± 0.58	97.83 ± 0.66
	Age+Time+Image					
LLFS	(ResNet)+Error	98.37 ± 1.49	84.62 ± 7.70	98.36 ± 1.48	97.98 ± 1.08	97.81 ± 1.30
	Age+Time+Image (ResNet)	96.51 ± 2.37	85.00 ± 22.36	96.75 ± 3.40	97.47 ± 1.93	96.56 ± 3.12
	Age+Time+Image					
FHS	(ResNet)+Error	95.79 ± 2.52	68.33 ± 20.75	97.24 ± 2.74	97.52 ± 1.53	96.80 ± 2.38
	Age+Time+Image (ResNet)	97.95 ± 1.02	73.68 ± 7.50	98.81 ± 0.52	98.14 ± 0.31	98.08 ± 0.39
Combine	Age+Time+Image					
d	(ResNet)+Error	97.76 ± 1.18	71.54 ± 15.60	98.86 ± 0.79	98.12 ± 0.74	98.07 ± 0.79

Time: completion time; Image (ResNet): probability score derived from ResNet. Error: the number of times participants lifted their pen from the paper) and the presence of obvious tremor.

90

## 91 **Appendix E: Clinical relevance of sensitivity improvements**

92 We conducted statistical analysis of our models using the Wilcoxon signed-rank test. While the  
93 improvements did not reach conventional statistical significance, we argue that in clinical diagnostic  
94 contexts, practical significance can be more relevant than statistical analysis—particularly with  
95 limited sample sizes. We therefore present four complementary metrics specifically designed to  
96 quantify clinical impact:

$$97 \quad \text{Absolute Improvement} = \textit{Fusion} - \textit{Baseline} + \textit{Time},$$

$$98 \quad \text{Relative Improvement} = \frac{\textit{Fusion} - \textit{Baseline} + \textit{Time}}{\textit{Baseline} + \textit{Time}} \times 100\%,$$

$$99 \quad \text{Clinical Significance Ratio (CSR)} = \frac{\textit{Fusion} - \textit{Baseline} + \textit{Time}}{100\% - \textit{Baseline} + \textit{Time}},$$

$$100 \quad \text{Number Needed to Diagnose (NND)} = \frac{1}{\textit{Fusion} - \textit{Baseline} + \textit{Time}}.$$

101 The CSR measures the percentage of previously missed cases now correctly identified, while the  
102 NND quantifies the screening effort required to benefit one additional patient—metrics widely  
103 accepted in clinical research for assessing practical utility.

104 As shown in Table A4, our analysis reveals substantial clinical relevance across all datasets, with  
105 particularly notable results for the FHS cohort. The FHS dataset shows remarkable improvements  
106 with an absolute sensitivity increase of 30 percentage points (from 55.00% to 85.00%). This  
107 translates to a 54.55% relative improvement, with the fusion model correctly identifying 66.67%  
108 of cases that would have been missed by the baseline model (CSR). Most importantly, the NND  
109 of only 3.33 indicates that for approximately every 3-4 patients screened, one additional case will

110 be correctly identified that would otherwise be missed—an efficiency level considered highly  
111 significant in clinical practice.

112 **Table A4:** Sensitivity performance comparison and clinical impact metrics for Baseline+Time vs.  
113 Fusion models.

<b>Dataset</b>	<b>Baseline+Time</b>	<b>Fusion</b>	<b>Absolute Improvement</b>	<b>Relative Improvement</b>	<b>Clinical Significance Ratio</b>	<b>Number Needed to Diagnose</b>
LLFS	81.54%	87.69%	6.15%	7.54%	33.29%	16.26
FHS	55.00%	85.00%	30.00%	54.55%	66.67%	3.33
Combined	70.15%	73.68%	3.53%	5.03%	11.83%	28.33

The Fusion model incorporates age, completion time, and image probability (ResNet).

114

## 115 **Appendix F: Comparison of cohort characteristics**

116 The comparison between the FHS and LLFS cohorts reveals notable differences that help explain  
117 the performance variations observed in our models (Table A5). Despite similar average ages  
118 (70.25 years in FHS vs. 71.74 years in LLFS), the LLFS cohort contains a substantially higher  
119 proportion of dementia cases (4.03% compared to just 1.52% in FHS), providing more balanced  
120 training examples for our deep learning models. Educational backgrounds also differ considerably,  
121 with FHS participants demonstrating higher educational attainment (81.14% with college education  
122 or above compared to 68.03% in LLFS). Additionally, the LLFS cohort shows longer average Trail  
123 Making Test (TMT) completion times (44.36 seconds vs. 39.14 seconds in FHS), suggesting  
124 potential underlying differences in cognitive performance distributions.

125 The FHS primarily focused on cardiovascular health with clinic-based assessments, while LLFS  
126 specifically selected families demonstrating exceptional longevity and conducted in-home  
127 assessments. This design difference may have resulted in the inclusion of participants with  
128 mobility issues in LLFS who might also exhibit higher rates of cognitive impairment.

129 These cohort differences, particularly the lower prevalence of dementia cases in FHS (only 19  
130 cases compared to 65 in LLFS), create a more challenging classification task due to fewer positive  
131 examples available for training. This disparity helps explain the higher standard deviations and  
132 lower performance metrics observed in our FHS models.

133

134 **Table A5:** Overall cohort comparison between FHS and LLFS datasets. This table highlights key  
 135 demographic differences between FHS and LLFS cohorts that may influence model performance.

Feature	FHS (N = 1, 252)	LLFS (N = 1, 613)
Age (mean $\pm$ SD)	70.25 $\pm$ 12.42	71.74 $\pm$ 10.74
Gender		
Female (%)	57.76%	54.00%
Male (%)	42.24%	46.00%
Education		
High school or lower (%)	18.86%	31.97%
College or above (%)	81.14%	68.03%
Completion time (mean $\pm$ SD)	39.14 $\pm$ 25.20	44.36 $\pm$ 33.05
Dementia cases (%)	1.52%	4.03%

136

137 **Appendix G: Evaluation of fusion models across the cognitive impairment**  
138 **spectrum**

139 During our initial experimental design phase, we systematically explored model performance  
140 across four distinct classification tasks representing different clinical objectives in cognitive  
141 assessment (Table A6).

142 *Normal vs. MCI*: This task focused on detecting early cognitive decline by differentiating  
143 individuals with normal cognition from those with MCI. This represents the earliest stage of  
144 detection, aimed at identifying subtle cognitive changes before dementia develops.

145 *Normal/MCI vs. Dementia*: This task aimed to identify established dementia regardless of whether  
146 individuals were cognitively normal or had intermediate impairment. This classification approach  
147 aligns with the clinical need to identify patients requiring dementia-specific interventions and care.

148 *Normal vs. MCI/Dementia*: This task focused on detecting any cognitive impairment, grouping  
149 together individuals with either MCI or dementia. This classification mirrors screening approaches  
150 that aim to identify patients requiring further clinical evaluation.

151 *Normal vs. Dementia*: This task specifically differentiated individuals with normal cognition from  
152 those with dementia, excluding the intermediate MCI category.

153 To systematically evaluate these classification approaches, we analyzed performance using our  
154 fusion model (Age + Time + Image (ResNet)) on the LLFS dataset, as shown in Table A6. Analysis  
155 of these performance metrics reveals several important patterns. The *Normal vs. Dementia* task  
156 demonstrated consistently superior performance across all evaluation metrics, with the highest  
157 AUC (98.62%), F1 score (97.97%), and accuracy (97.83%). In contrast, tasks involving MCI

158 classification showed relatively lower performance, particularly in terms of sensitivity. The  
159 *Normal vs. MCI* task presented the greatest challenge, with substantially lower sensitivity  
160 (60.98%) compared to other tasks. This indicates that approximately 39% of MCI cases were  
161 misclassified as normal cognition using our fusion model. Similarly, the *Normal vs. MCI/Dementia*  
162 task showed lower sensitivity (69.16%), primarily due to the difficulty in correctly identifying MCI  
163 cases. The *Normal/MCI vs. Dementia* task performed better than the MCI-specific classifications  
164 but still showed lower sensitivity (80.00%) compared to the *Normal vs. Dementia* task (87.69%).  
165 This suggests that while our fusion model can effectively identify dementia, the inclusion of MCI  
166 cases introduces significant classification challenges.

167 These performance differences likely stem from several factors. MCI represents a heterogeneous  
168 intermediate stage with subtle cognitive changes that may not be as consistently reflected in TMT  
169 performance as the more pronounced deficits seen in dementia. Additionally, the boundaries  
170 between normal aging and MCI are less distinct than those between normal cognition and dementia,  
171 leading to greater classification uncertainty. MCI detection typically requires more extensive  
172 neuropsychological assessment or longitudinal monitoring in clinical practice, suggesting that  
173 single timepoint TMT data may have inherent limitations for identifying this intermediate stage.

174 **Table A6:** Performance of Fusion model (Age + Time + Image) across different classification tasks  
175 on LLFS dataset.

<b>Task</b>	<b>AUC</b>	<b>Sensitivity</b>	<b>Specificity</b>	<b>F1</b>	<b>Accuracy</b>
<i>Normal vs. MCI</i>	86.38	60.98	94.69	92.59	92.22
<i>Normal/MCI vs. Dementia</i>	97.39	80.00	97.95	97.48	97.29
<i>Normal vs. MCI/Dementia</i>	90.20	69.16	95.20	92.53	92.39
<i>Normal vs. Dementia</i>	98.62	87.69	98.26	97.97	97.83

176

## 177 **Appendix H: Oversampling effects**

178 Table A7 presents the results of our ablation study examining the effects of oversampling on our  
179 best-performing non-image model (Age+Time). The results demonstrate that oversampling  
180 significantly improves model sensitivity for detecting the minority class (dementia). For the  
181 Combined dataset, oversampling increases recall from 13.02% to 71.32%, while maintaining high  
182 specificity (98.78% vs 99.89%). Similarly, for the FHS dataset, recall improves from 0% to 65%,  
183 albeit with increased variability, and for the LLFS dataset, recall enhances from 24.62% to 81.54%.  
184 These findings underscore the critical role of oversampling in addressing the inherent class  
185 imbalance in our dementia detection task. Without oversampling, models exhibit strong bias toward  
186 the majority class (normal cognition). The observed trade-off between improved recall and  
187 marginally reduced specificity aligns with our clinical objective of prioritizing dementia case  
188 identification. These results validate our methodological choice of employing oversampling  
189 techniques to enhance model performance.

190 **Table A7:** Performance comparison of the best-performing non-image model with and without  
 191 oversampling (mean  $\pm$  std across 5-fold cross validation).

<b>Dataset</b>	<b>Oversampling</b>	<b>AUC</b>	<b>Sensitivity</b>	<b>Specificity</b>	<b>Accuracy</b>
Combined	No	97.35 $\pm$ 1.00	13.02 $\pm$ 6.35	99.89 $\pm$ 0.16	97.38 $\pm$ 0.23
	Yes	97.56 $\pm$ 1.41	71.32 $\pm$ 15.51	98.78 $\pm$ 0.62	97.97 $\pm$ 0.45
FHS	No	95.36 $\pm$ 2.50	0.00 $\pm$ 0.00	100.00 $\pm$ 0.00	98.48 $\pm$ 0.18
	Yes	94.83 $\pm$ 2.30	65.00 $\pm$ 33.54	96.59 $\pm$ 3.99	96.08 $\pm$ 3.47
LLFS	No	98.21 $\pm$ 0.98	24.62 $\pm$ 6.44	99.81 $\pm$ 0.28	96.84 $\pm$ 0.35
	Yes	98.50 $\pm$ 0.59	81.54 $\pm$ 15.95	98.64 $\pm$ 0.48	97.95 $\pm$ 0.35

192