

# AI-based differential diagnosis of dementia etiologies on multimodal data

Chonghua Xue<sup>1,2,\*</sup>, Sahana S. Kowshik<sup>1,3,\*</sup>, Diala Lteif<sup>1,4</sup>, Shreyas Puducheri<sup>1</sup>, Olivia T. Zhou<sup>1</sup>, Anika S. Walia<sup>1</sup>, Osman B. Guney<sup>1,2</sup>, J. Diana Zhang<sup>1,5</sup>, Serena T. Pham<sup>6</sup>, Artem Kaliev<sup>6</sup>, V. Carlota Andreu-Arasa<sup>6†</sup>, Brigid C. Dwyer<sup>7†</sup>, Chad W. Farris<sup>6†</sup>, Honglin Hao<sup>8†</sup>, Sachin Kedar<sup>9†</sup>, Asim Z. Mian<sup>6†</sup>, Daniel L. Murman<sup>10†</sup>, Sarah A. O’Shea<sup>11†</sup>, Aaron B. Paul<sup>12†</sup>, Saurabh Rohatgi<sup>12†</sup>, Marie-Helene Saint-Hilaire<sup>7†</sup>, Emmett A. Sartor<sup>7†</sup>, Bindu N. Setty<sup>6†</sup>, Juan E. Small<sup>13†</sup>, Arun Swaminathan<sup>14†</sup>, Olga Taraschenko<sup>10†</sup>, Jing Yuan<sup>8†</sup>, Yan Zhou<sup>8†</sup>, Shuhan Zhu<sup>15†</sup>, Cody Karjadi<sup>16</sup>, Ting Fang Alvin Ang<sup>16,17</sup>, Sarah A. Bargal<sup>18</sup>, Bryan A. Plummer<sup>4</sup>, Kathleen L. Poston<sup>19</sup>, Meysam Ahangaran<sup>1</sup>, Rhoda Au<sup>1,7,16,17,20,21</sup> & Vijaya B. Kolachalama<sup>1,3,4,20,‡</sup>

<sup>1</sup>*Department of Medicine, Boston University Chobanian & Avedisian School of Medicine, Boston, MA, USA*

<sup>2</sup>*Department of Electrical & Computer Engineering, Boston University, MA, USA*

<sup>3</sup>*Faculty of Computing & Data Sciences, Boston University, MA, USA*

<sup>4</sup>*Department of Computer Science, Boston University, MA, USA*

<sup>5</sup>*School of Chemistry, University of New South Wales, Sydney, Australia*

<sup>6</sup>*Department of Radiology, Boston University Chobanian & Avedisian School of Medicine, Boston, MA, USA*

<sup>7</sup>*Department of Neurology, Boston University Chobanian & Avedisian School of Medicine, Boston, MA, USA*

<sup>8</sup>*Department of Neurology, Peking Union Medical College Hospital, Chinese Academy of Medical Sciences, Beijing, China*

<sup>9</sup>*Departments of Neurology & Ophthalmology, Emory University School of Medicine, Atlanta, GA, USA*

<sup>10</sup>*Department of Neurological Sciences, University of Nebraska Medical Center, Omaha, NE, USA*

<sup>11</sup>*Department of Neurology, Columbia University Irving Medical Center, New York, NY, USA*

<sup>12</sup>*Department of Radiology, Massachusetts General Hospital, Boston, MA, USA*

<sup>13</sup>*Department of Radiology, Lahey Hospital & Medical Center, Burlington, MA, USA*

<sup>14</sup>*Department of Neurology, SSM Health, Madison, WI, USA*

<sup>15</sup>*Department of Neurology, Brigham & Women’s Hospital, Boston, MA, USA*

<sup>16</sup>*The Framingham Heart Study, Boston University Chobanian & Avedisian School of Medicine, Boston, MA, USA*

<sup>17</sup>*Department of Anatomy and Neurobiology, Boston University Chobanian & Avedisian School of Medicine, Boston, MA, USA*

<sup>18</sup>*Department of Computer Science, Georgetown University, Washington DC, USA*

<sup>19</sup>*Department of Neurology, Stanford University, Palo Alto, CA, USA*

<sup>20</sup>*Boston University Alzheimer’s Disease Research Center, Boston, MA, USA*

<sup>21</sup>*Department of Epidemiology, Boston University School of Public Health, Boston, MA, USA*

\* These authors contributed equally to this work

† Listed in alphabetical order

‡ Corresponding author: Vijaya B. Kolachalama, PhD; Email: [vkola@bu.edu](mailto:vkola@bu.edu); ORCID: <https://orcid.org/0000-0002-5312-8644>

## 1 Abstract

2 Differential diagnosis of dementia, with its overlapping symptomatology, remains a significant challenge in  
3 neurology. Here we present an algorithmic framework employing state-of-the-art techniques such as trans-  
4 formers as well as self-supervised frameworks and harnessing a broad array of data including demographics,  
5 person-level and family medical history, medication use, neuropsychological exams, functional evaluations,  
6 and multimodal neuroimaging to identify the etiologies contributing to dementia in individuals. The study  
7 utilized 9 independent, geographically diverse datasets, including the National Alzheimer’s Coordinating  
8 Center with 45,349 participants, the Alzheimer’s Disease Neuroimaging Initiative encompassing 1,821  
9 participants, and the Frontotemporal Lobar Degeneration Neuroimaging Initiative comprising 253 partic-  
10 ipants. Additionally, the Parkinson’s Progression Marker Initiative with 198 participants, the Australian  
11 Imaging, Biomarker and Lifestyle Flagship Study of Ageing cohort including 661 participants, the Open  
12 Access Series of Imaging Studies dataset with 491 participants, and the 4 Repeat Tauopathy Neuroimaging  
13 Initiative comprising 80 participants were used. The study also included two in-house datasets: one from the  
14 Lewy Body Dementia Center for Excellence at Stanford University with 182 participants, and another from  
15 the Framingham Heart Study including 1,651 individuals. Our model traverses the intricate spectrum of  
16 dementia by mirroring real-world clinical settings, aligning diagnoses with similar management strategies,  
17 and delivering robust predictions, even in the face of incomplete data. On the testing cohort, our model  
18 achieved a micro-averaged area under the receiver operating characteristic curve (AUROC) of 0.93, and  
19 a micro-averaged area under precision-recall curve (AUPR) of 0.87, in classifying individuals with normal  
20 cognition, mild cognitive impairment and dementia. Also, the micro-averaged AUROC was 0.95 and micro-  
21 averaged AUPR was 0.68 in differentiating 10 distinct dementia etiologies, defined through a consensus  
22 among a team of neurologists. One key strength lies in our model’s capability to address mixed dementias, a  
23 prevalent challenge in clinical practice, and the incorporation of interpretability techniques further unveiled  
24 vital disease-specific patterns. On a randomly selected subset ( $n = 100$ ), our model differentiated true  
25 positive and true negative cases across 12 out of 13 categories ( $p < 0.01$ ), as opposed to the neurologists’  
26 expertise in identifying 9 out of these 13 categories ( $p < 0.01$ ). Furthermore, the model’s correlations with  
27 different proteinopathies were substantiated through postmortem analyses. This included a significant asso-  
28 ciation with the global Alzheimer’s disease neuropathologic change (ADNC) score ( $p < 0.001$ ), and notable  
29 correlations with TDP-pathology, the presence of old microinfarcts, arteriosclerosis, and Prion disease (all  
30 with  $p < 0.05$ ). Our framework has the potential to be integrated as a screening tool for dementia in various  
31 clinical settings and drug trials, with promising implications for person-level management.

## 1 **Research in context**

2 *Systematic review:* Previous studies have demonstrated that models utilizing multimodal data can differenti-  
3 ate individuals across the dementia spectrum, identifying those with normal cognition (NC), mild cognitive  
4 impairment (MCI), and dementia (DE). Some studies have also ventured beyond this tripartite classification,  
5 aiming to differentiate Alzheimer’s disease (AD) from other forms of non-AD dementia. Majority of these  
6 investigations have approached the task as a binary classification, primarily focusing on the distinction be-  
7 tween AD and other dementia types. Also, limited studies have effectively tackled the intricate challenge of  
8 diagnosing mixed dementia, which is a common and complex issue encountered in clinical practice.

9 *Methods and findings:* Employing multimodal data from 9 distinct cohorts, encompassing 50,686 partici-  
10 pants, we developed an algorithmic framework that leverages transformers and self-supervised learning to  
11 facilitate differential dementia diagnoses. This model adeptly classifies individuals into 13 curated diagnos-  
12 tic categories, each tailored to reflect real-world clinical needs. These categories comprehensively cover the  
13 cognitive spectrum, ranging from NC, MCI to DE, and extend to 10 distinct dementia types. Our model  
14 demonstrates the capability to accurately diagnose dementia, even with incomplete data, and efficiently  
15 manage cases involving multiple co-occurring dementia conditions, a common occurrence in clinical prac-  
16 tice. It has shown commendable performance, surpassing expert clinical assessments, and its predictions  
17 have been corroborated by postmortem data, particularly in relation to various proteinopathies.

18 *Interpretation:* Our work provides a robust and adaptable framework for comprehensive dementia screening  
19 for drug trials and in various clinical settings, ranging from primary care to memory clinics.

1 Dementia is one of the most pressing health challenges of our time. With nearly 10 million new cases  
2 reported annually, this syndrome, characterized by a progressive decline in cognitive function severe enough  
3 to impede daily life activities, continues to present considerable clinical and socio-economic challenges. In  
4 2017, the World Health Organization's global action plan on the public health response to dementia high-  
5 lighted the prompt and precise diagnosis of dementia as a pivotal strategic objective.<sup>1</sup> As such, diagnostic  
6 precision in the varied landscape of dementia remains an unmet need, even as the demands for such pre-  
7 cision escalate with the aging global population and the imperative for better participant screening in drug  
8 trials. The clinical presentation of different dementia forms often overlaps, further compounded by the het-  
9 erogeneity in findings on magnetic resonance imaging (MRI) scans. The necessity for improvements in the  
10 field becomes ever more pressing considering the projected shortage of specialists including neurologists,  
11 neuropsychologists and geriatric care providers,<sup>2-4</sup> emphasizing the urgency to innovate and evolve our  
12 diagnostic tools. While Alzheimer's disease (AD) is a leading cause, other forms such as vascular demen-  
13 tia (VD), Lewy body dementia (LBD), and frontotemporal dementia (FTD) are also prevalent, and often  
14 co-occur in most individuals. The significant intersection of symptoms among these disorders and other  
15 dementia etiologies, amplified by the varying intensity of symptom manifestation, frequently complicates  
16 the process of differential diagnosis.<sup>5</sup>

17 The imperative for enhanced diagnostic accuracy in AD and related neurodegenerative disorders is  
18 becoming increasingly critical in the context of significant advancements in medical diagnostics. While  
19 recent regulatory approvals have facilitated the transition of cerebrospinal fluid (CSF) biomarkers and ad-  
20 vanced imaging techniques like positron emission tomography (PET) from research environments to broader  
21 clinical applications, the development and potential clinical integration of blood-based biomarkers remain  
22 an area of active research.<sup>6,7</sup> However, accessibility to these diagnostic tools remains constrained, not only  
23 in remote and economically developing regions but also in urban healthcare centers, as exemplified by pro-  
24 longed waiting periods for specialist consultations. This challenge is compounded by a global shortage of  
25 specialists, such as behavioral neurologists and neuropsychologists, leading to overreliance on cognitive  
26 assessments that may not be culturally appropriate due to the lack of formal training programs in neuropsy-  
27 chology in many parts of the world. The recent advent of disease-modifying therapies marks a significant  
28 shift in the treatment landscape for AD,<sup>8,9</sup> further highlighting the necessity of early and accurate risk strat-  
29 ification in both primary care and general neurology settings. While conventional methods like clinical  
30 evaluations, neuropsychological testing, and MRI remain central to antemortem differential dementia di-  
31 agnosis, their effectiveness heavily relies on the diminishing pool of specialist clinicians. This situation  
32 underscores an urgent need for the healthcare system, particularly in primary care, to evolve and adapt to  
33 the rapidly changing dynamics of dementia diagnosis and treatment.

34 Machine learning holds promising potential for enhancing the accuracy and efficiency of dementia  
35 diagnosis.<sup>10-12</sup> While most previous research leveraging these advanced frameworks has concentrated on  
36 neuroimaging data,<sup>13,14</sup> few have ventured to combine imaging with non-imaging data,<sup>15</sup> such as demo-  
37 graphics, medical histories, and neuropsychological assessments to distinguish cognitively normal (NC)  
38 individuals from those with mild cognitive impairment (MCI) and dementia. AD being the most preva-  
39 lent type of dementia, has naturally been the primary focus of much of this research. Few studies have  
40 also attempted to discern neuroimaging signatures unique to AD by contrasting them with other dementia  
41 etiologies.<sup>16-22</sup> Recently, we proposed a nuanced approach to stratify individuals based on cognitive status  
42 and discern likely AD cases from other non-AD dementia types.<sup>23</sup> These investigations, coupled with other  
43 studies,<sup>24-26</sup> have begun to illuminate the complex matrix of factors contributing to dementia. Moreover,  
44 the challenge of differential diagnosis in dementia is further compounded by the inherent limitations in ex-  
45 isting dementia risk scoring systems.<sup>27-29</sup> Research indicates that individualized assessments based on these

46 generic scores frequently lead to significant error rates, thus diminishing their effectiveness in accurately  
47 identifying individuals suitable for targeted dementia prevention strategies.<sup>30</sup> These observations highlight  
48 the critical need for innovative solutions in differential dementia diagnosis, emphasizing the integration of  
49 diverse data sets to surpass the limitations of traditional risk assessment methods.

50 In this study, we propose a multimodal machine learning framework, harnessing a diverse array of  
51 data, including demographics, personal and family medical history, medication use, neuropsychological  
52 tests, functional evaluations, and multimodal neuroimaging to perform differential dementia diagnosis. The  
53 model incorporates state-of-the-art techniques such as transformers and self-supervised learning, enabling  
54 it to navigate the spectrum of dementia conditions with improved performance. Our model for differen-  
55 tial dementia diagnosis reflects real-world scenarios, with diagnostic categories designed for clinical rele-  
56 vance, aligning diagnoses with similar management strategies to aid neurologists and other practitioners in  
57 screening and treatment planning. The model's robustness is demonstrated through rigorous validation on  
58 independent, geographically diverse datasets, achieving parity with expert clinical diagnoses in comparative  
59 analysis. By employing advanced interpretability methods, our model elucidates disease-specific patterns  
60 critical for differential diagnosis, contributing significantly to our understanding of heterogeneous dementia  
61 phenotypes. The fidelity of these patterns was confirmed by postmortem data, underscoring our model's  
62 ability to dissect the intricate pathophysiology of dementia. Our algorithmic framework opens new avenues  
63 for dementia screening in various clinical settings, with significant implications for person-level manage-  
64 ment. This study underscores the potential of AI-driven tools in healthcare, paving the way for improved  
65 diagnostic accuracy, efficient resource utilization, and better outcomes.

## 1 Results

### Glossary 1

Acronym	Description
NC	Normal cognition
MCI	Mild cognitive impairment
DE	Dementia
AD	Alzheimer's disease including Down syndrome
LBD	Lewy body dementia including dementia with Lewy bodies and Parkinson's disease dementia
VD	Vascular dementia, vascular brain injury, and vascular dementia including stroke
PRD	Prion disease including Creutzfeldt-Jakob disease
FTD	Frontotemporal lobar degeneration and its variants, including primary progressive aphasia, corticobasal degeneration and progressive supranuclear palsy, and with or without amyotrophic lateral sclerosis
NPH	Normal pressure hydrocephalus
SEF	Systemic and environmental factors including infectious diseases (HIV included), metabolic, substance abuse / alcohol, medications, systemic disease, and delirium
PSY	Psychiatric conditions including schizophrenia, depression, bipolar disorder, anxiety, and post-traumatic stress disorder
TBI	Moderate/severe traumatic brain injury, repetitive head injury, and chronic traumatic encephalopathy
ODE	Other dementia conditions including neoplasms, multiple systems atrophy, Huntington's disease, seizures, etc.

2

3 Leveraging the power of routinely gathered clinical data, our model provides a nuanced approach to  
4 differential dementia diagnosis. This framework assigns individuals to one or more of the thirteen diagnostic  
5 categories (refer to Glossary 1), which were meticulously defined through consensus among a team of ex-  
6 pert neurologists. This practical and intuitive categorization is designed with clinical management pathways  
7 in mind, thereby echoing real-world scenarios. For instance, we have grouped dementia with Lewy bodies  
8 and Parkinson's disease dementia under the comprehensive category of Lewy body dementia (LBD). This  
9 classification stems from an understanding that the care for these conditions often follows a similar path,  
10 typically overseen by a multidisciplinary team of movement disorder specialists. In the context of vascular  
11 dementia (VD), we included individuals who exhibited symptoms of a stroke, possible or probable VD, or  
12 vascular brain injury. This encompassed cases with symptomatic stroke, cystic infarct in cognitive networks,  
13 extensive white matter hyperintensity, and/or executive dysfunction, where these conditions were identified  
14 as the primary contributors to the observed cognitive impairment. The inclusion criteria were based on  
15 the expectation that such patients would typically receive care from clinicians specializing in stroke and  
16 vascular diseases. Likewise, we have considered various psychiatric conditions - such as schizophrenia,  
17 depression, bipolar disorders, anxiety, and post-traumatic stress disorder - under one category (PSY), ac-  
18 knowledging that their management predominantly falls within the realm of psychiatric care providers. By  
19 aligning diagnostic categories with practical clinical care pathways, our model serves not only to demystify

20 an individual's condition but also to direct efficient and appropriate clinical management strategies. The  
21 result is a scientifically robust tool that is also intuitively aligned with the real-world scenarios encountered  
22 in clinical care.

23 Our model's development and testing leveraged a robust variety of data, encompassing thousands of  
24 individuals from population-level studies, community, clinical as well as research cohorts (Table 1, Table S1  
25 & Fig. 1). This broad, diverse set of data serves as a testament to our study's core strength, and the ensuing  
26 results speak to the model's strong accuracy and generalizability. The model accurately classified cognitive  
27 status across the full spectrum of clinical diagnosis tasks. Across three categories - normal cognition (NC),  
28 mild cognitive impairment (MCI), and dementia (DE) - our model achieved high performance as shown via  
29 accuracy (Acc.), sensitivity (Se.), specificity (Sp.), Matthews correlation coefficient (MCC), as well as the  
30 area under the receiver operating characteristic (AUC) and area under precision-recall (AP) curves (Fig. 2a).  
31 Further, the model delivered striking results across AD, LBD, VD, PRD, FTD, NPH, SEF, PSY, TBI, and  
32 ODE, reaffirming the model's high performance and reliability in diagnosing various forms of dementia.  
33 With these compelling metrics, we present robust empirical evidence of our model's ability to tackle differ-  
34 ential dementia diagnosis efficiently and effectively. Also, the utility of our model is notably demonstrated  
35 in its robustness and adaptability to various datasets acquired via different protocols at multiple institutions  
36 (Fig. 2a). An example of this is found in our tests on data from the National Alzheimer's Coordinating  
37 Center (NACC) cohort that was not used during model training. The model continued to exhibit its robust  
38 performance on external cohorts, including the Alzheimer's Disease Neuroimaging Initiative (ADNI), and  
39 the Framingham Heart Study (FHS), substantiating its consistency across diverse datasets. Despite the fact  
40 that the ADNI and FHS datasets comprised of a limited set of factors (Table S2), the model effectively  
41 accomplished all classification tasks. The model's ability to deliver reliable predictions with constrained  
42 data sets suggests that its performance, already robust under conditions of incomplete information, could be  
43 further optimized with access to more comprehensive datasets. To further evaluate the model's resilience  
44 to incomplete data, we artificially introduced varying levels of data missingness in the NACC cohort and  
45 assessed the impact on its predictive performance. By selectively removing portions of the data from the  
46 NACC cohort, we aimed to test our model's predictive performance under various constraints. As depicted  
47 in the chord diagram (Fig. 2b), even when confronted with missing data elements — whether it be MRI  
48 results, UPDRS, GDS, NPI-Q, FAQ, NP tests or other parameters — our model consistently produced high,  
49 reliable scores. This reinforces not only its predictive stability but also its potential applicability in diverse  
50 clinical scenarios where complete datasets might be unattainable.

51 We applied two-dimensional t-distributed stochastic neighbor embedding (tSNE) to evaluate the dis-  
52 tribution of our model's predictions within the cognitive spectrum encompassing NC, MCI and DE as de-  
53 picted in Fig. 2c. The visualization distinctly segregates the cognitive states using color codes: a densely  
54 populated blue pathway denotes NC, indicating a homogeneous assemblage of individuals with normal cog-  
55 nitive health. The points coded in yellow for MCI exhibit a sparser distribution, symbolizing a transitional  
56 and more varied cohort positioned between NC and DE. The pink stream earmarks DE, forming a dense  
57 aggregation that likely corresponds to individuals diagnosed with dementia. This graphical representation  
58 proves critical for distilling complex, multi-dimensional datasets into a comprehensible two-dimensional  
59 space, thus providing a graphical synopsis of how individual-level multimodal data processed via our mod-  
60 eling framework can demarcate these cognitive conditions.

61 Shapley analysis<sup>31</sup> was employed to determine feature importance on the NACC test set. The process  
62 began by categorizing cases according to their labels, focusing exclusively on those with correct predictions

63 for Shapley value calculations. This approach yielded feature-specific Shapley values for the chosen cases,  
64 illuminating each feature's role in influencing the model's final decision on the NC, MCI and DE cases, re-  
65 spectively (Figs. 2d-f). These features were arranged based on their mean impact, quantified by the average  
66 magnitude of their respective Shapley values. Additionally, the computation of Shapley values incorporated  
67 considerations for data missingness; Shapley values of the features that were unavailable were assigned a  
68 constant value of zero. The Shapley value distribution for each feature indicates the magnitude and direc-  
69 tion of the feature's impact on the model's prediction, with higher absolute Shapley values signifying greater  
70 influence. For NC, key features included cognitive status based on neuropsychological exam, MoCA and  
71 MMSE scores, memory-related tasks (like naming vegetables and animals), and daily functionality ques-  
72 tions. For MCI, similar cognitive and memory-related features were found to be impacting, in addition to  
73 questions about daily living activities and independence levels. For DE, the model placed high importance  
74 on cognitive status based on neuropsychological exams, the difficulty with daily tasks, and medication usage  
75 for AD symptoms. Across all conditions, the plots show that certain features (e.g., cognitive status based  
76 on neuropsychological exam) consistently have a higher impact on the model's predictions, as indicated by  
77 larger Shapley values. The distribution of Shapley values for each feature — spread across a spectrum from  
78 high to low impact — provides insights into the heterogeneity of feature influence on the prediction of each  
79 cognitive state. The consistency of cognitive assessments in predicting across all cognitive states suggests  
80 their fundamental role in the model's decision-making process. Conversely, the variability in the impact of  
81 daily functionality and memory-related tasks underscores the nuanced differences in their relevance to each  
82 cognitive state. Overall, the Shapley value plots offered a detailed and quantifiable visualization of how each  
83 feature contributes to the model's predictions, which could be crucial for understanding and improving the  
84 model's interpretability and accuracy in a clinical setting.

85 The receiver operating characteristic (ROC) and precision-recall (PR) curves reflected strong model  
86 performance across different averaging methods (Figs. 2g-h). On the test cohort, comprising the NACC  
87 dataset (unused in training) along with ADNI and FHS data, our model demonstrated strong classification  
88 abilities for NC, MCI, and DE, achieving a micro-averaged AUROC of 0.93 and a micro-averaged AUPR  
89 of 0.87. Additionally, the macro-averaged metrics showed an AUROC of 0.91 and a AUPR value of 0.80.  
90 The weighted-average AUROC and AUPR values further underscored the model's efficacy, standing at 0.92  
91 and 0.84, respectively. Of note, the micro-average approach consolidates true positives, true negatives, false  
92 positives, and false negatives from all classes into a unified curve, providing a global performance metric.  
93 In contrast, the macro-average calculates individual ROC/PR curves for each class before computing their  
94 unweighted mean, disregarding potential class imbalances. The weighted-average, while similar in approach  
95 to macro-averaging, assigns a weight to each class's ROC/PR curve proportionate to its representation in the  
96 dataset, thereby acknowledging class prevalence. The slight variations observed in the AUROC and AUPR  
97 values between the averaging methods might reflect the influence of class distribution and prevalence in  
98 the performance metric calculations. Overall, the micro-averaging method performed slightly better, which  
99 could be due to the fact that it takes the class imbalance into account by weighting the performance by the  
100 number of samples in each class.

101 We conducted a comparison between the model's predicted probability scores and the clinical de-  
102 mentia ratings (CDR) available for all participants in the NACC testing, and ADNI cohorts (Figs. 2i & 2j).  
103 Remarkably, despite not incorporating CDR scores as input during model training, our predictions exhibited  
104 a strong correlation with CDR scores. In the analysis of the NACC dataset, the model's predictions exhib-  
105 ited increasing heterogeneity as a function of ascending Clinical Dementia Rating (CDR), with statistically  
106 significant divergences manifest across the spectrum of cognitive impairment ( $p < 0.0001$ ). Notably, this  
107 pattern did not hold between CDR scores of 2.0 and 3.0, where no significant statistical difference was dis-



108 cerned. The ADNI dataset revealed a statistically significant demarcation ( $p < 0.0001$ ) in model-predicted  
109 probabilities between the baseline CDR rating and higher gradations. This underscores the model's sensi-  
110 tivity to incremental exacerbations in clinical dementia assessments. For the FHS dataset (Fig. 2k), which  
111 substitutes a consensus panel's diagnostic categorization (normal, impaired, and dementia) for CDR scores,  
112 a marked statistical significance ( $p < 0.0001$ ) was evident in the model's predicted probabilities across these  
113 diagnostic strata, with the exception of the distinction between normal cognition and impairment. This sug-  
114 gests a nuanced challenge for the model in discriminating between the initial stages of cognitive decline.  
115 Collectively, these findings illuminate the model's robust capacity to delineate differential cognitive states,  
116 showcasing its potential as a discerning tool for identifying levels of cognitive impairment across diverse  
117 clinical datasets.

118 Our methodology for conducting differential diagnosis led to the simultaneous prediction of probabili-  
119 ties associated with multiple potential causes of dementia, as visually represented in Fig. 3. To delve deeper  
120 into the intricate nature of real-world clinical scenarios, where persons frequently exhibit symptoms suggest-  
121 ing the coexistence of multiple dementia types, we created a comprehensive visualization,<sup>32</sup> which depicts  
122 probabilities that reflect co-occurring conditions. These plots provided a valuable tool for understanding the  
123 complex landscape of dementia diagnoses, where a substantial number of cases involve a combination of  
124 contributing factors. The box-and-whisker plot on the joint probabilities is particularly informative in illus-  
125 trating the range and variability of the model's predictions. For instance, the uppermost row corresponds  
126 to non-dementia cases. Here, the model's predictions are expressed as joint probabilities — the product of  
127 individual probabilities for each condition, symbolized as  $p(AD = 0) \times p(PSY = 0) \times p(ODE = 0) \times$   
128  $p(SEF = 0) \times p(LBD = 0) \times p(VD = 0) \times p(PRD = 0) \times p(FTD = 0) \times p(NPH = 0) \times p(TBI = 0)$ .  
129 The bottom row highlights cases that include the combination of AD, PSY, ODE and SEF, where the model's  
130 predictions are expressed as  $p(AD = 1) \times p(PSY = 1) \times p(ODE = 1) \times p(SEF = 1) \times p(LBD =$   
131  $0) \times p(VD = 0) \times p(PRD = 0) \times p(FTD = 0) \times p(NPH = 0) \times p(TBI = 0)$ . This level of quantification  
132 underscores the model's capacity to gauge the intricacies of dementia presentations that are multifactorial  
133 in nature, highlighting its utility in complex clinical assessments. Taken together, these findings demon-  
134 strate that our model effectively identified instances where multiple etiologies were at play, reflecting its  
135 ability to navigate the challenging task of recognizing mixed causes of dementia. This accomplishment un-  
136 derscores the significant contribution of our research, particularly in addressing the diagnostic complexities  
137 encountered in clinical practice when individuals present with overlapping symptoms indicative of multiple  
138 dementia types.

139 Similar to the Shapley value illustrations shown in Figs. 2d-f, we observed that Shapley value anal-  
140 yses revealed distinct feature importance across various dementia types (Figs. 4a-j). For instance, features  
141 such as 'cognitive status based on neuropsychological exams', 'level of independence', and 'difficulty in  
142 performing daily activities' were influential across several dementia categories. Such consistency of certain  
143 features across various dementia types underscores their potential as universal markers of neurodegenera-  
144 tion. Moreover, the analysis identified several features that are emblematic of well-known etiologies un-  
145 derlying dementia, as well as factors that are widely recognized as influential in driving the progression of  
146 specific diseases. In predicting AD, factors like the presence of ApoE alleles and the use of FDA-approved  
147 medication for AD symptoms emerged as influential. For LBD, a prior diagnosis of Parkinson's disease, the  
148 presence of gait disorders, and speech difficulties were among the key predictors. VD prediction was driven  
149 by the history of stroke, the impact of cerebrovascular disease on cognitive impairment, and the occurrence  
150 of multiple infarcts. Presence of depression or dysphoria in the last month as well as GDS score were signif-  
151 icant for predicting psychiatric conditions (PSY). These plots collectively demonstrate the model's reliance  
152 on a mix of features, reflecting the multifaceted nature of dementia disorders and the complexity of their

153 prediction. Additionally, the ROC and PR curves reflected strong model performance on the model's overall  
154 assessment on identifying dementia etiologies across different averaging methods (Figs. 4k-l). Our model  
155 attained impressive results, with micro-averaged values of AUROC and AUPR at 0.95 and 0.68, respectively.  
156 In macro-averaged terms, the AUROC and AUPR stood at 0.91 and 0.34. Moreover, the weighted-average  
157 values for AUROC and AUPR were 0.94 and 0.72, respectively.

158 Etiology-specific model probability scores revealed significant correlations with neuropathological  
159 evidence found in common dementia types. The composite violin and box plots in Fig. 5 illustrate the  
160 distribution and probability distributions and median tendencies for each cohort indicating that with in-  
161 creasing neuropathological severity, there is a corresponding elevation in the likelihood of neurodegener-  
162 ation according to the model. The first three plots (Figs. 5a-c) compare baseline stages of Thal phase for  
163  $A\beta$  plaques (Thal), Braak stage for neurofibrillary degeneration (Braak) and density of neocortical neu-  
164 ritic plaques (CERAD) against progressive Thal, Braak and CERAD stages (A1-A3, B1-B3 and C1-C3,  
165 respectively). Each demonstrated an upward shift in the median probability of AD and an expansion of  
166 the interquartile range as the stages advance, with statistical significance ( $p < 0.001$  for Thal stage and  
167  $p < 0.0001$  for Braak and CERAD stages, respectively). Also, we rejected the null hypothesis of there be-  
168 ing no significant differences in model-predicted AD probabilities between semi-quantitative scores of Thal,  
169 Braak and CERAD scores (Fig. S2). Furthermore, by contrasting individuals without AD against those with  
170 varying degrees of NIA-AA Alzheimer's disease neuropathologic change (ADNC), which unifies the Thal,  
171 Braak and CERAD scores (Fig. 5d), we observed a similar shift towards higher AD probabilities in the latter  
172 group ( $p < 0.001$ ). Collectively, these plots illustrate a clear trend where advancing stages of AD-related  
173 neurodegeneration are associated with increased probabilities of AD. Finally, we rejected the null hypothe-  
174 sis of there being no significant differences in FTD probability between the presence or absence of TDP-43  
175 pathology ( $p < 0.05$ ), VD probability between the presence or absence of old microinfarcts ( $p < 0.05$ ) and  
176 arteriosclerosis ( $p < 0.05$ ), as well as PRD probability between the presence or absence of Prion disease  
177 ( $p < 0.05$ )(Figs. 5e-h). The results are consistent with the well-documented association between TDP-43  
178 protein aggregation and its prevalence in FTD as well as other neurodegenerative diseases.<sup>33,34</sup> Additionally,  
179 the clear linkage between cerebrovascular pathologies and the incidence of VD is reinforced by our data.  
180 Crucially, these outcomes highlight the capability of our AI-driven framework to align model-generated  
181 probability scores with a range of neuropathological states beyond AD, supporting its potential utility in  
182 broader neurodegenerative disease research.

183 The incorporation of independent confidence scores for diagnostic tasks improved the model's inter-  
184 pretability, and facilitated comparison of model performance with the clinicians. Neurologists reviewed 100  
185 randomly selected cases, including various dementia subtypes, with comprehensive data including demo-  
186 graphics, medical history, neuropsychological tests, and multi-sequence MRI scans. The Shapiro-Wilk test  
187 revealed non-normal distributions in the confidence scores, both from the model and the experts. Conse-  
188 quently, the Brunner-Munzel test was applied to compare the differences between these sets of confidence  
189 scores. Pearson correlation analysis was used to measure the interrater reliability among the confidence  
190 scores assigned by different evaluators. Notably, in instances where the diagnosis was confirmed (true pos-  
191 itives), the neurologists' confidence scores across categories such as NC, MCI, DE, AD, LBD, VD, FTD,  
192 NPH, and PSY were statistically significant in comparison to cases deemed non-diagnostic (true negatives)  
193 ( $p < 0.01$ ) (Fig. 6a). Similarly, the model's probabilities in true positive cases across an extended range of  
194 conditions — including those aforementioned plus PRD, SEF, and TBI — were statistically significant when  
195 contrasted with true negatives ( $p < 0.01$ ). However, for certain etiologies, including PRD, SEF, TBI, and  
196 ODE, there was no apparent statistical distinction in the confidence ratings provided by neurologists between  
197 true positive and true negative cases. This absence of statistical significance was also observed in the model's

198 predicted probabilities for ODE between true positive and true negative cases. Several potential explanations may elucidate this inconsistency. Despite providing neurologists with comprehensive information, 199 including demographics, medical history, medication usage, neuropsychological test results, functional assessments, and multimodal neuroimaging data for all cases, this information might have been insufficient to 200 confirm the impact of these etiologies. Additionally, the model's limitations could be attributed to a smaller 201 training dataset specific to certain etiologies and a potential deficiency in the breadth of features needed for 202 the identification of etiology-specific contributions, especially concerning ODE. Pearson correlation coefficients 203 offered a comprehensive overview of the interrelationships among the neurologists and the model 204 across a spectrum of neurological conditions (Fig. 6b). NC was distinguished by the highest correlation coefficient 205 between the model-predicted probabilities and mean neurologist confidence scores ( $0.92 \pm 0.02$ ), 206 indicating robust, consistent associations in this group. DE showed notable correlations ( $0.89 \pm 0.02$ ). In 207 contrast, MCI, AD, LBD, VD, FTD and PSY exhibited modest correlations, highlighting the potential heterogeneity 208 within these disorders. The lower correlations observed on certain etiologies (PRD, NPH, SEF, 209 TBI and ODE) underscore the diverse and complex nature of these conditions compounded by the lack of 210 extensive features to tease out their unique signatures. These observations underscore the model's potential 211 in complementing neurologist expertise, yet also highlights the complexities and limitations in accurately 212 diagnosing more diverse and less represented conditions. 213 214

215 In a separate assessment, neuroradiologists evaluated a randomly selected set of 70 clinically diagnosed 216 dementia cases, concentrating on MRI findings and demographic information. They provided ratings 217 for atrophy and pathological changes within specific brain regions, including the temporal lobes, frontal 218 lobe, insula, limbic systems, fusiform gyrus, and overall brain considerations (Fig. S3). The calculated 219 pairwise Pearson correlation coefficients, representing interrater reliability among seven neuroradiologists, 220 revealed a moderate overall agreement with a mean coefficient of ( $0.39 \pm 0.02$ ). Radiologists were mostly 221 in agreement on the whole brain-level assessments related to the presence of hyperintensities ( $0.67 \pm 0.03$ ) 222 and identifying prior infarcts ( $0.59 \pm 0.08$ ). Within the temporal lobe, the highest concordance was found 223 in ratings of the anterior temporal lobe ( $0.68 \pm 0.09$ ), while the posterior temporal lobe showed the least 224 agreement ( $0.18 \pm 0.17$ ). An interesting pattern emerged, indicating a trend towards greater agreement in 225 assessments of the right-sided brain regions compared to the left-sided ones. We stratified cases based on 226 expert consensus on atrophy and pathological changes and compared the distributions of etiology-specific 227 model probability scores between the groups (Fig. S4). We found that the presence of infarcts, widely recognized 228 as a diagnostic marker for VD,<sup>35</sup> was associated with an increase in model probabilities of VD. 229 The presence of infarcts was also associated with elevated model probabilities for TBI. This association can 230 be attributed to common secondary injuries in TBI, such as cranial hypoperfusion, subsequent ischemia, 231 post-traumatic cerebral infarctions, and encephalomalacia.<sup>36-40</sup> Furthermore, atrophy in the anterior cingulate 232 gyrus was linked with higher model probabilities for FTD, supporting its recognized predictive value 233 for behavioral variant FTD.<sup>41-43</sup> Overall, these observations suggest the model's ability to mirror expert radiological 234 evaluations, particularly in recognizing key brain changes indicative of different dementia types.

## 1 Discussion

2 The key contributions of our study are highlighted as follows. By leveraging a transformer architecture as  
3 the backbone and using principles of self-supervised learning, we developed an algorithmic framework that  
4 can process flexible combinations of multimodal data and perform differential diagnosis of dementia. This  
5 approach allows it to effectively navigate the spectrum of cognitive states, assigning probability scores to  
6 normal cognition (NC), mild cognitive impairment (MCI), and dementia (DE), as well as to the specific eti-  
7 ologies underlying DE. Unlike our previous work,<sup>15,23</sup> a unique capability of our model is its adeptness in  
8 quantifying co-existing dementia conditions within individuals. The model's robustness was further estab-  
9 lished through its training and validation across a diverse set of independent cohorts, demonstrating its broad  
10 applicability. Additionally, our model can identify key features associated with various cognitive states and  
11 dementia etiologies, providing insights that align with established pathologies. In a comparative evaluation  
12 on a randomly selected case subset, our model's proficiency in distinguishing true positive and true neg-  
13 ative cases notably surpassed that of neurologist-level assessments. These results underscore our model's  
14 significant potential in enhancing the precision and efficacy of diagnosing dementia-related disorders.

15 Our modeling framework encompasses a diverse set of elements, each with unique strengths. For  
16 instance, the backbone transformer architecture, serving as the foundation, excels at learning long-range de-  
17 pendencies in data, a crucial asset for tasks like medical diagnosis and prediction. Furthermore, our model  
18 integrates feature-specific embedding modules, facilitating the absorption of richer representations from  
19 various data types, including nominal, ordinal, numerical features, and complex imaging data. Notably,  
20 our framework's capability to handle diverse MRI sequences is an advantage. SwinUNETR, a fusion of U-  
21 shaped network design with Swin transformer encoder and CNN-based decoder connected through skip con-  
22 nections, efficiently generates image embeddings. These embeddings, when integrated into our backbone  
23 transformer, empower the processing of multimodal data for differential dementia diagnosis. Moreover,  
24 the backbone transformer framework has the potential for extension to incorporate text data from electronic  
25 health records and emerging digital data types, such as voice recordings as well as sensor data from wearable  
26 technologies.

27 Accommodating the intricacies of dementia, our model assigns distinct probabilities to each potential  
28 diagnosis. We note this as a key contribution in our work because these values present a nuanced frame-  
29 work for clinicians to create a rank-ordered list of plausible etiologies. In essence, the model processes  
30 available multimodal data and outputs probabilities that outline the person's primary condition and account  
31 for the multifactorial and often overlapping nature of dementias. The visualizations that we generated were  
32 instrumental in quantifying the variability and distribution of the model's predictions, facilitating an un-  
33 derstanding of the likelihood of various diagnostic intersections. Effectively, our model acknowledges and  
34 characterizes the complexity of mixed dementias, a clinical scenario frequently encountered in practice,  
35 and often confirmed via postmortem evidence.<sup>44-46</sup> Such capability fosters a detailed understanding of the  
36 person's condition, but also operationalizes the process of differential diagnosis and could facilitate clinical  
37 uptake of software-based assistive tools.

38 The utility of our modeling framework is founded on its robust processing of diverse input types and  
39 its adept handling of incomplete datasets, properties that are essential for clinicians requiring immediate and  
40 accurate diagnostic information in environments with variable data availability. For example, when a gen-  
41 eral practitioner records clinical observations and cognitive test results for an elderly person with possible  
42 cognitive decline, our model can calculate a probability score indicative of MCI or DE. This function facili-

43 tates early medical intervention and more informed decisions regarding specialist referrals. At a specialized  
44 memory clinic, the addition of extensive neuroimaging data and in-depth neuropsychological battery to the  
45 model may increase the precision of the diagnosis, which, in turn, enhances the formulation of individual  
46 management strategies with a revised probability score. Such capacity to tailor its output to the scope of  
47 input data exemplifies our modeling framework’s role in different healthcare settings, including those where  
48 swift and resource-efficient diagnosis is paramount. The generation of specific, quantifiable probability  
49 scores by the model augments its utility, establishing it as a useful component in the healthcare delivery pro-  
50 cess. Displaying diagnostic accuracy using varied training data — ranging from demographic information  
51 to clinical signs, neuroimaging findings, and neurological test results — the model’s versatility facilitates its  
52 adaptation to varied clinical operations without necessitating a fundamental overhaul of existing workflows.  
53 Consequently, our model fosters a seamless transition across the different levels of dementia care, enabling  
54 general practitioners to perform preliminary cognitive screenings and specialists to conduct thorough exami-  
55 nations. Its inclusive functionality assures an accessible and comprehensive tool ensuring fail-safe operation  
56 in early detection, continuous monitoring, and the fine-tuning of differential diagnoses, thereby elevating the  
57 standard of dementia care.

58 Shapley value analysis served as an indispensable interpretive mechanism in our study, elucidating the  
59 specific variables’ influence on our deep learning model’s predictions and bridging computational forecasts  
60 with clinical implications. By quantifying the influence of individual factors on the outcomes for NC,  
61 MCI, and DE as well as for the dementia etiologies, Shapley values not only bolster the transparency of  
62 our models but also reinforce the validity of our approach by aligning with established medical evidence.  
63 Such corroboration with recognized diagnostic standards was crucial for embedding machine learning into  
64 healthcare, ensuring trust in its predictive capabilities and foundational logic. Additionally, through model-  
65 specific associations with postmortem data, our study robustly validated the alignment of our model with  
66 dementia-related neurodegeneration. In essence, our model’s capacity to link probability predictions with  
67 semi-quantitative postmortem scores paves the way for the integration of deep learning methodologies with  
68 well-established clinical evidence.

69 While our study has the potential to advance the field of differential dementia diagnosis, it does  
70 have certain limitations that warrant consideration. Our model was developed and validated on multiple  
71 cohorts from numerous studies, and its full generalizability across diverse populations and clinical settings  
72 remains to be determined. Moving forward, we see potential in evaluating the model’s efficacy across the  
73 care continuum, encompassing primary care facilities, geriatric and general neurology practices, family  
74 medicine, and specialized clinics in tertiary medical centers. The datasets used in our study predominantly  
75 feature AD cases, which could potentially introduce a bias towards better recognition of this particular  
76 dementia subtype. Although we incorporated various dementia etiologies, the imbalanced representation  
77 might affect the model’s generalizability and sensitivity towards less frequent types. Also, we chose to  
78 amalgamate mild, moderate, and severe dementia cases into a single category. We acknowledge that this  
79 categorization method might not completely reflect the nuanced individual staging practiced in specific  
80 healthcare settings, where varying degrees of dementia severity carry distinct implications for treatment and  
81 management strategies. Our focus was primarily on differential diagnosis rather than disease staging, which  
82 motivated this decision. Future enhancements to our model could potentially include disease staging as an  
83 additional dimension, thereby augmenting its granularity and relevance.

84 The evidence collected from this study signals a convergence between advanced computational meth-  
85 ods and the nuanced task of differential diagnosis in dementia, crucial for scenarios with scarce resources and

86 the multifaceted realm of mixed dementia, a condition frequently encountered yet diagnostically complex.  
87 Our adaptable model efficiently integrates multimodal data, showing strong performance across diverse set-  
88 tings. Future validations, encompassing a wider demographic and geographical expanse, will be pivotal  
89 to substantiate the model's robustness and enhance its diagnostic utility in dementia care. Our pragmatic  
90 investigation accentuates the potential of neural networks to refine the granularity of diagnostic evaluations  
91 in neurocognitive disorders.

## 1 Methods

2 **Study population** We collected demographics, personal and family history, laboratory results, findings  
3 from the physical/neurological exams, medications, neuropsychological tests, and functional assessments as  
4 well as multi-sequence magnetic resonance imaging (MRI) scans from 9 distinct cohorts, totaling 50,686  
5 participants. There were 19,462 participants with normal cognition (NC), 9,209 participants with mild  
6 cognitive impairment (MCI), and 22,015 participants with dementia (DE). We further identified 10 primary  
7 and contributing causes of dementia: 17,298 participants with Alzheimer's disease (AD), 2,003 partici-  
8 pants with dementia with Lewy bodies and Parkinson's disease dementia (LBD), 2,032 participants with  
9 vascular brain injury or vascular dementia including stroke (VD), 114 participants with Prion disease in-  
10 cluding Creutzfeldt-Jakob disease (PRD), 3,076 participants with frontotemporal lobar degeneration and its  
11 variants, which includes corticobasal degeneration (CBD) and progressive supranuclear palsy (PSP), and  
12 with or without amyotrophic lateral sclerosis (FTD), 138 participants with normal pressure hydrocephalus  
13 (NPH), 808 participants suffering from dementia due to infections, metabolic disorders, substance abuse  
14 including alcohol, medications, delirium and systemic disease - a category termed as systemic and external  
15 factors (SEF), 2,700 participants suffering from psychiatric diseases including schizophrenia, depression,  
16 bipolar disorder, anxiety, and post-traumatic stress disorder (PSY), 265 participants with dementia due to  
17 traumatic brain injury (TBI), and 1,234 participants with dementia due to other causes which include neo-  
18 plasms, multiple systems atrophy, essential tremor, Huntington's disease, and seizures (ODE).

19 The cohorts include the National Alzheimer's Coordinating Center (NACC) dataset ( $n = 45,349$ ),<sup>47</sup>  
20 the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset ( $n = 1,821$ ),<sup>48</sup> the frontotemporal lo-  
21 bar degeneration neuroimaging initiative (NIFD) dataset ( $n = 253$ ),<sup>49</sup> the Parkinson's Progression Marker  
22 Initiative (PPMI) dataset ( $n = 198$ ),<sup>50</sup> the Australian Imaging, Biomarker and Lifestyle Flagship Study  
23 of Ageing (AIBL) dataset ( $n = 661$ ),<sup>51</sup> the Open Access Series of Imaging Studies-3 (OASIS) dataset  
24 ( $n = 491$ ),<sup>52</sup> the 4 Repeat Tauopathy Neuroimaging Initiative (4RTNI) dataset ( $n = 80$ ),<sup>53</sup> and three  
25 in-house datasets maintained by the Lewy Body Dementia Center for Excellence at Stanford University  
26 (LBDSU) ( $n = 182$ ),<sup>54</sup> and the Framingham Heart Study (FHS) ( $n = 1,651$ ).<sup>55</sup> Since its inception in 1948,  
27 FHS has been dedicated to identifying factors contributing to cardiovascular disease, monitoring multiple  
28 generations from Framingham, Massachusetts. Over time, the study has pinpointed major cardiovascular  
29 disease risk factors and explored their effects, while also investigating risk factors for conditions like de-  
30 mentia and analyzing the relationship between physical traits and genetics. Additional details on the study  
31 population are presented in Tables 1 & S1.

32 **Inclusion and exclusion criterion** Individuals from each cohort were eligible for study inclusion if  
33 they were diagnosed with normal cognition (NC), mild cognitive impairment (MCI), or dementia (DE). We  
34 used the National Alzheimer's Coordinating Center (NACC) dataset,<sup>47</sup> which is based on the Uniform Data  
35 Set (UDS) 3.0 dictionary,<sup>56</sup> as the baseline for our study. To ensure data consistency, we organized the data  
36 from the other cohorts according to the UDS dictionary. For individuals from the NACC cohort who had  
37 multiple clinical visits, we initially prioritized the visits at which the person received the diagnostic label  
38 of dementia. We then selected the visit with the most data features available prioritizing the availability of  
39 neuroimaging information. If multiple visits met all the above criteria, we chose the most recent visit among  
40 them. This approach maximized the sample sizes of dementia cases, as well as ensured that each individual  
41 had the latest record included in the study while maximizing the utilization of available neuroimaging and  
42 non-imaging data. We included participants from the 4RTNI dataset<sup>53</sup> with frontotemporal lobar degenera-

43 tion (FTD)-related disorders like progressive supranuclear palsy (PSP) or corticobasal syndrome (CBS). For  
44 other cohorts (NIFD,<sup>49</sup> PPMI,<sup>50</sup> LBDSU,<sup>54</sup> AIBL,<sup>51</sup> ADNI,<sup>48</sup> and OASIS<sup>52</sup>), participants were included if  
45 they had at least one MRI scan within 6 months of an officially documented diagnosis. From the FHS,<sup>55</sup> we  
46 utilized data from the Original Cohort (Gen 1) enrolled in 1948, and the Offspring Cohort (Gen 2) enrolled  
47 in 1971. For these participants, we selected available data including demographics, history, clinical exam  
48 scores, neuropsychological test scores, and MRI within 6 months of the date of diagnosis. We did not ex-  
49 clude cases based on the absence of features (including imaging) or diagnostic labels. Instead, we employed  
50 our innovative model training approach to address missing features or labels (See below).

51 **Data processing and training strategy** Various non-imaging features (n=391) corresponding to sub-  
52 ject demographics, medical history, laboratory results, medications, neuropsychological tests, and functional  
53 assessments were included in our study. We combined data from 4RTNI, AIBL, LBDSU, NACC, NIFD,  
54 OASIS, and PPMI to train the model. We used a portion of the NACC dataset for internal testing, while the  
55 ADNI, and FHS cohorts served for external validation (Tables 1, S1, S2, & S3). We used a series of steps  
56 such as standardizing the data across all cohorts and formatting the features into numerical or categorical  
57 variables before using them for model training. We used stratified sampling at the person-level to create  
58 the training, validation, and testing splits. As we pooled the data from multiple cohorts, we encountered  
59 challenges related to missing features and labels. To address these issues and enhance the robustness of our  
60 model against data unavailability, we incorporated several strategies such as random feature masking and  
61 masking of missing labels (see below).

62 Our investigation harnessed the potential of multi-sequence magnetic resonance imaging (MRI) vol-  
63 umetric scans sourced from diverse cohorts (Table S3). The majority of these scans encompassed a range of  
64 sequences, including T1-weighted, T2-weighted, diffusion-weighted imaging (DWI), susceptibility-weighted  
65 imaging (SWI), and fluid attenuated inversion recovery (FLAIR). The collected imaging data were stored  
66 in the NIFTI file format, categorized by participant and the date of their visit. The MRI scans underwent  
67 a singular pre-processing step, which involved skull stripping using SynthStrip,<sup>57</sup> a computational tool de-  
68 signed for extracting brain voxels from various image types. No registration procedures were applied to the  
69 resulting scans. To ensure the purity of the dataset, we excluded calibration, localizer, and 2D scans from  
70 the downloaded data before initiating model training.

71 **Backbone architecture** Our modeling framework harnesses the power of the transformer architecture  
72 to interpret and process a vast array of diagnostic parameters, including person-level demographics, medical  
73 history, neuroimaging, functional assessments, and neuropsychological test scores. Each of these distinct  
74 features is initially transformed into a fixed-length vector using a modality-specific strategy, forming the  
75 initial layer of input for the transformer model. Following this, the transformer acts to aggregate these  
76 vector inputs, decoding them into a series of predictions. A distinguishing strength of this framework lies  
77 in its integration of the transformer's masking mechanism,<sup>58,59</sup> strategically deployed to emulate missing  
78 features. This capability enhances the model's robustness and predictive power, allowing it to adeptly handle  
79 real-world scenarios characterized by incomplete data.

80 **Multimodal data embeddings** Transformers use a uniform representation for all input tokens, typi-  
81 cally in the form of fixed-length vectors. However, the inherent complexity of medical data, with its variety



82 of modalities, poses a challenge to this requirement. Therefore, medical data needs to be adapted into a  
83 unified embedding that our transformer model can process. The data we accessed falls into three primary  
84 categories: numerical data, categorical data, and imaging data. Each category requires a specific method of  
85 embedding. Numerical data typically encompasses those data types where values are defined in an ordinal  
86 manner that holds distinct real-world implications. For instance, chronological age fits into this category, as  
87 it serves as an indicator of the aging process. To project numerical data into the input space of the trans-  
88 former, we employed a single linear layer to ensure an appropriate preservation of the structure inherent to  
89 the original data space. Categorical data encompasses those inputs that can be divided into distinct cate-  
90 gories yet lack any implicit order or priority. An example of this is gender, which can be categorized as  
91 ‘male’ or ‘female’. We utilized a lookup table to translate categorical inputs into corresponding embeddings.  
92 It’s noteworthy that this approach is akin to a linear transformation when the data is one-hot vectorized,  
93 but is computationally efficient, particularly when dealing with a vast number of categories. Imaging data,  
94 which includes MRI scans in medical applications, can be seen as a special case of numerical data. How-  
95 ever, due to their high dimensionality and complexity, it is difficult to compress raw imaging data into a  
96 significantly lower-dimensionality vector using a linear transformation, while still retaining essential infor-  
97 mation. We leveraged the advanced capabilities of modern deep learning architectures to extract meaningful  
98 imaging embeddings (see below). Once these embeddings were generated, they were treated as numerical  
99 data, undergoing linear projection into vectors of suitable length, thus enabling their integration with other  
100 inputs to the transformer.

101 **Imaging feature extraction** We harnessed the Swin UNETR (Fig. S1),<sup>60,61</sup> a three-dimensional (3D)  
102 transformer-based architecture, to extract embeddings from a multitude of brain MRI scans, encompassing  
103 various sequences including T1-weighted (T1w), T2-weighted (T2w), diffusion-weighted (DWI), susceptibility-  
104 weighted (SWI), and fluid-attenuated inversion recovery (FLAIR) imaging sequences. The Swin UNETR  
105 model consists of a Swin Transformer encoder, designed to operate on 3D patches, seamlessly connected  
106 to a convolutional neural network (CNN)-based decoder through multi-resolution skip connections. Com-  
107 mencing with an input volume  $X \in \mathbb{R}^{H \times W \times D}$ , the encoder segmented  $X$  into a sequence of 3D tokens  
108 with dimensions  $\frac{H}{H'} \times \frac{W}{W'} \times \frac{D}{D'}$ , and projected them into a  $C$ -dimensional space via an embedding layer.  
109 It employed a patch size of  $2 \times 2 \times 2$  with a feature dimension of  $2 \times 2 \times 2 \times 1$  and an embedding space  
110 dimension of  $C = 48$ . The Swin UNETR encoder was subsequently interconnected with a CNN-based  
111 decoder at various resolutions through skip connections, collectively forming a ‘U-shaped’ network. This  
112 decoder amalgamated the encoder’s outputs at different resolutions, conducted upsampling via deconvolu-  
113 tions, ultimately generating a reconstruction of the initial input volume. The pre-trained weights were the  
114 product of self-supervised pre-training of the Swin UNETR encoder, primarily conducted on 3D volumes  
115 encompassing the chest, abdomen, and head/neck.<sup>60,61</sup>

116 The process of obtaining imaging embeddings began with several transformations applied to the MRI  
117 scans. These transformations included resampling the scans to standardized pixel dimensions, foreground  
118 cropping, and spatial resizing, resulting in the creation of sub-volumes with dimensions of  $128 \times 128 \times$   
119  $128$ . Subsequently, these sub-volumes were input into the Swin UNETR model, which in turn extracted  
120 encoder outputs sized at  $768 \times 4 \times 4 \times 4$ . These extracted embeddings underwent downsampling via a  
121 learnable embedding module, consisting of four convolutional blocks, to align with the input token size of  
122 the downstream transformer. As a result, the MRI scans were effectively embedded into one-dimensional  
123 vectors, each of size 128. These vectors were then combined with non-imaging features and directed into  
124 the downstream transformer for further processing. The entire process utilized a dataset comprising 11,438  
125 MRI volumes, which were allocated for model training, validation, and testing (Table S3).

126 **Random feature masking** To enhance the robustness of the backbone transformer in handling data  
127 incompleteness, we leveraged the masking mechanism<sup>58,59</sup> to emulate arbitrary missing features during  
128 training. The masking mechanism, when paired with the attention mechanism, effectively halts the informa-  
129 tion flow from a given set for input tokens, ensuring that certain features are concealed during prediction. A  
130 practical challenge arises when considering the potential combinations of input features, which increase ex-  
131 ponentially. With hundreds of features in play, capturing every potential combination is intractable. Inspired  
132 by the definition of Shapley values, we deployed an efficient strategy for feature dropout. Given a sample  
133 with feature set  $S$ ,  $S$  is randomly permuted as  $\sigma$ ; simultaneously, an integer  $i$  is selected independently from  
134 the range  $[1, |S|]$ . Subsequent to this, the features  $\sigma_{i+1}, \sigma_{i+2}, \dots, \sigma_{|S|}$  are masked out from the backbone  
135 transformer. It's noteworthy that the dropout process was applied afresh across different training batches or  
136 epochs to ensure that the model gets exposed to a diverse array of missing information even within a single  
137 sample.

138 **Handling missing labels** The backbone transformer was trained by amalgamating data from multiple  
139 different cohorts, each focused on distinct etiologies, which introduced the challenge of missing labels  
140 in the dataset. While most conventional approaches involve discarding records with incomplete output  
141 labels during training, we chose a more inclusive strategy to maximize the utility of the available data.  
142 Our approach framed the task as a multi-label classification problem, introducing thirteen separate binary  
143 heads, one for each target label. With this design, for every training sample, we generated a binary mask  
144 indicating the absence of each label. We then masked the loss associated with samples lacking specific  
145 labels before backpropagation. This method ensured optimal utilization of the dataset, irrespective of label  
146 availability. The primary advantage of this approach lies in its adaptability. By implementing this label-  
147 masking strategy, our model can be evaluated against datasets with varying degrees of label availability,  
148 granting us the flexibility to address a wide spectrum of real-world scenarios.

149 **Loss function** Our model was trained by minimizing the loss function ( $\mathcal{L}$ ) composed of two loss terms:  
150 “Focal Loss (FL)”<sup>62</sup> ( $\mathcal{L}_{\text{FL}}$ ) and “Ranking Loss (RL)” ( $\mathcal{L}_{\text{RL}}$ ), along with the standard L2 regularization term.  
151 FL is a variant of standard cross-entropy loss that addresses the issue of class imbalance. It assigns low  
152 weight to easy (well-classified) instances and employs a balance parameter. This loss function was used for  
153 each of the diagnostic categories (a total of 13, see Glossary 1). Therefore, our  $\mathcal{L}_{\text{FL}}$  term was:

$$\mathcal{L}_{\text{FL}} = \frac{1}{N} \sum_{k=1}^N \sum_{i=1}^{13} -y_{k,i} \alpha_i (1 - p_{k,i})^\gamma \log(p_{k,i}) - (1 - y_{k,i}) (1 - \alpha_i) (p_{k,i})^\gamma \log(1 - p_{k,i}),$$

154 where  $N$  was the batch size (i.e.,  $N = 64$ ), and other parameters and variables were as defined. The  
155 focusing parameter  $\gamma$  was set to 2, which had been reported to work well in most of the experiments in the  
156 original paper.<sup>62</sup> Moreover,  $\alpha_i \in [0, 1]$  was the balancing parameter that influenced the weights of positive  
157 and negative instances. It was set as the square of the complement of the fraction of samples labeled  
158 as 1, varying for each  $i$  due to the differing level of class imbalance across diagnostic categories (refer to  
159 Table 1). The FL term did not take inter-class relationships into account. To address these relationships  
160 in our overall loss function, we also incorporated the RL term that induced loss if the sigmoid outputs for  
161 diagnostic categories labeled as 0 were not lower than those labeled as 1 by a predefined margin of  $\epsilon$ , for  
162 any training sample  $k$ . We defined the RL term for any pair of diagnostic categories  $i$  and  $j$ , as follows:

$$\mathcal{L}_{\text{RL}}^{(i,j)}(\mathbf{p}_k, \mathbf{y}_k) = \max(0, (p_{k,i} - p_{k,j})(y_{k,j} - y_{k,i}) + \epsilon),$$

163 Overall, the RL term was:

$$\mathcal{L}_{\text{RL}} = \frac{1}{N} \sum_{k=1}^N \sum_{i=1}^{13} \sum_{j=i+1}^{13} \mathcal{L}_{\text{RL}}^{(i,j)}(\mathbf{p}_k, \mathbf{y}_k).$$

164 Combining all terms, our overall loss function ( $\mathcal{L}$ ) was:

$$\mathcal{L} = \mathcal{L}_{\text{FL}} + \lambda \mathcal{L}_{\text{RL}} + \beta \|\mathbf{w}\|^2,$$

165 where  $\lambda$  and  $\beta$  were the weights that controlled the importance of  $\mathcal{L}_{\text{RL}}$  and the L2 regularization terms,  
166 respectively. The training was done using the mini-batch strategy with the AdamW optimizer,<sup>63</sup> an improved  
167 version of the Adam optimizer,<sup>64</sup> with a learning rate of 0.0005 for a total of 256 epochs. Additionally, we  
168 utilized a cosine learning rate scheduler with warm restarts,<sup>65</sup> initiating the first restart after 64 batches and  
169 extending the restart period by a factor of 2 for each subsequent restart. The values of  $\epsilon$ ,  $\lambda$ , and  $\beta$  were  
170 determined to be  $\epsilon = 0.25$ ,  $\lambda = 0.025$ , and  $\beta = 0.01$ , respectively, based on an evaluation of the overall  
171 model performance on the validation set. During training, the model performance was evaluated on the  
172 validation set at the end of each epoch, and the model with the highest performance was selected.

173 **Interpretability analysis** The primary goal of interpretability analysis is to demystify machine learning  
174 models by providing clear insights into how various features influence predictions. At the heart of this field  
175 lies the Shapley value,<sup>31</sup> originally a game theory concept, now repurposed to evaluate feature significance  
176 in machine learning models. In this context, each instance is considered a unique ‘game’, where features  
177 act as players contributing to the outcome. The model’s output is analogous to the game’s payoff, with  
178 the Shapley value quantifying each feature’s contribution towards this outcome. Calculating exact Shapley  
179 values necessitates evaluating the model with every possible combination of missing features. Given the  
180 extensive size of our input features, this process becomes computationally prohibitive. To address this  
181 challenge, we employed the permutation sampling method for Shapley value estimation.<sup>66</sup> This approach  
182 involves randomly sampling permutations of the features and approximating the contribution of each feature  
183 to the prediction by averaging its marginal contributions across these permutations. It significantly reduces  
184 computational load while still providing a reasonable estimate of the Shapley values.

185 **Head-to-head expert validation** We evaluated our model’s predictive power against the diagnostic  
186 acumen of clinicians who are directly involved in dementia diagnosis and care. A group of neurologists and  
187 neuroradiologists were invited to participate in diagnostic tasks using a select subset of NACC cases (see  
188 ‘Data processing and training strategy’). Neurologists were presented with 100 cases – 15 cases each of  
189 NC and MCI, and 7 cases each of the dementia subtypes. The data encompassed person-level demograph-  
190 ics, medical history, social history, neuropsychological tests, functional assessments, and multi-sequence  
191 MRI scans where possible (i.e., T1-weighted, T2-weighted, FLAIR, DWI and SWI sequences). They were  
192 asked to provide their diagnostic impressions. Similarly, neuroradiologists were provided with the same  
193 multi-sequence MRI scans, along with details on age, gender, race, and education status from 70 clinically  
194 diagnosed DE cases. They were tasked with providing diagnostic impressions concerning the origin of de-  
195 mentia (Refer to Glossary 1). Additionally, neuroradiologists completed a REDCap questionnaire to rate  
196 atrophy and pathological changes in each brain sub-region on a scale from 0 to 3, with higher scores indi-  
197 cating more severe degeneration, and the presence or absence of specific disease markers such as infarcts.  
198 Case samples and questionnaires supplied to the neurologists and neuroradiologists can be found in the  
199 Supplementary Information.

200 **Neuropathologic validation** The model’s predictive capacity for various dementia etiologies was sub-  
201 stantiated through alignment with neuropathological evaluations sourced from the NACC, FHS and ADNI  
202 cohorts (Table S4). We included participants who conformed to the study’s inclusion criteria, had undergone  
203 MRI scans no more than three years prior to death, and for whom neuropathological data were available.  
204 Standardization of data was conducted in accordance with the Neuropathology Data Form Version 10 pro-  
205 tocols from the National Institute on Aging.<sup>67</sup> We pinpointed neuropathological indicators that influence the  
206 pathological signature of each dementia subtype, such as arteriolosclerosis, the presence of neurofibrillary  
207 tangles and amyloid plaques, cerebral amyloid angiopathy, and markers of tauopathy. These indicators were  
208 carefully chosen to reflect the complex pathological terrain that defines each form of dementia. To examine  
209 the Thal phase for amyloid plaques (A score), subjects were categorized into two groups: one encompassing  
210 Phase 0, indicative of no amyloid plaque presence, and a composite group merging Phases 1-5, reflect-  
211 ing varying degrees of amyloid pathology. The model’s predictive performance was then compared across  
212 these groupings. For the Braak stage of neurofibrillary degeneration (B score), we consolidated stages I-VI  
213 into a single collective, representing the presence of AD-type neurofibrillary pathology, whereas stage 0  
214 was designated for cases devoid of AD-type neurofibrillary degeneration. With respect to the density of  
215 neocortical neuritic plaques, assessed by the (CERAD or C score), individuals without neuritic plaques con-  
216 stituted one group, while those with any manifestation of neuritic plaques—sparse, moderate, or frequent  
217 (C1-C3)—were aggregated into a separate group for comparative analysis of the model’s predictive out-  
218 comes. The Thal, Braak, and CERAD scores were integrated into a composite ABC score, delineating the  
219 National Institute on Aging-Alzheimer’s Association (NIA-AA) criteria for AD neuropathological change  
220 (ADNC).<sup>68</sup> This summation resulted in two distinct groups for analysis: one encapsulating cases with no  
221 neuropathological evidence of Alzheimer’s disease, and another amalgamating cases classified across the  
222 spectrum of low, intermediate, and high ADNC. Furthermore, to evaluate the model’s concordance with  
223 non-AD pathologies, we analyzed the association between the model-generated probabilities and the pres-  
224 ence or absence of TDP-43 pathology, old microinfarcts, arteriosclerosis, and Prion disease.

225 **Statistical analysis** We used one-way ANOVA and the  $\chi^2$  test for continuous and categorical vari-  
226 ables, respectively to assess the overall differences in the population characteristics between the diagnostic  
227 groups across the study cohorts. We applied the Kruskal-Wallis H-test for independent samples and subse-  
228 quently conducted post-hoc Dunn’s testing with Bonferroni correction to evaluate the relationship between  
229 clinical dementia rating scores and the model-predicted probabilities, as well as the relationship between  
230 neuropathologic scores and the model-predicted probabilities. We opted for non-parametric tests because  
231 the Shapiro-Wilk test indicated significant deviations from normality. For comparing model predictions  
232 with expert-driven assessments, we used the Brunner Munzel test to identify statistically significant in-  
233 creases in the mean disease probability scores between the levels of scoring categories. We conducted a  
234 Shapiro-Wilk test on the distributions of the true negative and true positive cases for each etiology. The  
235 Brunner-Munzel test was then used to compare the expert and model confidence scores for the true negative  
236 and true positive cases for each etiology. To evaluate the interrater reliability of label-specific confidence  
237 scores, we performed pairwise Pearson correlation analyses between clinicians’ scores and those gener-  
238 ated by the model.<sup>69</sup> We calculated the average correlation coefficient across pairs and determined its 95%  
239 confidence interval. In addition, we assessed the correlation between the aggregated confidence score of  
240 neurologists and the model’s score for each diagnostic label. Using a bootstrapping approach with resam-  
241 pling, we created 1,000 iterations of the consensus score from the pool of individual neurologist scores. We  
242 then calculated the Pearson correlation for each iteration against the model’s scores, from which we derived  
243 the average correlation coefficient and its 95% confidence interval. All statistical analyses were conducted  
244 at a significance level of 0.05.

245 **Performance metrics** We generated receiver operating characteristic (ROC) and precision-recall (PR)  
246 curves from predictions on both the NACC test data and other datasets. From each ROC and PR curve,  
247 we further derived the area under the curve values. Also, we evaluated the model's accuracy, sensitivity,  
248 specificity, and Matthews correlation coefficient, with the latter being a balanced measure of quality for  
249 classes of varying sizes in a binary classifier. We also gauged inter-expert consensus using Cohen's kappa  
250 ( $\kappa$ ), which measures the degree to which two experts agree on a diagnosis. For each subgroup task, we  
251 computed the average pairwise  $\kappa$  as a comprehensive measure of agreement between the expert clinicians.

252 **Computational hardware and software** All MRI and non-imaging data were processed on a worksta-  
253 tion equipped with an Intel i9 14-core 3.3 GHz processor and 4 NVIDIA RTX 2080Ti GPUs. Our software  
254 development utilized Python (version 3.7.7) and the models were developed using PyTorch (version 1.5.1).  
255 We used several other Python libraries to support data analysis, including pandas (version 1.0.3), scipy (ver-  
256 sion 1.3.1), tensorboardX (version 1.9), torchvision (version 0.6), and scikit-learn (version 0.22.1). Training  
257 the model on a single Quadro RTX8000 GPU on a shared computing cluster had an average runtime of 7  
258 minutes per epoch, while the inference task took less than a minute per instance. All clinicians reviewed  
259 MRIs using 3D Slicer (version 4.10.2) and logged their findings in REDCap (version 11.1.3).

260 **Data and code availability** Data from ADNI, AIBL, NACC, NIFD, OASIS, PPMI and 4RTNI can be  
261 downloaded from publicly available resources. Data from FHS and LBDSU can be obtained upon request,  
262 subject to institutional approval. The Python scripts used in this study can be found on the Kolachalama  
263 Laboratory's GitHub page (<https://github.com/vkola-lab>).

## 1 Acknowledgements

2 This project was supported by grants from the Karen Toffler Charitable Trust (VBK), National Institute  
3 on Aging's Artificial Intelligence and Technology Collaboratories (P30-AG073014, VBK), the American  
4 Heart Association (20SFRN35460031, VBK & RA), Gates Ventures (RA & VBK), the Michael J. Fox  
5 Foundation (KLP), and the National Institutes of Health (R01-HL159620 [VBK], R21-CA253498 [VBK],  
6 R43-DK134273 [VBK], RF1-AG062109 [RA & VBK], U19-AG068753 [RA], P20-GM130447 [OT], K23-  
7 NS075097 [KLP], P50-AG047366 [KLP], and R01-NS115114 [KLP]). We acknowledge grant support from  
8 Boston University, CTSI 1UL1TR001430, for the REDCap Survey. We acknowledge the efforts of several  
9 individuals from the ADNI, AIBL, FHS, LBDSU, NACC, NIFD, OASIS, PPMI, and 4RTNI for providing  
10 access to data. Finally, we thank Drs. Shangran Qiu, Joyce C. Lee, Courtney E. Takahashi, Andrew M.  
11 Stern and Jesse B. Mez for several useful discussions.

## 12 Contributions

13 C.X. and S.S.K. contributed equally to this work. S.S.K., D.L., S.P., O.T.Z., A.S.W., A.K., C.K., and  
14 T.F.A.A. performed data collection. C.X. and S.S.K. designed and developed the machine learning frame-  
15 work. C.X., S.S.K., D.L., and O.B.G. performed model training and validation. S.S.K., S.P. and M.A.  
16 performed statistical analysis. C.X., S.S.K., D.L., S.P., O.T.Z., A.S.W., O.B.G., J.D.Z., S.T.P. and M.A.  
17 generated the figures and tables. V.C.A.A., B.C.D., C.W.F., H.H., S.K., A.Z.M., D.L.M., S.O., A.B.P., S.R.,  
18 M-H.S-H., E.A.S., B.N.S., J.E.S., A.S., O.T., J.Y., Y.Z. and S.Z. are practicing clinicians who reviewed the  
19 cases. S.A.B. and B.A.P. provided guidance on the modeling framework. K.L.P. and R.A. provided access  
20 to data. V.B.K. wrote the manuscript. All authors reviewed and approved the manuscript. V.B.K. conceived,  
21 designed and directed the study.

## 22 Ethics declarations

23 V.B.K. is on the scientific advisory board for Altoida Inc., and serves as a consultant to AstraZeneca. S.K.  
24 serves as consultant to AstraZeneca. C.W.F. is a consultant to Boston Imaging Core Lab. K.L.P. is a mem-  
25 ber of the scientific advisory boards for Curasen, Biohaven, and Neuron23, receiving consulting fees and  
26 stock options, and for Amprion, receiving stock options. R.A. is a scientific advisor to Signant Health and  
27 NovoNordisk. She also serves as a consultant to Davos Alzheimer's Collaborative. The remaining authors  
28 declare no competing interests.

## 1 References

- 3 1. Cahill, S. Who's global action plan on the public health response to dementia: some challenges and  
4 opportunities. *Aging & Mental Health* **24**, 197–199 (2019).
- 5 2. Dall, T. M. *et al.* Supply and demand analysis of the current and future us neurology workforce.  
6 *Neurology* **81**, 470–478 (2013).
- 7 3. Burton, A. How do we fix the shortage of neurologists? *The Lancet Neurology* **17**, 502–503 (2018).
- 8 4. Lester, P. E., Dharmarajan, T. S. & Weinstein, E. The looming geriatrician shortage: Ramifications and  
9 solutions. *J Aging Health* **32**, 1052–1062 (2020). Epub 2019 Oct 4.
- 10 5. Knopman, D. S. *et al.* Practice parameter: Diagnosis of dementia (an evidence-based review). *Neurol-*  
11 *ogy* **56**, 1143–1153 (2001).
- 12 6. Thijssen, E. H. & Rabinovici, G. D. Rapid progress toward reliable blood tests for alzheimer disease.  
13 *JAMA Neurology* **78**, 143–145 (2021).
- 14 7. Teunissen, C. E. *et al.* Blood-based biomarkers for alzheimer's disease: towards clinical im-  
15 plementation. *Lancet Neurology* **21**, 66–77 (2022). URL [https://doi.org/10.1016/](https://doi.org/10.1016/S1474-4422(21)00361-6)  
16 [S1474-4422\(21\)00361-6](https://doi.org/10.1016/S1474-4422(21)00361-6).
- 17 8. Sevigny, J. *et al.* The antibody aducanumab reduces abeta plaques in alzheimer's disease. *Nature* **537**,  
18 50–56 (2016).
- 19 9. van Dyck, C. H. *et al.* Lecanemab in early alzheimer's disease. *New England Journal of Medicine* **388**,  
20 9–21 (2023).
- 21 10. Martin, S. A., Townend, F. J., Barkhof, F. & Cole, J. H. Interpretable machine learning for dementia: A  
22 systematic review. *Alzheimer's & Dementia* **19**, 2135–2149 (2023).
- 23 11. Myszczyńska, M. A. *et al.* Applications of machine learning to diagnosis and treatment of neurodegener-  
24 ative diseases. *Nature Reviews Neurology* **16**, 440–456 (2020).
- 25 12. Borchert, R. J. *et al.* Artificial intelligence for diagnostic and prognostic neuroimaging in dementia: A  
26 systematic review. *Alzheimer's & Dementia* (2023). Online ahead of print.
- 27 13. Ahmed, M. R., Mahmood, A. N., Huq, M. A., Funk, P. & Mafi, A. Neuroimaging and machine learn-  
28 ing for dementia diagnosis: Recent advancements and future prospects. *IEEE Reviews in Biomedical*  
29 *Engineering* **12**, 19–33 (2019).
- 30 14. Bron, E. E. *et al.* Ten years of image analysis and machine learning competitions in dementia. *Neu-*  
31 *roImage* **253** (2022).
- 32 15. Qiu, S. *et al.* Development and validation of an interpretable deep learning framework for alzheimer's  
33 disease classification. *Brain* **143**, 1920–1933 (2020).
- 34 16. Nemoto, K. *et al.* Differentiating dementia with lewy bodies and alzheimer's disease by deep learning  
35 to structural mri. *Journal of Neuroimaging* **31**, 579–587 (2021).
- 36 17. Zheng, Y., Zhang, Y., Zhang, Y., Wang, Y. & Zheng, B. Machine learning-based framework for dif-  
37 ferential diagnosis between vascular dementia and alzheimer's disease using structural mri features.  
38 *Frontiers in Neurology* **10** (2019).

- 39 18. Castellazzi, G. *et al.* A machine learning approach for the differential diagnosis of alzheimer and  
40 vascular dementia fed by mri selected features. *Frontiers in Neuroinformatics* **14** (2020).
- 41 19. Hu, J. *et al.* Deep learning-based classification and voxel-based visualization of frontotemporal demen-  
42 tia and alzheimer's disease. *Frontiers in Neuroscience* **14** (2021).
- 43 20. Kim, J. *et al.* Machine learning based hierarchical classification of frontotemporal dementia and  
44 alzheimer's disease. *NeuroImage: Clinical* **23** (2019).
- 45 21. Nguyen, H.-D. *et al.* Multimodal deep learning for alzheimer's disease dementia assessment. In *Medical*  
46 *Image Computing and Computer Assisted Intervention–MICCAI 2022*, 55–65 (2022).
- 47 22. Moguilner, S. *et al.* Visual deep learning of unprocessed neuroimaging characterises dementia subtypes  
48 and generalises across non-stereotypic samples. *EBioMedicine* **90**, 104540 (2023).
- 49 23. Qiu, S., Miller, M., Joshi, P. *et al.* Multimodal deep learning for alzheimer's disease dementia assess-  
50 ment. *Nature Communications* **13**, 3404 (2022).
- 51 24. Vemuri, P. *et al.* Antemortem differential diagnosis of dementia pathology using structural mri:  
52 Differential-stand. *NeuroImage* **55**, 522–531 (2011).
- 53 25. Chagué, P. *et al.* Radiological classification of dementia from anatomical mri assisted by machine  
54 learning-derived maps. *Journal of Neuroradiology* **48**, 412–418 (2021).
- 55 26. Burgos, N. *et al.* Machine learning for classification and prediction of brain diseases: recent advances  
56 and upcoming challenges. *Current Opinion in Neurology* **33**, 439–450 (2020).
- 57 27. Barnes, D. E. *et al.* Development and validation of a brief dementia screening indicator for primary  
58 care. *Alzheimers Dement* **10**, 656–665.e1 (2014). Author manuscript; available in PMC 2015 Nov 1.  
59 Published in final edited form as: *Alzheimers Dement*. 2014 Nov; 10(6): 656–665.e1. Published online  
60 2014 Feb 1.
- 61 28. Anstey, K. J. *et al.* A self-report risk index to predict occurrence of dementia in three independent  
62 cohorts of older adults: the anu-adri. *PLoS One* **9**, e86141 (2014). ECollection 2014.
- 63 29. Sindi, S. *et al.* The caide dementia risk score app: The development of an evidence-based mobile  
64 application to predict the risk of dementia. *Alzheimers Dement (Amst)* **1**, 328–333 (2015). ECollection  
65 2015 Sep.
- 66 30. Kivimäki, M. *et al.* Estimating dementia risk using multifactorial prediction models. *JAMA Netw Open*  
67 **6**, e2318132 (2023).
- 68 31. Shapley, L. S. A value for n-person games. *Contributions to the Theory of Games* **2**, 307–317 (1953).
- 69 32. Lex, A., Gehlenborg, N., Strobel, H., Vuillemot, R. & Pfister, H. Upset: Visualization of intersecting  
70 sets. *IEEE Transactions on Visualization and Computer Graphics* **20**, 1983–1992 (2014).
- 71 33. Jo, M. *et al.* The role of tdp-43 propagation in neurodegenerative diseases: integrating insights from  
72 clinical and experimental studies. *Experimental Molecular Medicine* **52**, 1652–1662 (2020). Epub  
73 2020 Oct 13.
- 74 34. Cairns, N. J. *et al.* Tdp-43 in familial and sporadic frontotemporal lobar degeneration with ubiquitin  
75 inclusions. *The American Journal of Pathology* **171**, 227–240 (2007).



- 76 35. O'Brien, J. T. & Thomas, A. Vascular dementia. *The Lancet* **386**, 1698–1706 (2015). URL <https://www.sciencedirect.com/science/article/pii/S0140673615004638>.  
77
- 78 36. Le, T. H. & Gean, A. D. Neuroimaging of traumatic brain injury. *Mount Sinai Journal of Medicine: A*  
79 *Journal of Translational and Personalized Medicine* **76**, 145–162 (2009).
- 80 37. Mirvis, S., Wolf, A., Numaguchi, Y., Corradino, G. & Joslyn, J. Posttraumatic cerebral infarction  
81 diagnosed by ct: prevalence, origin, and outcome. *American Journal of Neuroradiology* **11**, 355–360  
82 (1990).
- 83 38. Tian, H.-L. *et al.* Risk factors for posttraumatic cerebral infarction in patients with moderate or severe  
84 head trauma. *Neurosurgical Review* **31**, 431–437 (2008).
- 85 39. Haber, M. *et al.* Vascular abnormalities within normal appearing tissue in chronic traumatic brain injury.  
86 *Journal of Neurotrauma* **35**, 2250–2258 (2018).
- 87 40. Latronico, N. *et al.* Impact of a posttraumatic cerebral infarction on outcome in patients with TBI:  
88 The Italian Multicenter cohort INCEPT study. *Crit Care* **24**, 33 (2020). URL <https://link.springer.com/article/10.1186/s13054-020-2746-5>.  
89
- 90 41. Hornberger, M. *et al.* In vivo and post-mortem memory circuit integrity in frontotemporal dementia and  
91 alzheimer's disease. *Brain : a journal of neurology* (2012).
- 92 42. Rabinovici, G. *et al.* Distinct mri atrophy patterns in autopsy-proven alzheimer's disease and frontotem-  
93 poral lobar degeneration. *American Journal of Alzheimer's Disease & Other Dementias*® **22**, 474–488  
94 (2008).
- 95 43. Chu, M. *et al.* Investigating the roles of anterior cingulate in behavioral variant frontotemporal demen-  
96 tia: A pet/mri study. *Journal of Alzheimer's Disease* **84**, 1771–1779 (2021).
- 97 44. Armstrong, R. A., Lantos, P. L. & Cairns, N. J. Overlap between neurodegenerative disorders. *Neu-*  
98 *ropathology* **25**, 111–124 (2005). 15875904.
- 99 45. Schneider, J. A., Arvanitakis, Z., Bang, W. & Bennett, D. A. Mixed brain pathologies account for most  
100 dementia cases in community-dwelling older persons. *Neurology* **69**, 2197–2204 (2007).
- 101 46. Rahimi, J. & Kovacs, G. G. Prevalence of mixed pathologies in the aging brain. *Alzheimer's Research*  
102 *& Therapy* **6**, 82 (2014).
- 103 47. Beekly, D. L. *et al.* The national alzheimer's coordinating center (nacc) database: an alzheimer disease  
104 database. *Alzheimer Disease & Associated Disorders* **18**, 270–277 (2004).
- 105 48. Mueller, S. G. *et al.* Ways toward an early diagnosis in alzheimer's disease: The alzheimer's disease  
106 neuroimaging initiative (adni). *Alzheimer's & Dementia* **1**, 55–66 (2005).
- 107 49. Boxer, A. L. *et al.* Frontotemporal degeneration, the next therapeutic frontier: Molecules and ani-  
108 mal models for frontotemporal degeneration drug development. *Alzheimer's & Dementia* **9**, 176–188  
109 (2013).
- 110 50. Marek, K. *et al.* The parkinson progression marker initiative (ppmi). *Progress in Neurobiology* **95**,  
111 629–635 (2011). Biological Markers for Neurodegenerative Diseases.
- 112 51. Ellis, K., Ames, D., Martins, R., Hudson, P. & Masters, C. The australian biomarkers lifestyle and  
113 imaging flagship study of ageing. *Acta Neuropsychiatrica* **18**, 285–285 (2006).

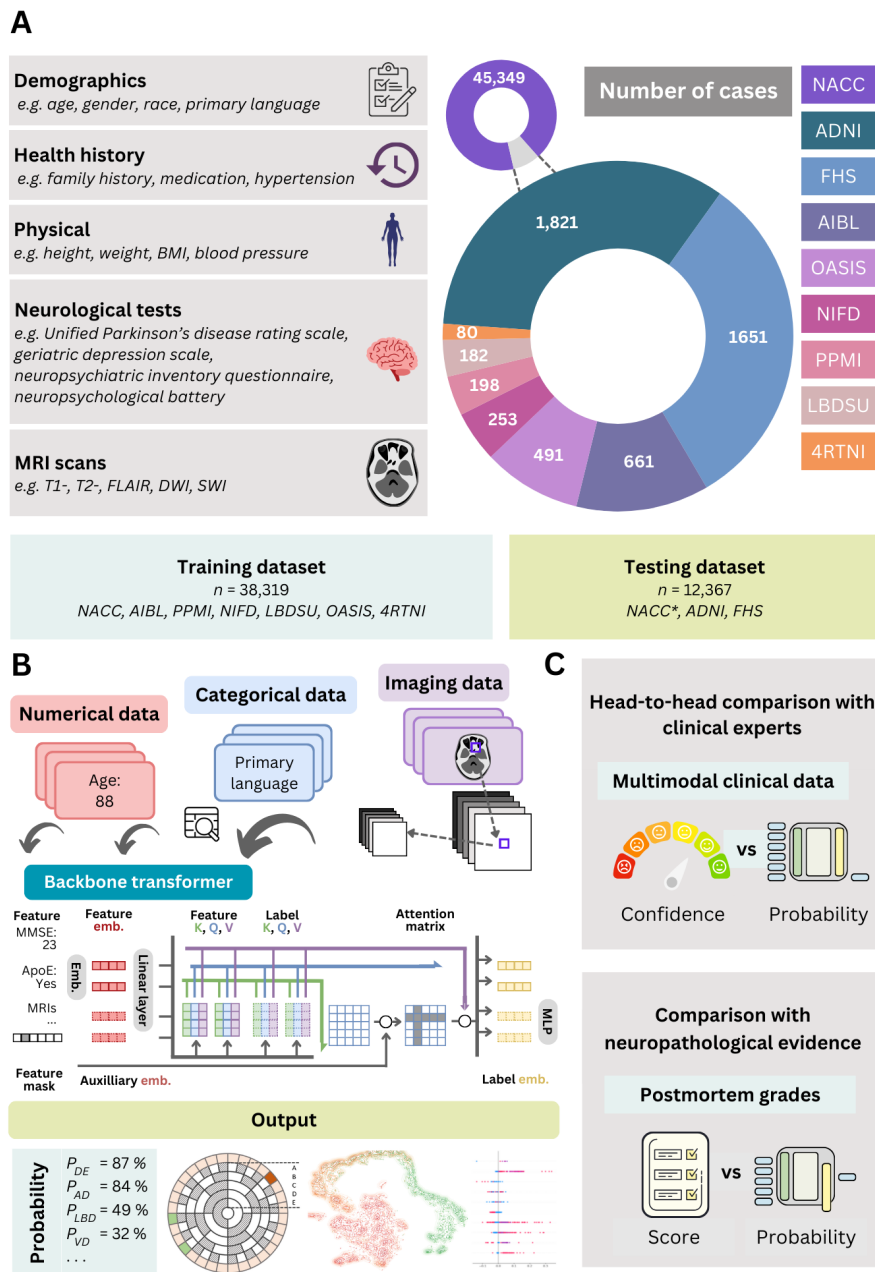
- 114 52. Marcus, D. S., Fotenos, A. F., Csernansky, J. G., Morris, J. C. & Buckner, R. L. Open Access Series  
115 of Imaging Studies: Longitudinal MRI Data in Nondemented and Demented Older Adults. *Journal of*  
116 *Cognitive Neuroscience* **22**, 2677–2684 (2010).
- 117 53. Dutt, S. *et al.* Progression of brain atrophy in psp and cbs over 6 months and 1 year. *Neurology* **87**,  
118 2016–2025 (2016).
- 119 54. Linortner, P. *et al.* White matter hyperintensities related to parkinson’s disease executive function.  
120 *Movement Disorders Clinical Practice* **7**, 629–638 (2020).
- 121 55. Yang, J. *et al.* Establishing cognitive baseline in three generations: Framingham heart study.  
122 *Alzheimer’s & Dementia: Diagnosis, Assessment & Disease Monitoring* **15** (2023).
- 123 56. Beekly, D. L. *et al.* The national alzheimer’s coordinating center (nacc) database: the uniform data set.  
124 *Alzheimer Disease & Associated Disorders* **21**, 249–258 (2007).
- 125 57. Hoopes, A., Mora, J. S., Dalca, A. V., Fischl, B. & Hoffmann, M. Synthstrip: skull-stripping for  
126 any brain image. *NeuroImage* **260**, 119474 (2022). URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S1053811922005900)  
127 [science/article/pii/S1053811922005900](https://www.sciencedirect.com/science/article/pii/S1053811922005900).
- 128 58. Vaswani, A. *et al.* Attention is all you need. In *Advances in neural information processing systems*,  
129 5998–6008 (2017).
- 130 59. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional transformers  
131 for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- 132 60. Hatamizadeh, A. *et al.* Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri  
133 images. arxiv 2022. *arXiv preprint arXiv:2201.01266*.
- 134 61. Tang, Y. *et al.* Self-supervised pre-training of swin transformers for 3d medical image analysis. In  
135 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 20730–20740  
136 (2022).
- 137 62. Lin, T.-Y., Goyal, P., Girshick, R., He, K. & Dollár, P. Focal loss for dense object detection. In *2017*  
138 *IEEE International Conference on Computer Vision (ICCV)*, 2999–3007 (2017).
- 139 63. Loshchilov, I. & Hutter, F. Decoupled weight decay regularization. In *International Confer-*  
140 *ence on Learning Representations* (2019). URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Bkg6RiCqY7)  
141 [Bkg6RiCqY7](https://openreview.net/forum?id=Bkg6RiCqY7).
- 142 64. Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv* (2014). 1412.6980.
- 143 65. Loshchilov, I. & Hutter, F. SGDR: Stochastic gradient descent with warm restarts. In *International*  
144 *Conference on Learning Representations* (2017). URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Skq89Scxx)  
145 [Skq89Scxx](https://openreview.net/forum?id=Skq89Scxx).
- 146 66. Mitchell, R., Cooper, J., Frank, E. & Holmes, G. Sampling permutations for shapley value estimation.  
147 *Journal of Machine Learning Research* **23**, 1–46 (2022).
- 148 67. National alzheimer’s coordinating center. neuropathology data form version 10, January 2014. URL  
149 <https://naccddata.org/data-collection/forms-documentation/np-10>.

- 150 68. Montine, T. J. *et al.* National institute on aging-alzheimer's association guidelines for the neuropathologic assessment of alzheimer's disease: A practical approach. *Acta Neuropathologica* **123**, 1–11  
151 (2012).  
152
- 153 69. de Raadt, W. M. B. R. e. a., A. A comparison of reliability coefficients for ordinal rating scales.  
154 *Journal of Classification* **38**, 519–543 (2021). URL [https://link.springer.com/article/](https://link.springer.com/article/10.1007/s00357-021-09386-5)  
155 [10.1007/s00357-021-09386-5](https://link.springer.com/article/10.1007/s00357-021-09386-5).

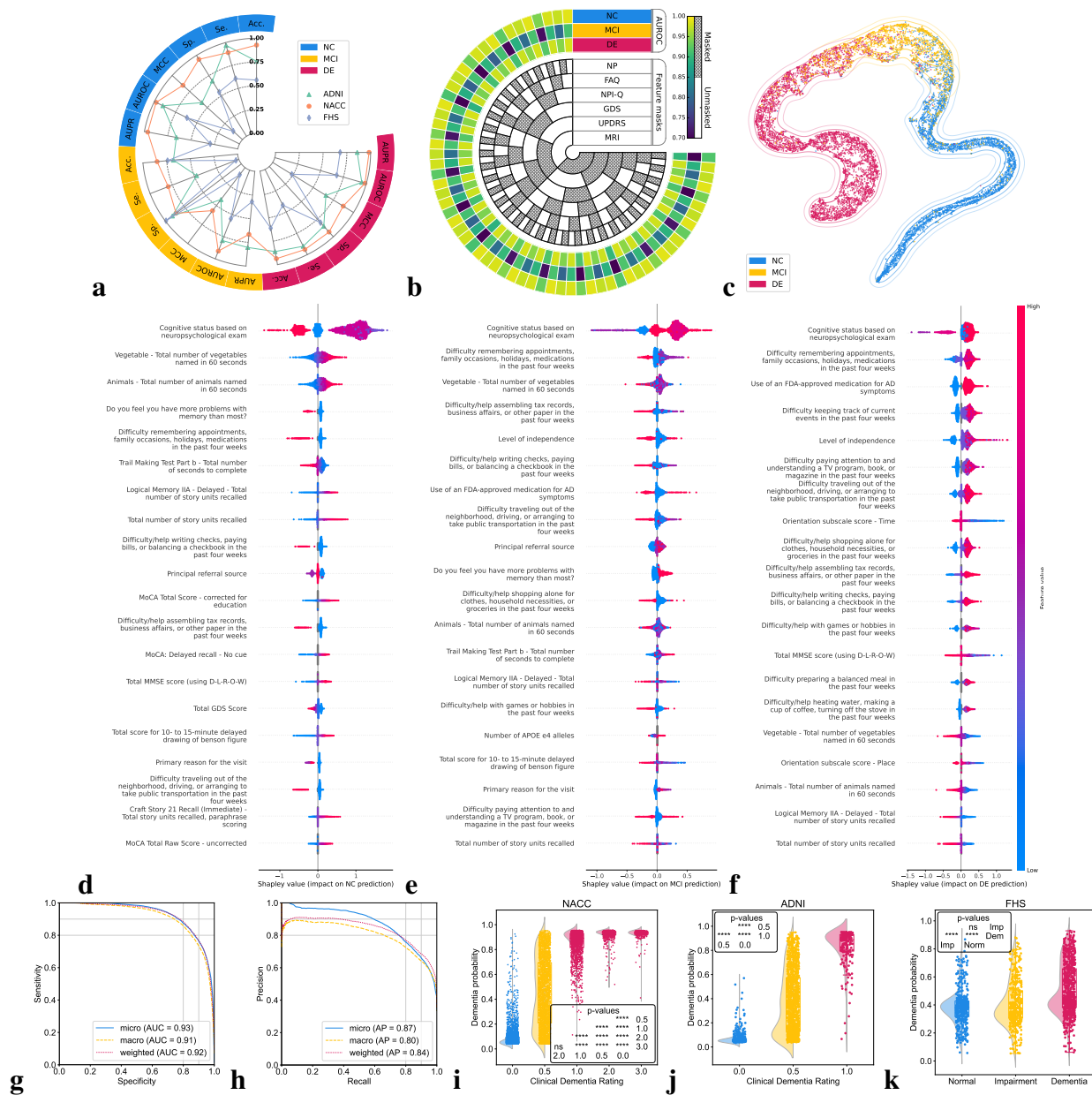
Dataset (group)	Age mean $\pm$ std	Male gender (percentage)	Education in years mean $\pm$ std	Race (White; Black; Asian; American Indian; Pacific; Multi-race)	CDR mean $\pm$ std
<b>NACC</b>					
NC [n = 17242]	71.25 $\pm$ 11.16	6009, 34.85%	15.83 $\pm$ 2.98 <sup>^</sup>	(13266, 2541, 528, 109, 10, 575) <sup>^</sup>	0.05 $\pm$ 0.15
MCI [n = 7582]	73.72 $\pm$ 9.81	3615, 47.68%	15.16 $\pm$ 3.45 <sup>^</sup>	(5708, 1185, 231, 53, 5, 276) <sup>^</sup>	0.45 $\pm$ 0.18
AD [n = 16131]	76.0 $\pm$ 10.31	7234, 44.85%	14.52 $\pm$ 3.74 <sup>^</sup>	(13161, 1702, 354, 92, 10, 458) <sup>^</sup>	1.2 $\pm$ 0.73
LBD [n = 1913]	75.01 $\pm$ 8.55	1365, 71.35%	15.12 $\pm$ 3.63 <sup>^</sup>	(1659, 128, 39, 17, 0, 37) <sup>^</sup>	1.29 $\pm$ 0.78
VD [n = 1919]	80.32 $\pm$ 8.76	947, 49.35%	14.15 $\pm$ 4.22 <sup>^</sup>	(1394, 332, 67, 2, 1, 68) <sup>^</sup>	1.22 $\pm$ 0.74
PRD [n = 114]	60.07 $\pm$ 10.36	62, 54.39%	14.8 $\pm$ 3.33 <sup>^</sup>	(93, 5, 5, 0, 1, 1) <sup>^</sup>	1.95 $\pm$ 0.95
FTD [n = 2898]	65.86 $\pm$ 9.36	1603, 55.31%	15.45 $\pm$ 3.09 <sup>^</sup>	(2664, 69, 73, 4, 5, 39) <sup>^</sup>	1.2 $\pm$ 0.83
NPH [n = 138]	79.1 $\pm$ 9.24	69, 50.0%	15.0 $\pm$ 3.28 <sup>^</sup>	(119, 10, 4, 0, 0, 4) <sup>^</sup>	1.18 $\pm$ 0.71
SEF [n = 808]	76.3 $\pm$ 11.15	413, 51.11%	14.6 $\pm$ 3.77 <sup>^</sup>	(646, 95, 15, 5, 2, 31) <sup>^</sup>	1.11 $\pm$ 0.7
PSY [n = 2700]	73.74 $\pm$ 10.78	1102, 40.81%	14.13 $\pm$ 4.12 <sup>^</sup>	(2163, 238, 59, 14, 5, 87) <sup>^</sup>	1.1 $\pm$ 0.64
TBI [n = 265]	72.87 $\pm$ 11.23	192, 72.45%	14.42 $\pm$ 4.13 <sup>^</sup>	(212, 27, 3, 2, 1, 11) <sup>^</sup>	1.11 $\pm$ 0.69
ODE [n = 1234]	72.94 $\pm$ 12.14	654, 53.0%	14.5 $\pm$ 3.78 <sup>^</sup>	(1046, 93, 28, 5, 4, 36) <sup>^</sup>	1.2 $\pm$ 0.76
<i>p-value</i>	<1.0e-200	<1.0e-200	<1.0e-200	8.341e-145	<1.0e-200
<b>NIFD</b>					
NC [n = 124]	63.21 $\pm$ 7.27	56, 45.16%	17.48 $\pm$ 1.87 <sup>^</sup>	(89, 0, 0, 0, 0, 3) <sup>^</sup>	0.03 $\pm$ 0.12 <sup>^</sup>
FTD [n = 129]	63.66 $\pm$ 7.33	75, 58.14%	16.18 $\pm$ 3.29 <sup>^</sup>	(109, 1, 1, 0, 0, 4) <sup>^</sup>	0.82 $\pm$ 0.54 <sup>^</sup>
<i>p-value</i>	6.266e-01	5.246e-02	2.606e-04	6.531e-01	4.333e-28
<b>PPMI</b>					
NC [n = 171]	62.74 $\pm$ 10.12	109, 63.74%	15.82 $\pm$ 2.93	(163, 3, 2, 0, 0, 1) <sup>^</sup>	N.A.
MCI [n = 27]	68.04 $\pm$ 7.32	22, 81.48%	15.52 $\pm$ 3.08	(24, 1, 1, 0, 0, 1)	N.A.
<i>p-value</i>	1.006e-02	1.115e-01	6.194e-01	2.910e-01	N.A.
<b>AIBL</b>					
NC [n = 480]	72.45 $\pm$ 6.22	203, 42.29%	N.A.	N.A.	0.03 $\pm$ 0.12
MCI [n = 102]	74.73 $\pm$ 7.11	53, 51.96%	N.A.	N.A.	0.47 $\pm$ 0.14
AD [n = 79]	73.34 $\pm$ 7.77	33, 41.77%	N.A.	N.A.	0.93 $\pm$ 0.54
<i>p-value</i>	5.521e-03	1.887e-01	N.A.	N.A.	4.542e-158
<b>OASIS</b>					
NC [n = 424]	71.34 $\pm$ 9.43	164, 38.68%	15.79 $\pm$ 2.62 <sup>^</sup>	(53, 18, 1, 0, 0, 0) <sup>^</sup>	0.0 $\pm$ 0.02
MCI [n = 27]	75.04 $\pm$ 7.25	14, 51.85%	15.19 $\pm$ 2.76	(4, 1, 0, 0, 0, 0) <sup>^</sup>	0.52 $\pm$ 0.09
AD [n = 32]	77.44 $\pm$ 7.42	20, 62.5%	15.19 $\pm$ 2.8	(8, 1, 0, 0, 0, 0) <sup>^</sup>	0.86 $\pm$ 0.44
LBD [n = 4]	74.75 $\pm$ 5.67	4, 100.0%	16.0 $\pm$ 2.83	N.A.	1.0 $\pm$ 0.0
FTD [n = 4]	64.25 $\pm$ 8.61	3, 75.0%	16.5 $\pm$ 2.96	(4, 0, 0, 0, 0, 0)	1.25 $\pm$ 0.75
<i>p-value</i>	7.789e-04	3.239e-03	5.507e-01	8.735e-01	2.855e-169
<b>LBDSU</b>					
NC [n = 134]	68.77 $\pm$ 7.62	61, 45.52%	17.27 $\pm$ 2.47 <sup>^</sup>	N.A.	N.A.
MCI [n = 35]	70.16 $\pm$ 8.41	26, 74.29%	16.6 $\pm$ 2.58	N.A.	N.A.
LBD [n = 13]	73.42 $\pm$ 7.81	8, 61.54%	16.77 $\pm$ 2.15	N.A.	N.A.
<i>p-value</i>	1.033e-01	7.863e-03	3.243e-01	N.A.	N.A.
<b>4RTNI</b>					
NC [n = 12]	68.08 $\pm$ 4.92	5, 41.67%	15.45 $\pm$ 2.57 <sup>^</sup>	(12, 0, 0, 0, 0, 0)	0.0 $\pm$ 0.0
MCI [n = 31]	67.61 $\pm$ 7.0	11, 35.48%	16.68 $\pm$ 4.02	(25, 1, 2, 0, 1, 1) <sup>^</sup>	0.55 $\pm$ 0.15
FTD [n = 37]	69.14 $\pm$ 7.43	20, 54.05%	16.46 $\pm$ 4.21	(31, 1, 0, 0, 1, 2) <sup>^</sup>	1.27 $\pm$ 0.55
<i>p-value</i>	6.691e-01	2.992e-01	6.843e-01	7.620e-01	5.700e-16
<b>ADNI</b>					
NC [n = 481]	74.26 $\pm$ 6.0	235, 48.86%	16.34 $\pm$ 2.67	(432, 36, 8, 1, 0, 3) <sup>^</sup>	0.0 $\pm$ 0.0 <sup>*</sup>
MCI [n = 971]	72.84 $\pm$ 7.71	572, 58.91%	15.94 $\pm$ 2.81	(903, 34, 15, 2, 2, 12) <sup>^</sup>	0.5 $\pm$ 0.04
AD [n = 369]	74.91 $\pm$ 7.84	203, 55.01%	15.18 $\pm$ 2.97	(343, 15, 7, 0, 0, 4)	0.77 $\pm$ 0.26
<i>p-value</i>	2.565e-06	1.364e-03	1.872e-08	1.132e-01	<1.0e-200
<b>FHS</b>					
NC [n = 394]	74.9 $\pm$ 10.22 <sup>*</sup>	206, 52.28%	N.A.	(394, 0, 0, 0, 0, 0)	0.0 $\pm$ 0.0
MCI [n = 434]	79.92 $\pm$ 8.8 <sup>*</sup>	203, 46.77%	N.A.	(434, 0, 0, 0, 0, 0)	0.49 $\pm$ 0.07
AD [n = 687]	82.99 $\pm$ 7.87 <sup>*</sup>	211, 30.71%	N.A.	(687, 0, 0, 0, 0, 0)	2.04 $\pm$ 0.88
LBD [n = 73]	79.34 $\pm$ 9.37 <sup>*</sup>	46, 63.01%	N.A.	(73, 0, 0, 0, 0, 0)	1.84 $\pm$ 0.84
VD [n = 113]	81.74 $\pm$ 7.3 <sup>*</sup>	48, 42.48%	N.A.	(113, 0, 0, 0, 0, 0)	1.85 $\pm$ 0.8
FTD [n = 8]	85.67 $\pm$ 5.91 <sup>*</sup>	4, 50.0%	N.A.	(8, 0, 0, 0, 0, 0)	2.0 $\pm$ 0.87
<i>p-value</i>	1.316e-31	7.905e-14	N.A.	1.0	<1.0e-200

Table 1: **Study population.** Nine independent datasets were used for this study, including ADNI, NACC, NIFD, PPMI, OASIS, LBDSU, 4RTNI, and FHS. Data from NACC, NIFD, PPMI, OASIS, LBDSU, and 4RTNI were used for model training. Data from ADNI, FHS, and a held-out set from NACC were used for model testing. The p-value for each dataset indicates the statistical significance of inter-group differences per column. We used one-way ANOVA and  $\chi^2$  tests for continuous and categorical variables, respectively. Please refer to Glossary 1 for more information on the acronyms. Here N.A. denotes not available. The symbol <sup>^</sup> indicates that data was not available for some subjects.

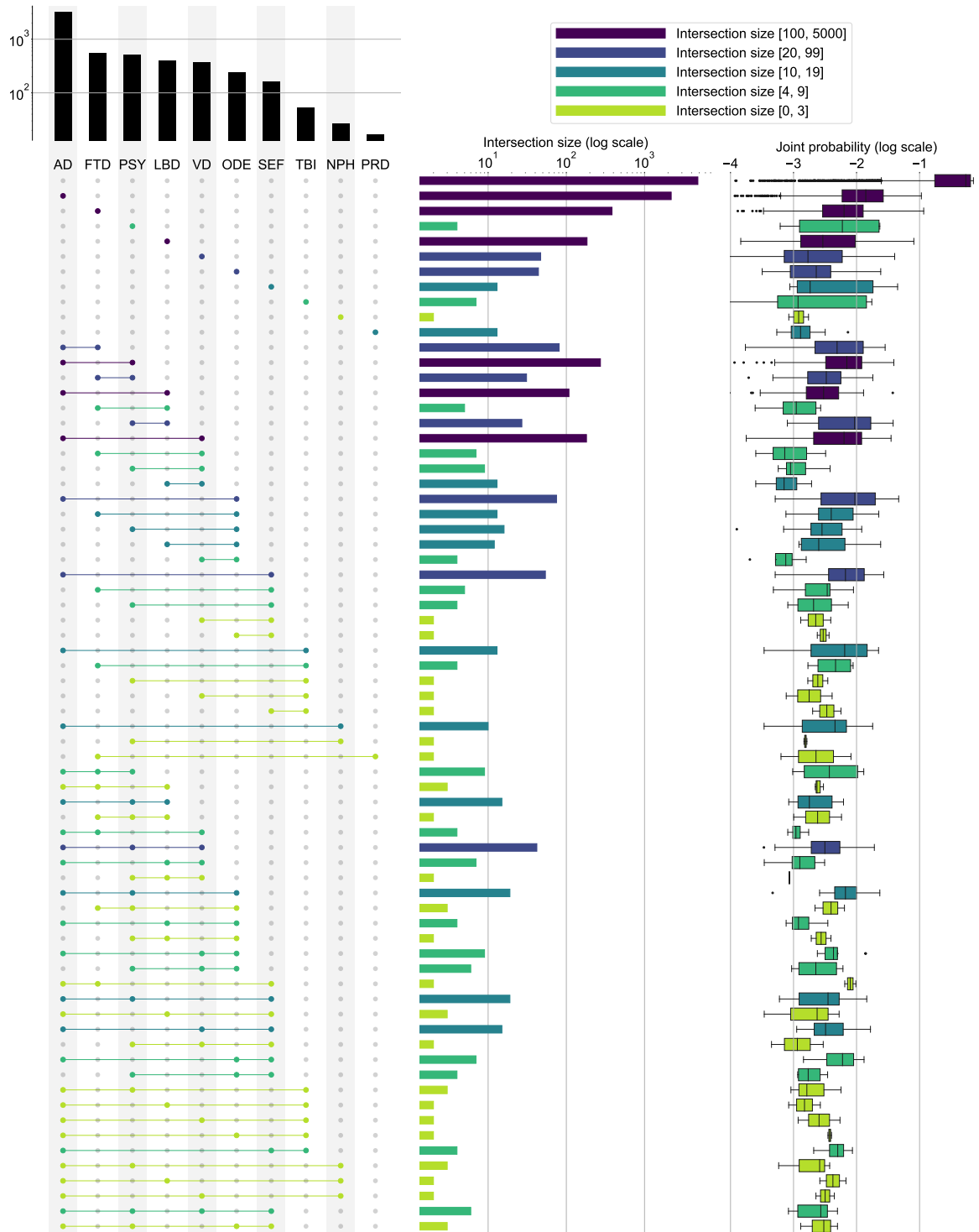
\* Due to the absence of CDR scores in the FHS dataset, we used the following definition: 0.0 - normal cognition, 0.5 - cognitive impairment, 1.0 - mild dementia, 2.0 - moderate dementia, 3.0 - severe dementia.



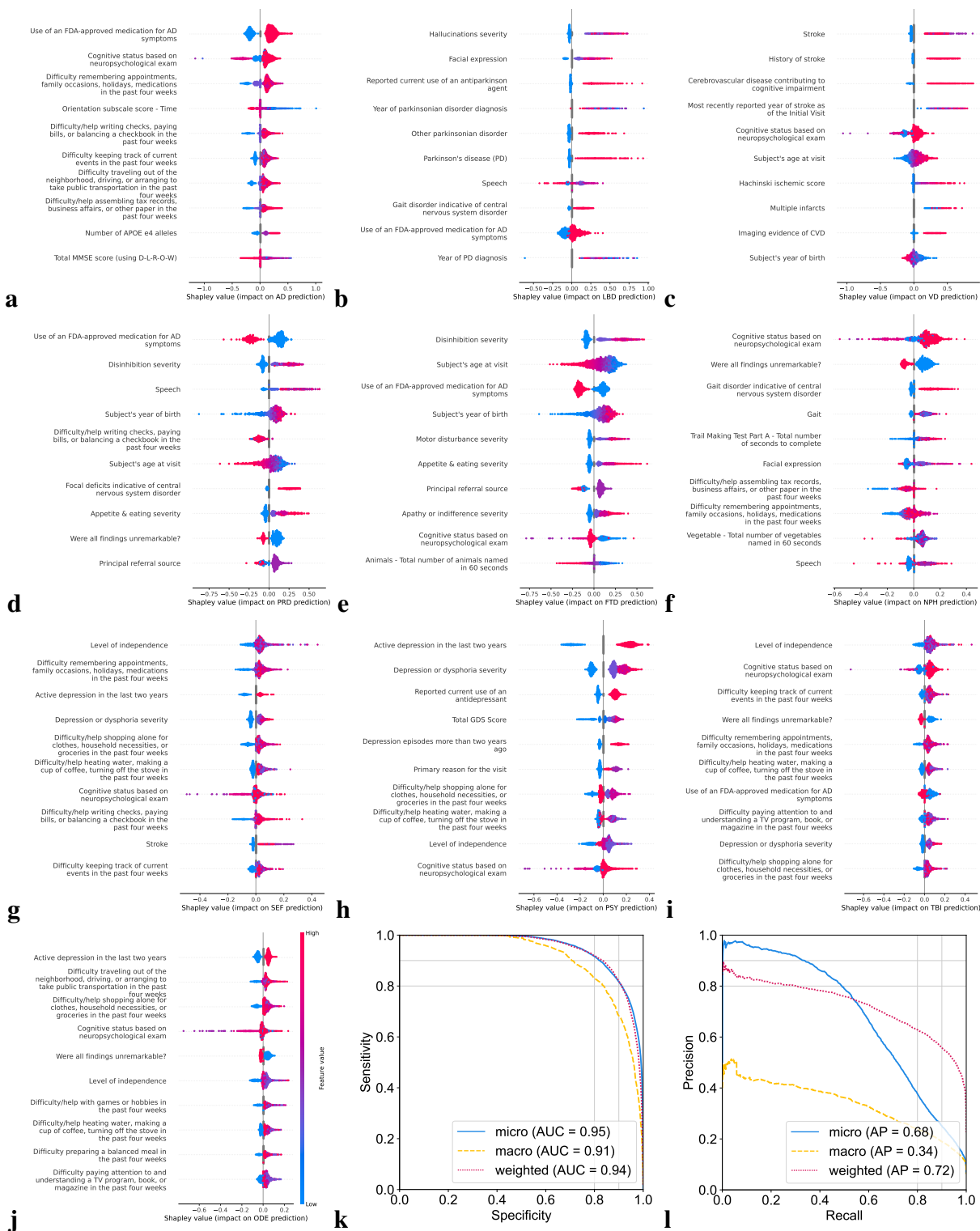
**Figure 1: Data, model architecture and modeling strategy.** (a) Our model for differential dementia diagnosis was developed using diverse data modalities, including individual-level demographics, health history, neurological testing, physical/neurological exams, and multi-sequence MRI scans. These data sources whenever available were aggregated from nine independent cohorts: 4RTNI, ADNI, AIBL, FHS, LBDSU, NACC, NIFD, OASIS, and PPMI (Tables 1 & S1). For model training, we merged data from NACC, AIBL, PPMI, NIFD, LBDSU, OASIS and 4RTNI. We employed a subset of the NACC dataset for internal testing. For external validation, we utilized the ADNI and FHS cohorts. (b) A transformer served as the scaffold for the model. Each feature was processed into a fixed-length vector using a modality-specific embedding strategy and fed into the transformer as input. A linear layer was used to connect the transformer with the output prediction layer. (c) A distinct portion of the NACC dataset was randomly selected to enable a comparative analysis of the model's performance against practicing neurologists. Furthermore, we conducted a direct comparison between the model and a team of practicing neuroradiologists using a random sample of cases with confirmed dementia from the NACC testing cohort. For both these evaluations, the model and clinicians had access to the same set of multimodal data. Finally, we assessed the model's predictions by comparing them with pathology grades available from the NACC, ADNI, and FHS cohorts.



**Figure 2: Model performance on individuals along the cognitive spectrum.** (a) Radar plot illustrating the performance of the model on individuals with normal cognition (NC), mild cognitive impairment (MCI), and dementia (DE) is shown. We present a range of metrics including mean values along with their standard deviations, for model accuracy, sensitivity, specificity, precision, area under the receiver operating characteristic curve, area under the precision-recall curve, F1-score, and Matthews correlation coefficient. (b) Chord diagram indicating varied levels of model performance in the presence of missing data. The inner concentric circles represent various scenarios in which particular test information was either omitted (masked) or included (unmasked). The three outer concentric rings depict the model's performance as measured by the area under the receiver operating characteristic curve (AUROC) for the NC, MCI and DE labels. (c) Two-dimensional t-distributed stochastic neighbor embeddings obtained from the penultimate layer of the model are shown. The legend in the lower-left corner indicates the color coding representing NC, MCI and DE, respectively. (d, e, f) Beeswarm plots visualize Shapley values for subjects classified as NC, MCI, and DE, respectively. (g, h) The receiver operating characteristic (ROC) and precision-recall (PR) curves are presented, with their respective micro-average, macro-average, and weighted-average calculations based on the labels for NC, MCI, and DE. These averaging techniques consolidated the model's performance across the spectrum of cognitive states. (i, j, k) Raincloud plots with overlying violin and box diagrams are shown to denote the distribution of clinical dementia rating scores (horizontal axis) versus model-predicted probability of dementia (vertical axis), on the NACC, ADNI and FHS cohorts, respectively. For (a,g,h), cases from the NACC testing, ADNI and FHS were used. Significance levels are denoted as 'ns' (not significant) for  $p \geq 0.05$ ; \* for  $p < 0.05$ ; \*\* for  $p < 0.01$ ; \*\*\* for  $p < 0.001$ ; and \*\*\*\* for  $p < 0.0001$  based on Kruskal-Wallis H-test for independent samples followed by post-hoc Dunn's testing with Bonferroni correction.

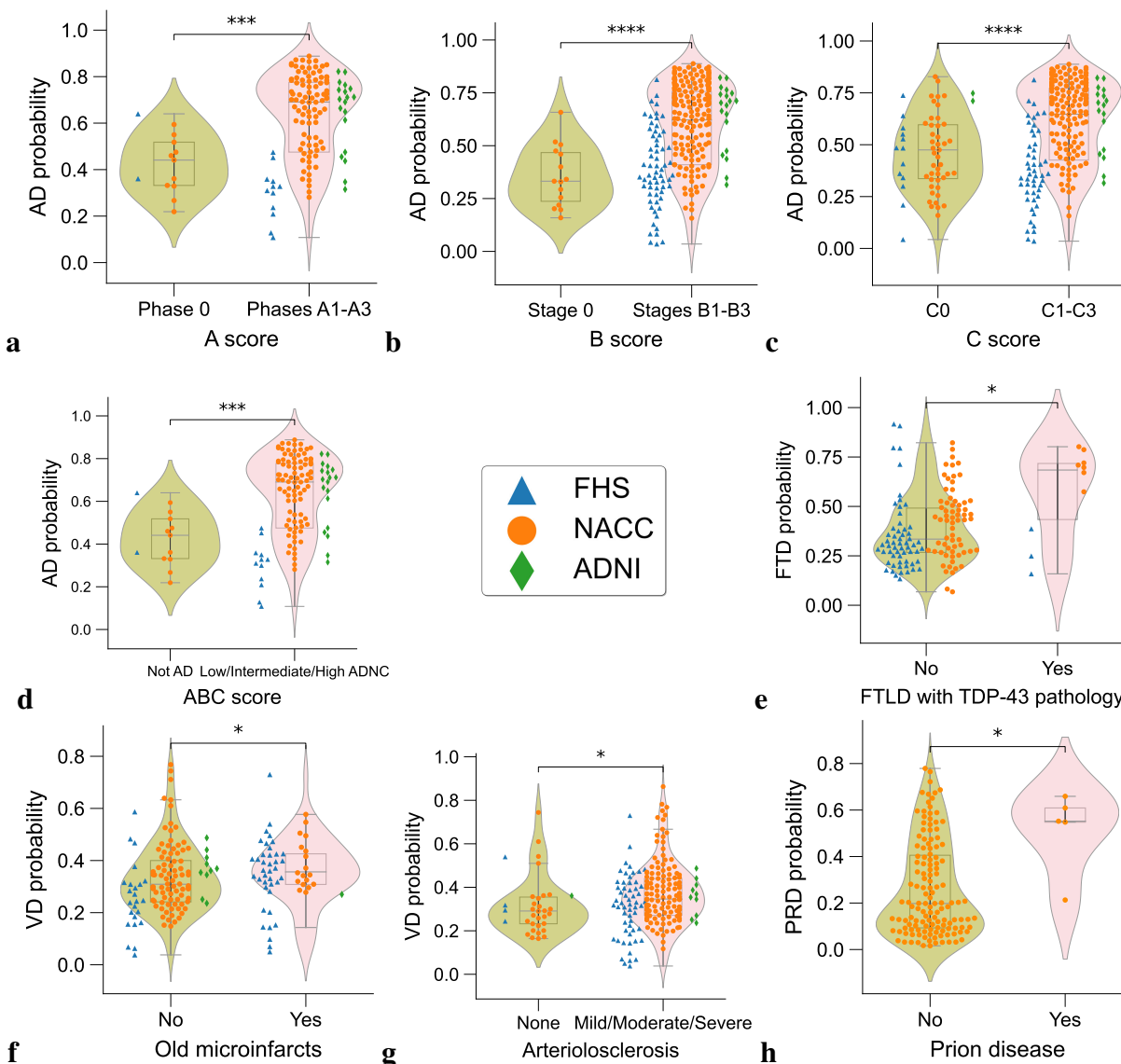


**Figure 3: Model assessment on mixed dementias.** Visualization of the distribution and model-predicted probabilities of various etiological categories found within the NACC testing dataset. The left segment enumerates both single and co-occurring diagnostic categories, offering a tally of each condition's frequency within the dataset. In the center, a logarithmic scale is used to delineate the overlap among these categories, shedding light on their relative commonality and the extent of their coexistence. This method grants a refined perspective on the prevalence of comorbid conditions. Additionally, the legend in the upper right interprets the counts within the central panel, providing a reference for the logarithmic data representation. The panel on the far right features a box-and-whisker plot, delineating the spread and central tendency of the model's predicted probabilities for each combination of diagnostic categories.

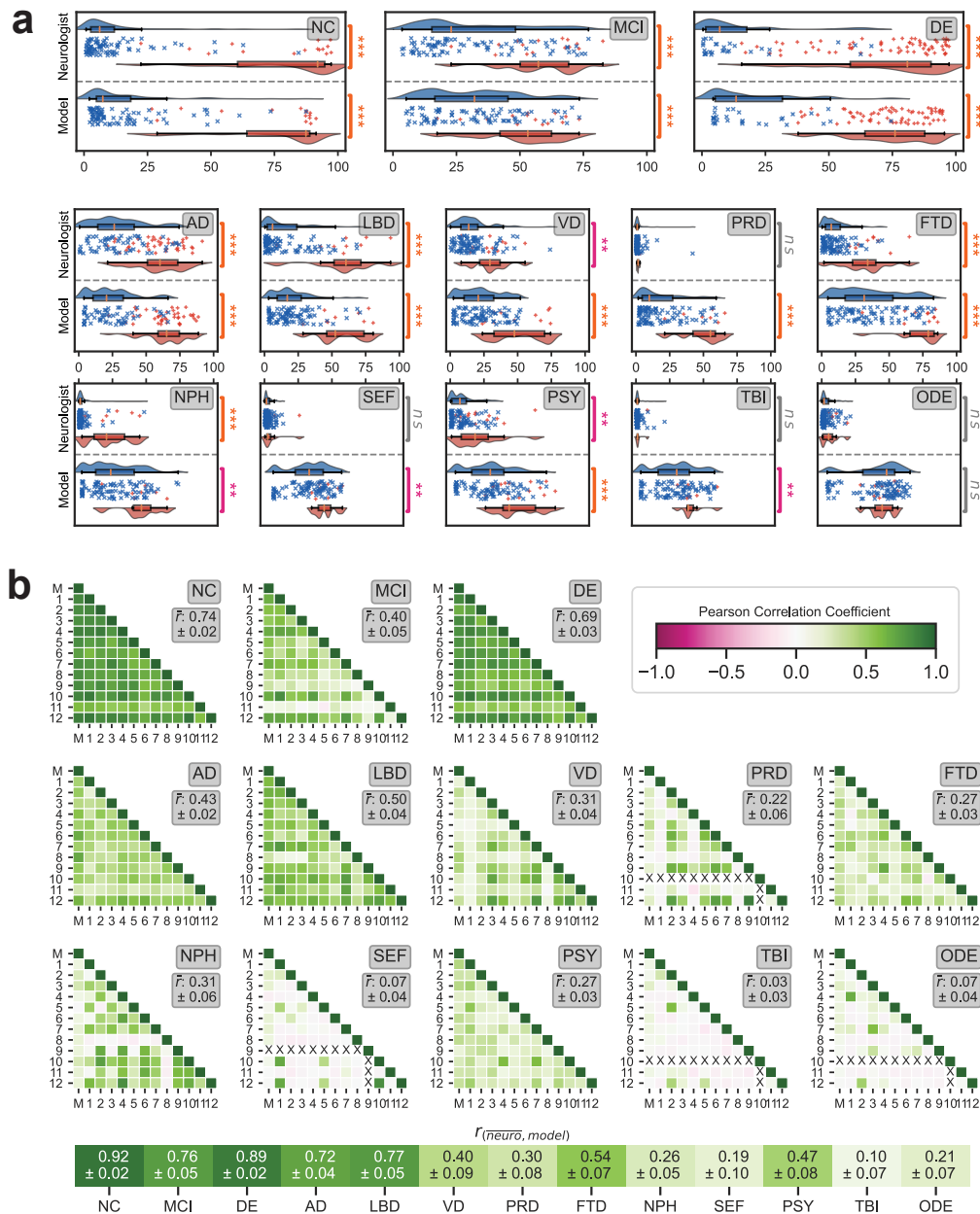


**Figure 4: Model assessment on dementia etiologies.** (a-j) Beeswarm plots illustrating Shapley values for subjects classified with AD, LBD, VD, PRD, FTD, NPH, SEF, PSY, TBI, and ODE are depicted. Adjacent to plot (j) is a colorbar that delineates the range of feature values for these conditions. (k, l) Receiver operating characteristic (ROC) and precision-recall (PR) curves are provided, utilizing micro-average, macro-average, and weighted-average methods across all the dementia diagnostic labels. These averages were computed to synthesize the performance metrics across all the dementia etiologies.





**Figure 5: Neuropathological validation.** Array of violin plots with integrated box plots, delineating the probability distributions as predicted by the model for different neuropathological grades. The analysis encompasses data from three distinct cohorts: the Framingham Heart Study (FHS), the National Alzheimer’s Coordinating Center (NACC), and the Alzheimer’s Disease Neuroimaging Initiative (ADNI), each denoted by unique markers (triangles, circles, and diamonds, respectively). Statistical significance is encoded using asterisks, determined by Dunn-Bonferroni post-hoc test: one asterisk (\*) for a p-value less than 0.05; two asterisks (\*\*) for p-values less than 0.01; three asterisks (\*\*\*) for p-values less than 0.001; and four asterisks (\*\*\*\*) for p-values less than 0.0001, reflecting increasing levels of statistical significance.



**Figure 6: Neurologist-level validation.** (a) Comparison between model-predicted probability scores and the assessments provided by practicing neurologists is shown. For the analysis, neurologists were given 100 randomly selected cases encompassing individual-level demographics, health history, neurological tests, physical as well as neurological examinations, and multi-sequence MRI scans. The neurologists were then tasked with assigning confidence scores for NC, MCI, DE, and the 10 dementia etiologies: AD, LBD, VD, PRD, FTD, NPH, SEF, PSY, TBI, and ODE (see Glossary 1). In the visual representation, the boxplot in blue indicates the distribution of confidence scores for true negative cases, while the boxplot in red signifies true positive cases. The symbol '+' represents true positive cases, and 'x' denotes true negative cases. Significance levels are denoted as: ns (not significant) for  $p \geq 0.05$ ; \* for  $p < 0.05$ ; \*\* for  $p < 0.01$ ; \*\*\* for  $p < 0.001$ ; and \*\*\*\* for  $p < 0.0001$ . These levels were determined using pairwise comparisons via the Brunner-Munzel test. (b) The figure presents the Pearson correlation coefficient across different diagnostic categories, comparing assessments from the neurologists ( $n = 12$ ) and the model, marked as 'M'. Each diagnostic category from NC to ODE includes a matrix reflecting correlation coefficient values between individual neurologists and the model. Shades of green signify positive correlation, indicating agreement between the model and neurologists, whereas magenta shades suggest negative correlations, indicating potential discrepancies in assessments. The mean pairwise Pearson correlation coefficient for each etiology is presented along with a 95% confidence interval. The symbol 'X' denotes rater pairs where the Pearson correlation was not calculable, due to one or both raters giving label-specific confidence scores with no variance. The heatmap at the bottom shows the Pearson correlation coefficients between model probabilities and mean neurologist confidence scores for each label.