
Language-Based Dementia Classification Should Consider Model Cognition for Interpretability

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Current approaches to dementia detection in machine learning based on language
2 often treat the task as an end-to-end binary classification problem, directly classi-
3 fying a person’s audio or transcripts to a final diagnostic label (e.g., dementia or
4 cognitively normal). While these prior techniques can be effective for accuracy,
5 because predictions are made in a single step, this binary approach overlooks the
6 progressive nature of cognitive decline and lacks the interpretive analyses that
7 clinicians rely on in clinical settings. Furthermore, these approaches face inter-
8 pretability limitations, particularly in terms of the linguistic features the models
9 are focusing on. **This position paper argues that language-based dementia**
10 **detection research should be re-framed to include cognition-based reasoning,**
11 **specifically by probing model cognition through structured, fine-grained in-**
12 **put–output analyses, which allow clearer understanding of how ML systems**
13 **reason in this task.** This new direction will help advance current ML models
14 towards dementia detection frameworks that are more interpretable and clinically
15 trustworthy.

16 1 Introduction: Diagnosing Dementia Isn’t Just a Classification Task

17 Machine learning (ML) and deep learning have become important tools in healthcare for assisting
18 in the analysis of medical records, accelerating testing, and predicting disease risks [46, 42]. ML
19 could be particularly useful for screening for cognitive impairment, including dementia, which is
20 estimated to affect 57 million people worldwide [39]. Dementia research using ML has assessed
21 input modalities such as demographic data, neuro-imaging, bio-fluid biomarkers, genetic and medical
22 information, and cognitive assessments [8, 14, 52, 43]. Among these, language-based inputs offer a
23 non-invasive and sensitive source of dementia-related cues.

24 The most common task used in language-based dementia detection is the picture description task
25 [11, 19, 41, 26], wherein a participant describes everything they see in an image (e.g., Cookie Theft
26 Task [16]) as their audio is recorded. The audio is then analyzed through audio signal processing [60,
27 26, 7, 21] and natural language processing (NLP) of the transcripts [19, 41, 9].

28 Current ML research on dementia detection is commonly framed as a binary classification task,
29 typically distinguishing Alzheimer’s disease (AD) or mild cognitive impairment (MCI) from normal
30 cognition (NC) [54, 44, 4], as seen in Fig. 1. Existing methods range from traditional classifiers built
31 on handcrafted features to large language models (LLMs) applied directly to transcripts, but in both
32 cases models are trained to map inputs to a single diagnostic label [22, 26, 21, 34, 1, 50]. Datasets
33 often reinforce this binary framing, with recordings annotated only with the final diagnosis [53, 5, 30–
34 32], or with cognitive test score (e.g., Mini Mental State Exam (MMSE) [33])(for an overview of
35 datasets, see [12] and Appendix A). While effective at optimizing categorical or score accuracy, it

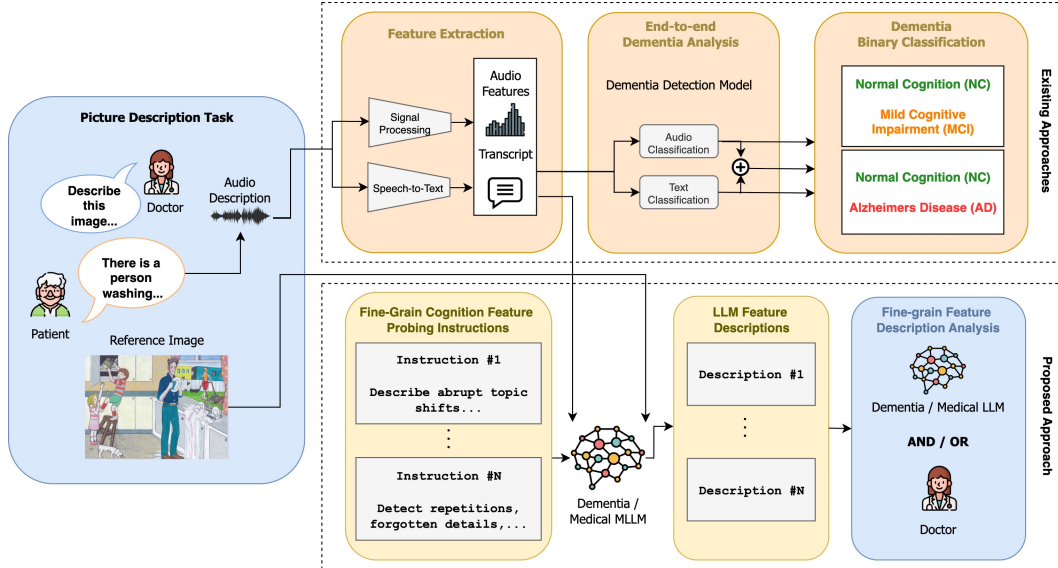


Figure 1: Overview of picture description task for dementia detection. Existing approaches classify audio or text, or both, into NC/MCI or NC/AD, often as a binary task. Our proposed approach probes model cognition through fine-grained feature analysis to assess reasoning and better support clinical diagnosis, offering more interpretable diagnosis. See Appendix B for more details.

36 collapses the progressive and multidimensional nature of dementia into a single outcome and often
 37 lacks interpretability of the model itself.

38 This raises a fundamental question for the field: *can we improve interpretability of dementia detection*
 39 *models if they reason through the linguistic cues?* Clinical cognitive assessment is inherently a
 40 staged process, and clinicians assess narrative completeness and structure, evaluate semantic richness,
 41 and integrate these with other observations before reaching a conclusion. When we reduce the
 42 problem to a single binary decision, *what aspects of cognition are we ignoring, and what biases are*
 43 *introduced—both by our task design and by the model’s own reasoning?*

44 **We argue that detection of cognitive impairment should be reformulated as a structured**
 45 **cognitive inference problem.** Rather than collapsing information into a single prediction, LLMs
 46 can serve as a cognitive analysis tool: first extracting fine-grained linguistic features from transcripts,
 47 then assessing narrative coherence and task faithfulness, and separately analyzing information with
 48 acoustic biomarkers before any diagnostic decision for a description of the patient’s cognition.
 49 Producing these intermediate outputs provides clinicians, who are ultimately responsible for making a
 50 diagnosis, the full context and rationale for why the model identified an individual as having cognitive
 51 impairment. The clinicians can then decide to include – or not include – aspects of this information
 52 in their assessment. To provide more clinically meaningful information, this will require datasets
 53 with richer annotations and evaluation methods that include clinically relevant features, to produce
 54 dementia detection systems that are interpretable and aligned with real-world healthcare applications.

55 2 Related Work: The Gap Between ML and Medical Experts

56 **Traditional Methods.** Traditional ML and NLP techniques have been widely adopted for language-
 57 based dementia detection tasks [26, 7]. Signal processing techniques are used to extract acoustic
 58 properties such as pitch, silent pauses, speech rate, harmonic-to-noise ratio and intensity from audio
 59 recordings [60, 26, 21]. Traditional NLP methods (e.g., part-of-speech tagging, syntactic constituent
 60 length) have been used to extract semantic, lexical, and syntactic properties from transcriptions of the
 61 speech [19, 41, 45]. Combined with the rise of traditional classifiers like Random Forest [21] and
 62 deep learning models such as convolutional neural networks (CNNs) to extract acoustic embeddings
 63 [26, 7], these approaches can help identify complex patterns in detecting cognitive impairment.
 64 Furthermore, feature importance (e.g., Random Forest, Support Vector Machines) can shed light
 65 on the features the model is relying on [40, 21]. However, while these techniques may achieve
 66 strong prediction accuracy, they are limited in their clinical interpretability. For example, [40] found

67 that Mel-Frequency Cepstral Coefficients (MFCCs) were the most important feature in dementia
68 prediction models, but critically MFCCs currently do not have clinical meaning. Interpretability
69 barriers and challenges with selecting clinically relevant features are important limitations in ML
70 dementia models’ ability to contribute in medical diagnostic processes (for more discussion on
71 explainable AI for medical applications, see [3] and [36]).

72 **LLM-Based Methods.** With their capability for language understanding, LLMs have introduced
73 a new wave of research in ML for dementia detection. Early works have utilized LLMs to extract
74 embeddings from transcribed speech for downstream classification [62, 48, 23, 35]. More recently,
75 LLMs have been explored through fine-tuning and transfer learning for direct inference of dementia
76 detection tasks [60, 61]. Many of these approaches have shown promise through high classification
77 accuracy, but face notable challenges such as bias [24] and hallucinations, where LLMs generate
78 outputs that sound plausible but are not grounded with accurate information [24, 13, 63]. Recent
79 efforts to develop medical-domain LLMs [57, 56, 27, 28, 51, 49], offer progress, but these models are
80 often trained on general clinical text and are not yet specialized in understanding human cognition
81 for dementia [63]. Currently, state-of-the-art LLMs also lack the embedded clinical knowledge
82 of established clinical diagnostic standards required to understand the dementia task [63, 6]. As a
83 result, they do not properly rely on medically validated markers of cognitive impairment [15, 55].
84 This causes LLMs to overlook subtle but critical linguistic features that are required to align model
85 cognition with clinical reasoning.

86 3 Position: Bridging the Gap

87 We argue that language-based dementia detection research should be re-framed to include cognition-
88 based reasoning, specifically by probing model cognition through structured, fine-grained input-output
89 analyses, which will enable clearer understanding of how ML models reason in this task. This can
90 be achieved by focusing on three key points: **(1) Probing Model Cognition to Better Understand**
91 **Input Cues:** systematically examining which linguistic features models rely on and how these
92 align with clinical reasoning. **(2) Moving Beyond Binary Classification to Progression-Aware**
93 **Modeling:** shifting from single-step prediction toward fine-grained, multi-stage outputs that reflect
94 the progressive nature of cognitive decline. **(3) Developing Fine-Grained Datasets for Deeper**
95 **Understanding:** developing richer, expert-annotated datasets that capture clinically meaningful
96 features to better support interpretable models that are clinically applicable. We believe that these
97 considerations will help design dementia detection ML models that are more aligned with human
98 cognition and clinically trustworthy.

99 3.1 Probing Alignment between Clinician and Model Cognition to Understand Input Cues

100 We propose that by systematically probing LLMs and ML models by testing feature importance
101 in the language-based dementia detection task, we can reveal what cues models actually rely on.
102 Clinicians pay attention to features in patients’ speech such as pauses, narrative coherence, speech
103 rate, word choice, and limited descriptions of images or recent events [58, 20, 37]. Recent work [63]
104 shows that structured probing in language-based dementia detection by varying prompt types (both
105 human-defined and LLM-generated) can reveal how different inputs affect model accuracy. They
106 demonstrate that by directing LLMs to focus on linguistic features such as places, objects, and actions,
107 and even generate intermediate steps in the detection process, we can understand which aspects of
108 cognition the model is focusing on. Building on this idea, we argue that systematically testing inputs
109 (e.g. masking speech features or removing demographic data), can help serve as interpretability
110 probes to test whether models perform detection through the same cues that clinicians focus on. In
111 parallel, recent trends in multi-agent systems, visual question answering (VQA) [10, 29, 18], and
112 human-in-the-loop [2, 59] show how LLMs can be guided through intermediate steps, either by
113 auxiliary models or by human experts, to probe what the model is focusing on when producing an
114 output. For example, in the case of a picture description task, the model can first analyze the input
115 image and then answer structured questions about the accompanying transcription (see Fig. 1). This
116 approach can enable us to directly test how well model cognition aligns with clinical reasoning, since
117 clinicians themselves can understand the intermediate reasoning steps linking the reference image,
118 patient speech, and model outputs. Understanding model cognition through careful input-output
119 analysis can help make the internal decision-making process of models more transparent and clinically
120 useful by showing cognitive reasoning and process.

121 3.2 Moving Beyond Binary Classification to Progression-Aware Modeling

122 To better align model cognition with clinical reasoning, we argue that language-based dementia
123 detection should move beyond end-to-end binary classification, to a continuum that is more consistent
124 with the neurodegenerative process. While many existing datasets include graded measures such as
125 MMSE scores, prior work commonly collapse these into binary categories (e.g. Healthy vs Dementia
126 or NC vs. MCI), losing clinically meaningful granularity. For example, the DAIC-WOZ dataset for
127 depression [17] uses Patient Health Questionnaire-8 (PHQ-8) scores ranging from 0 to 23 in the train
128 set to model severity, showing how scaled outputs can reveal deeper insights into model responses than
129 binary labels. This progression-aware framing for the dementia task will help produce more clinically
130 aligned predictions and also provide understanding of model cognition, showing whether internal
131 representations capture clinically meaningful details in a way clinicians may observe. Furthermore,
132 we suggest using LLMs to provide a clinician-friendly summary to convey these features to improve
133 interpretability and utility in making a clinical evaluation. Relevant demographic and contextual
134 information can then be incorporated into LLM summaries of the intermediate steps (e.g., "For an
135 81-year-old woman from Sacramento, California reading a passage her speech rate range, (1.8-2.3
136 syllables/second) is slower than average"). To support this framing, we propose building on existing
137 studies that leverage ML and NLP features by incorporating them into interpretability frameworks
138 that examine how models capture clinically meaningful cues semantic-pragmatic features, narrative
139 coherence, and acoustic features across the full spectrum from cognitively normal to dementia. Rather
140 than focusing solely on accuracy metrics such as accuracy, ROC AUC, or F-1 score, this combined
141 approach enables us to examine how models weigh specific cues across the NC-MCI-AD spectrum.
142 Techniques such as reinforcement learning, fine-tuning of LLMs, and Chain-of-Thought prompting
143 can further help reveal how models understand cognition and can be applied to focus on clinically
144 relevant features and structured reasoning steps in diagnosis.

145 3.3 Annotating Datasets with Clinically Relevant Features for Deeper Understanding

146 Progress in dementia detection from audio is also limited in the availability of datasets due to the
147 sensitive nature of medical data. Most publicly available datasets of this modality contain only
148 basic demographic information (e.g., sex, gender), clinical scores such as MMSE, and raw audio
149 recordings, and are only sometimes accompanied by ground truth transcriptions [33]. While this
150 minimal data is valuable, it lacks the depth needed for models to understand the nuanced reasoning of
151 medical experts. We believe that significant advances could be made with the addition of fine-grain
152 annotations in datasets, including annotating existing datasets with expert-labeled features such as
153 markers of emotion in voice, semantic-pragmatic cues, lexico-syntactic features, and other fine-
154 grained indicators a medical professional might consider during diagnosis. These enriched datasets
155 would not only support model training but also serve as interpretability probes, which will enable
156 exploration on whether models observe the same clinically meaningful signals that guide expert
157 reasoning. In this way, dataset development can become integral in understanding model cognition,
158 which further provides the foundation for interpretable and trustworthy dementia detection models.

159 4 Conclusion: Call for a Fresh Outlook

160 In this work, we argue that despite advances in applications of ML models for language-based demen-
161 tia detection, prior techniques have largely optimized for accuracy at the expense of interpretability
162 and alignment with clinical reasoning. Current classification techniques tend to lose the intermediate
163 cognitive assessment steps that human clinicians rely on, which limits the reliability of LLM usage in
164 clinical settings. We assert the need to re-frame dementia detection as a cognitive inference process –
165 one that systematically probes model cognition to understand what linguistic cues models focus on
166 and how these align with clinical reasoning. Proper integration of ML frameworks in real-world
167 clinics will require fine-grained, expert-annotated datasets that capture clinically meaningful features,
168 guiding models to focus on cues that matter most in practice. By combining these ideas with the
169 strengths of existing approaches, we can develop diagnostic tools that are more clinically trustworthy
170 and interpretable.

171 **Impact Statement:** This position paper calls for a fundamental shift in language-based dementia
172 detection research toward probing model cognition through structured, fine-grained input-output
173 analyses, enabling clearer understanding of how ML systems reason in this task and paving the way
174 for more interpretable, clinically trustworthy models. See Appendix C for a discussion on limitations.

References

- 175
- 176 [1] Randa Ben Ammar and Yassine Ben Ayed. Speech processing for early alzheimer disease
177 diagnosis: machine learning based approach. In *2018 IEEE/ACS 15th International Conference*
178 *on Computer Systems and Applications (AICCSA)*, pages 1–8. IEEE, 2018.
- 179 [2] Suzanne Bakken. Ai in health: keeping the human in the loop, 2023.
- 180 [3] Shahab S Band, Atefeh Yarahmadi, Chung-Chian Hsu, Meghdad Biyari, Mehdi Sookhak, Rasoul
181 Ameri, Iman Dehzangi, Anthony Theodore Chronopoulos, and Huey-Wen Liang. Application
182 of explainable artificial intelligence in medical health: A systematic review of interpretability
183 methods. *Informatics in Medicine Unlocked*, 40:101286, 2023.
- 184 [4] Benjamin Barrera-Altuna, Daeun Lee, Zaima Zarnaz, Jinyoung Han, and Seungbae Kim. The
185 interspeech 2024 taukadiial challenge: Multilingual mild cognitive impairment detection with
186 multimodal approach. *group*, 25:26, 2024.
- 187 [5] James T Becker, François Boiler, Oscar L Lopez, Judith Saxton, and Karen L McGonigle. The
188 natural history of alzheimer’s disease: description of study cohort and accuracy of diagnosis.
189 *Archives of neurology*, 51(6):585–594, 1994.
- 190 [6] Mondher Bouazizi, Chuheng Zheng, Siyuan Yang, and Tomoaki Ohtsuki. Dementia detection
191 from speech: What if language models are not the answer? *Information*, 15(1):2, 2023.
- 192 [7] Karol Chlasta and Krzysztof Wołk. Towards computer-based automated screening of dementia
193 through spontaneous speech. *Frontiers in Psychology*, 11:623237, 2021.
- 194 [8] Hongyoon Choi, Kyong Hwan Jin, Alzheimer’s Disease Neuroimaging Initiative, et al. Predict-
195 ing cognitive decline with deep learning of brain metabolism and amyloid imaging. *Behavioural*
196 *brain research*, 344:103–109, 2018.
- 197 [9] Natasha Clarke, Thomas R Barrick, and Peter Garrard. A comparison of connected speech tasks
198 for detecting early alzheimer’s disease and mild cognitive impairment using natural language
199 processing and machine learning. *Frontiers in Computer Science*, 3:634360, 2021.
- 200 [10] Fuze Cong, Shibiao Xu, Li Guo, and Yinbing Tian. Caption-aware medical vqa via semantic
201 focusing and progressive cross-modality comprehension. In *Proceedings of the 30th ACM*
202 *International Conference on Multimedia*, pages 3569–3577, 2022.
- 203 [11] Patricia V Cooper. Discourse production and normal aging: Performance on oral picture
204 description tasks. *Journal of gerontology*, 45(5):P210–P214, 1990.
- 205 [12] Kewen Ding, Madhu Chetty, Azadeh Noori Hoshyar, Tanusri Bhattacharya, and Britt Klein.
206 Speech based detection of alzheimer’s disease: a survey of ai techniques, datasets and challenges.
207 *Artificial Intelligence Review*, 57(12):325, 2024.
- 208 [13] Chengfeng Dou, Ying Zhang, Yanyuan Chen, Zhi Jin, Wenpin Jiao, Haiyan Zhao, and Yu Huang.
209 Detection, diagnosis, and explanation: A benchmark for chinese medial hallucination evaluation.
210 In *Proceedings of the 2024 Joint International Conference on Computational Linguistics,*
211 *Language Resources and Evaluation (LREC-COLING 2024)*, pages 4784–4794, 2024.
- 212 [14] Xinyang Feng, Zachary C Lipton, Jie Yang, Scott A Small, Frank A Provenzano, Alzheimer’s
213 Disease Neuroimaging Initiative, Frontotemporal Lobar Degeneration Neuroimaging Initiative,
214 et al. Estimating brain age based on a uniform healthy population with deep learning and
215 structural magnetic resonance imaging. *Neurobiology of aging*, 91:15–25, 2020.
- 216 [15] Ravindra Kumar Garg, Vijeth L Urs, Akshay Anand Agarwal, Sarvesh Kumar Chaudhary,
217 Vimal Paliwal, and Sujita Kumar Kar. Exploring the role of chatgpt in patient care (diagnosis
218 and treatment) and medical research: A systematic review. *Health Promotion Perspectives*, 13
219 (3):183, 2023.
- 220 [16] Harold Goodglass, Edith Kaplan, and Sandra Weintraub. *BDAE: The Boston diagnostic aphasia*
221 *examination*. Lippincott Williams & Wilkins Philadelphia, PA, 2001.

- 222 [17] Jonathan Gratch, Ron Artstein, Gale Lucas, Giota Stratou, Stefan Scherer, Angela Nazarian,
223 Rachel Wood, Jill Boberg, David DeVault, Stacy Marsella, David Traum, Skip Rizzo, and Louis-
224 Philippe Morency. The distress analysis interview corpus of human and computer interviews.
225 In Nicoletta Calzolari, Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard,
226 Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of*
227 *the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages
228 3123–3128, Reykjavik, Iceland, May 2014. European Language Resources Association (ELRA).
229 URL <https://aclanthology.org/L14-1421/>.
- 230 [18] Kaisi Guan, Zhengfeng Lai, Yuchong Sun, Peng Zhang, Wei Liu, Kieran Liu, Meng Cao, and
231 Ruihua Song. Etva: Evaluation of text-to-video alignment via fine-grained question generation
232 and answering. *arXiv preprint arXiv:2503.16867*, 2025.
- 233 [19] Melisa Gumus, Morgan Koo, Christa M. Studzinski, Aparna Bhan, Jessica Robin, and San-
234 dra E. Black. Linguistic changes in neurodegenerative diseases relate to clinical symp-
235 toms. *Frontiers in Neurology*, 15, March 2024. ISSN 1664-2295. doi: 10.3389/fneur.2024.
236 1373341. URL [https://www.frontiersin.org/journals/neurology/articles/10.](https://www.frontiersin.org/journals/neurology/articles/10.3389/fneur.2024.1373341/full)
237 [3389/fneur.2024.1373341/full](https://www.frontiersin.org/journals/neurology/articles/10.3389/fneur.2024.1373341/full).
- 238 [20] Phillip Hamrick, Victoria Sanborn, Rachel Ostrand, and John Gunstad. Lexical speech features
239 of spontaneous speech in older persons with and without cognitive impairment: reliability
240 analysis. *JMIR aging*, 6(1):e46483, 2023.
- 241 [21] Lior Hason and Sri Krishnan. Spontaneous speech feature analysis for alzheimer’s disease
242 screening using a random forest classifier. *Frontiers in Digital Health*, 4:901419, 2022.
- 243 [22] Emtiaz Hussain, Mahmudul Hasan, Syed Zafrul Hassan, Tanzina Hassan Azmi, Md Anisur Rah-
244 man, and Mohammad Zavid Parvez. Deep learning based binary classification for alzheimer’s
245 disease detection using brain mri images. In *2020 15th IEEE Conference on Industrial Elec-*
246 *tronics and Applications (ICIEA)*, pages 1115–1120. IEEE, 2020.
- 247 [23] Zerine Jahan, Surbhi Bhatia Khan, and Mo Saraee. Early dementia detection with speech analysis
248 and machine learning techniques. *Discover Sustainability*, 5(1):65, 2024.
- 249 [24] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
250 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation.
251 *ACM computing surveys*, 55(12):1–38, 2023.
- 252 [25] Edith Kaplan. *Boston diagnostic aphasia examination booklet*. Lea & Febiger Philadelphia, PA,
253 1983.
- 254 [26] M Rupesh Kumar, Susmitha Vekkot, S Lalitha, Deepa Gupta, Varasiddhi Jayasuryaa Govindraj,
255 Kamran Shaukat, Yousef Ajami Alotaibi, and Mohammed Zakariah. Dementia detection from
256 speech using machine learning and deep learning architectures. *Sensors*, 22(23):9311, 2022.
- 257 [27] Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier,
258 and Richard Dufour. Biomistral: A collection of open-source pretrained large language models
259 for medical domains. *arXiv preprint arXiv:2402.10373*, 2024.
- 260 [28] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan
261 Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision
262 assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36:
263 28541–28564, 2023.
- 264 [29] Yunyi Liu, Zhanyu Wang, Dong Xu, and Luping Zhou. Q2atransformer: Improving medical
265 vqa via an answer querying decoder. In *International conference on information processing in*
266 *medical imaging*, pages 445–456. Springer, 2023.
- 267 [30] Saturnino Luz, Fasih Haider, Sofia de la Fuente Garcia, Davida Fromm, and Brian MacWhinney.
268 Alzheimer’s dementia recognition through spontaneous speech, 2021.

- 269 [31] Saturnino Luz, Fasih Haider, Davida Fromm, Ioulietta Lazarou, Ioannis Kompatsiaris, and Brian
270 MacWhinney. Multilingual alzheimer’s dementia recognition through spontaneous speech: a
271 signal processing grand challenge. In *ICASSP 2023-2023 IEEE International Conference on*
272 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–2. IEEE, 2023.
- 273 [32] Saturnino Luz, Sofia De La Fuente Garcia, Fasih Haider, Davida Fromm, Brian MacWhinney,
274 Alyssa Lanzi, Ya-Ning Chang, Chia-Ju Chou, and Yi-Chien Liu. Connected speech-based
275 cognitive assessment in chinese and english. *arXiv preprint arXiv:2406.10272*, 2024.
- 276 [33] Saturnino Luz, Fasih Haider, Davida Fromm, Ioulietta Lazarou, Ioannis Kompatsiaris, and Brian
277 MacWhinney. An overview of the adress-m signal processing grand challenge on multilingual
278 alzheimer’s dementia recognition through spontaneous speech. *IEEE open journal of signal*
279 *processing*, 5:738–749, 2024.
- 280 [34] Pranav Mahajan and Veeky Baths. Acoustic and language based deep learning approaches for
281 alzheimer’s dementia detection from spontaneous speech. *Frontiers in Aging Neuroscience*, 13:
282 623607, 2021.
- 283 [35] Amit Meghanani, Chandran Savithri Anoop, and Angarai Ganesan Ramakrishnan. Recognition
284 of alzheimer’s dementia from the transcriptions of spontaneous speech using fasttext and cnn
285 models. *Frontiers in Computer Science*, 3:624558, 2021.
- 286 [36] Munib Mesinovic, Peter Watkinson, and Tingting Zhu. Explainability in the age of large
287 language models for healthcare. *Communications Engineering*, 4(1):128, 2025.
- 288 [37] Shamila Nasreen, Morteza Rohanian, Julian Hough, and Matthew Purver. Alzheimer’s dementia
289 recognition from spontaneous speech using disfluency and interactional features. *Frontiers in*
290 *Computer Science*, 3:640669, 2021.
- 291 [38] Linda E Nicholas and Robert H Brookshire. A system for quantifying the informativeness and
292 efficiency of the connected speech of adults with aphasia. *Journal of Speech, Language, and*
293 *Hearing Research*, 36(2):338–350, 1993.
- 294 [39] Emma Nichols, Jaimie D Steinmetz, Stein Emil Vollset, Kai Fukutaki, Julian Chalek, Foad
295 Abd-Allah, Amir Abdoli, Ahmed Abualhasan, Eman Abu-Gharbieh, Tayyaba Tayyaba Akram,
296 et al. Estimation of the global prevalence of dementia in 2019 and forecasted prevalence in
297 2050: an analysis for the global burden of disease study 2019. *The Lancet Public Health*, 7(2):
298 e105–e125, 2022.
- 299 [40] Marko Niemelä, Mikaela von Bonsdorff, Sami Äyrämö, and Tommi Kärkkäinen. Dementia
300 classification using acoustic speech and feature selection. *arXiv preprint arXiv:2502.03484*,
301 2025.
- 302 [41] Rachel Ostrand and John Gunstad. Using automatic assessment of speech production to predict
303 current and future cognitive function in older adults. *Journal of Geriatric Psychiatry and Neu-*
304 *rology*, 34(5):357–369, September 2021. ISSN 0891-9887. doi: 10.1177/0891988720933358.
305 URL <https://doi.org/10.1177/0891988720933358>.
- 306 [42] Dong Jin Park, Min Woo Park, Homin Lee, Young-Jin Kim, Yeongsic Kim, and Young Hoon
307 Park. Development of machine learning model for diagnostic disease prediction based on
308 laboratory tests. *Scientific reports*, 11(1):7567, 2021.
- 309 [43] Zhao Pei, Zhiyang Wan, Yanning Zhang, Miao Wang, Chengcai Leng, and Yee-Hong Yang.
310 Multi-scale attention-based pseudo-3d convolution neural network for alzheimer’s disease
311 diagnosis using structural mri. *pattern recognition*, 131:108825, 2022.
- 312 [44] Paula Andrea Pérez-Toro, Tomás Arias-Vergara, Philipp Klumpp, Tobias Weise, Maria Schuster,
313 Elmar Noeth, Juan Rafael Orozco-Aroyave, and Andreas Maier. Multilingual speech and
314 language analysis for the assessment of mild cognitive impairment: Outcomes from the taukadiadial
315 challenge. In *Proc. Interspeech 2024*, pages 982–986, 2024.

- 316 [45] Aurélie Pistono, Mélanie Jucla, Emmanuel J. Barbeau, Laure Saint-Aubert, Béatrice Lemesle,
317 Benjamin Calvet, Barbara Köpke, Michèle Puel, and Jérémie Pariente. Pauses during autobio-
318 graphical discourse reflect episodic memory processes in early alzheimer’s disease. *Journal of*
319 *Alzheimer’s disease: JAD*, 50(3):687–698, 2016. ISSN 1875-8908. doi: 10.3233/JAD-150408.
- 320 [46] Alvin Rajkomar, Eyal Oren, Kai Chen, Andrew M Dai, Nissan Hajaj, Michaela Hardt, Peter J
321 Liu, Xiaobing Liu, Jake Marcus, Mimi Sun, et al. Scalable and accurate deep learning with
322 electronic health records. *NPJ digital medicine*, 1(1):18, 2018.
- 323 [47] N Rockwell. Going and coming (oil on canvas). *Norman Rockwell Art Collection Trust*,
324 *Indianapolis, IN, United States*, 1947.
- 325 [48] Yamanki Santander-Cruz, Sebastián Salazar-Colores, Wilfrido Jacobo Paredes-García, Hum-
326 berto Guendulain-Arenas, and Saúl Tovar-Arriaga. Semantic feature extraction using sbert for
327 dementia detection. *Brain sciences*, 12(2):270, 2022.
- 328 [49] Andrew Sellergren, Sahar Kazemzadeh, Tiam Jaroensri, Atilla Kiraly, Madeleine Traverse,
329 Timo Kohlberger, Shawn Xu, Fayaz Jamil, Cían Hughes, Charles Lau, et al. Medgemma
330 technical report. *arXiv preprint arXiv:2507.05201*, 2025.
- 331 [50] Zehra Shah, Jeffrey Sawalha, Mashrura Tasnim, Shi-ang Qi, Eleni Stroulia, and Russell Greiner.
332 Learning language and acoustic models for identifying alzheimer’s dementia from speech.
333 *Frontiers in Computer Science*, 3:624659, 2021.
- 334 [51] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou,
335 Kevin Clark, Stephen R Pfohl, Heather Cole-Lewis, et al. Toward expert-level medical question
336 answering with large language models. *Nature Medicine*, 31(3):943–950, 2025.
- 337 [52] Xuegang Song, Feng Zhou, Alejandro F Frangi, Jiuwen Cao, Xiaohua Xiao, Yi Lei, Tianfu
338 Wang, and Baiying Lei. Multicenter and multichannel pooling gcn for early ad diagnosis based
339 on dual-modality fused brain network. *IEEE Transactions on Medical Imaging*, 42(2):354–367,
340 2022.
- 341 [53] Muhammad Shehram Shah Syed, Zafi Sherhan Syed, Margaret Lech, and Elena Pirogova.
342 Automated screening for alzheimer’s dementia through spontaneous speech. In *Interspeech*,
343 volume 2020, pages 2222–6, 2020.
- 344 [54] Khandaker Mohammad Mohi Uddin, Mir Jafikul Alam, Md Ashraf Uddin, and Sunil Aryal.
345 A novel approach utilizing machine learning for the early diagnosis of alzheimer’s disease.
346 *Biomedical Materials & Devices*, 1(2):882–898, 2023.
- 347 [55] Harriet Louise Walker, Shahi Ghani, Christoph Kueimmerli, Christian Andreas Nebiker, Beat Pe-
348 ter Müller, Dimitri Aristotle Raptis, and Sebastian Manuel Staubli. Reliability of medical
349 information provided by chatgpt: assessment against clinical guidelines and patient information
350 quality instrument. *Journal of Medical Internet Research*, 25:e47479, 2023.
- 351 [56] Xidong Wang, Nuo Chen, Junyin Chen, Yidong Wang, Guorui Zhen, Chunxian Zhang, Xiangbo
352 Wu, Yan Hu, Anningzhe Gao, Xiang Wan, et al. Apollo: A lightweight multilingual medical
353 llm towards democratizing medical ai to 6b people. *arXiv preprint arXiv:2403.03640*, 2024.
- 354 [57] Qianqian Xie, Qingyu Chen, Aokun Chen, Cheng Peng, Yan Hu, Fongci Lin, Xueqing Peng,
355 Jimin Huang, Jeffrey Zhang, Vipina Keloth, et al. Me-llama: Foundation large language models
356 for medical applications. *Research square*, pages rs–3, 2024.
- 357 [58] Anthony Yeung, Andrea Iaboni, Elizabeth Rochon, Monica Lavoie, Calvin Santiago, Maria
358 Yancheva, Jekaterina Novikova, Mengdan Xu, Jessica Robin, Liam D Kaufman, et al. Correl-
359 ating natural language processing and automated speech analysis with clinician assessment
360 to quantify speech-language changes in mild cognitive impairment and alzheimer’s dementia.
361 *Alzheimer’s research & therapy*, 13(1):109, 2021.
- 362 [59] Han Yuan, Lican Kang, Yong Li, and Zhenqian Fan. Human-in-the-loop machine learning for
363 healthcare: current progress and future opportunities in electronic health records. *Medicine*
364 *Advances*, 2(3):318–322, 2024.

- 365 [60] Jiahong Yuan, Yuchen Bian, Xingyu Cai, Jiayi Huang, Zheng Ye, and Kenneth Church. Dis-
366 fluencies and fine-tuning pre-trained language models for detection of alzheimer’s disease. In
367 *Interspeech*, volume 2020, pages 2162–6, 2020.
- 368 [61] Youxiang Zhu, Xiaohui Liang, John A Batsis, and Robert M Roth. Exploring deep transfer
369 learning techniques for alzheimer’s dementia detection. *Frontiers in computer science*, 3:
370 624683, 2021.
- 371 [62] Youxiang Zhu, Abdelrahman Obyat, Xiaohui Liang, John A Batsis, and Robert M Roth. Wavbert:
372 Exploiting semantic and non-semantic speech using wav2vec and bert for dementia detection.
373 In *Interspeech*, volume 2021, page 3790, 2021.
- 374 [63] Youxiang Zhu, Nana Lin, Kiran Balivada, Daniel Haehn, and Xiaohui Liang. Adversarial text
375 generation using large language models for dementia detection. In *Proceedings of the 2024*
376 *Conference on Empirical Methods in Natural Language Processing*, pages 21918–21933, 2024.

377 **A Alzheimer’s/Dementia Datasets**

378 Picture description datasets are widely used in ML research for dementia detection. Participants are
 379 asked to describe reference/stimuli images such as the well-known “Cookie Theft” picture, or in some
 380 cases multiple images like “Cat Rescue” and “Coming and Going” (see Fig. 2). Table 1 summarizes
 381 key datasets and subject counts. Among these, ADReSSo and TAUkADIAL have become especially
 382 popular datasets, with the former derived from the Pitt corpus.



Figure 2: Stimuli/reference images used in the ADReSS [30] and TAUkADIAL [32] datasets: “Cookie Theft” [25], “Cat Rescue” [38], and “Coming and Going” [47]. ADReSS includes only the Cookie Theft task, while TAUkADIAL requires each participant to describe all three images.

Dataset	Language	# of Subjects			
		NC	MCI	AD	Total
Pitt [5]	English	243	-	306	549
ADReSSo [30]	English	115	-	122	237
ADReSS-M [31]	English	115	-	122	237
	Greek	24	-	22	46
TAUKADIAL [32]	English	31	51	-	82
	Chinese	43	44	-	87

Table 1: Overview of existing dementia classification datasets for the picture description task. ADReSSo and TAUkADIAL are widely used datasets [26, 21, 34, 50]. Pitt Corpus and ADReSS rely solely on the Cookie Theft image, while TAUkADIAL employs three images (see Fig. 2), providing three transcripts per subject. ADReSS is a curated subset of the Pitt Corpus. NC: Normal Cognition, MCI: Mild Cognitive Impairment, AD: Alzheimer’s Disease.

383 **ADReSSo.** The ADReSSo [30] dataset contains speech samples collected from patients diagnosed
 384 with *Probable Alzheimer’s Disease (AD)* and *non-AD controls*. It captures both conversational
 385 and task-oriented speech, with a focus on distinguishing acoustic, lexical, semantic, and syntactic
 386 markers of Alzheimer’s progression. The dataset provides several hours of audio along with manual
 387 transcripts, allowing multimodal dementia detection. Fig. 3 shows an example for the two labels,
 388 Probable Alzheimer’s Disease (AD) and control.

389 **TAUKADIAL.** The TAUkADIAL [32] dataset consists of speech samples from individuals with
 390 *Mild Cognitive Impairment (MCI)* and *Normal Cognition (NC)* controls, in both English and Chinese.
 391 Participants completed the picture description task by answering questions about three stimuli images
 392 (see Fig. 2), providing three transcripts per subject. Fig. 4 shows example transcripts for the two
 393 labels, MCI and normal control.



Figure 3: Example samples from the ADReSSo [30] dataset, showing a **probable AD** patient (top) and a cognitively normal **control** (bottom). Speaker roles (patient vs. clinician) are not distinguished in the dataset; they are presented here for readability.

394 B Fine-Grain Feature Probing for Interpretability

395 As discussed in Sec. 3, we argue that dementia detection should be repositioned from a direct
 396 classification problem to a structured cognitive reasoning process. The picture description task is
 397 well-suited for fine-grained analysis, since transcripts can be broken down into clinically meaningful
 398 cues such as word repetitions, story coherence, image alignment, narrative flow, etc. Fig. 5 describes
 399 some of these fine-grain prompting categories that can be used to assess model cognition. Fig. 6
 400 shows how structured prompting can extract these features explicitly, in contrast to vague outputs
 401 from zero-shot prompting. Standardizing such fine-grained extractions, and providing ground-truth
 402 annotations from clinicians, would allow future models to reason through intermediate features
 403 before arriving at a diagnosis. This process mirrors clinical assessment more closely and offers
 404 interpretability by showing intermediate reasoning for features.

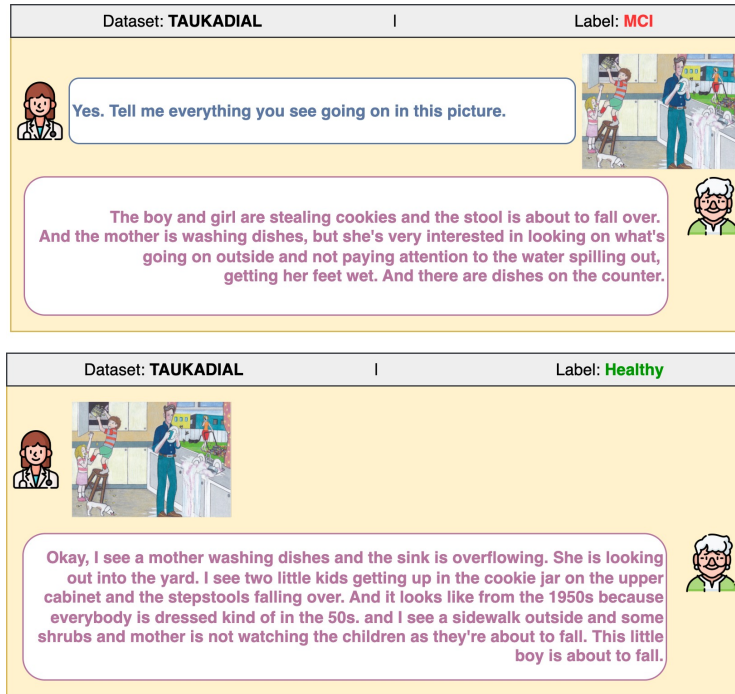


Figure 4: Example samples from the TAUADIAL [32] dataset, showing an **MCI** patient (top) and a cognitively **normal** control (bottom). Transcripts are ASR-generated; speaker roles are not always distinguished, and clinician speech may be absent.

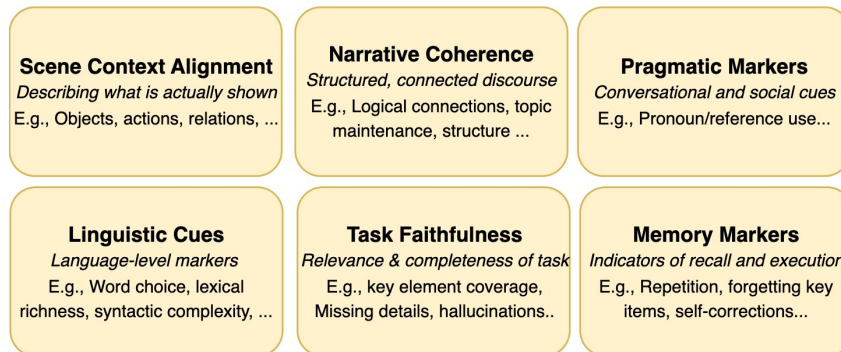


Figure 5: Examples of proposed structured probing tasks applied to dementia models, showing how varying prompt types can reveal attention to fine grain linguistic features and aid explainability.

405 C Limitations and Ethical Considerations

406 **Dataset sizes and lack of diversity.** While picture description datasets have enabled significant
 407 progress, their small scale and limited demographic diversity restrict generalizability. Most corpora
 408 are dominated by English-speaking participants or a small handful of languages, and subject pools
 409 underrepresent age, and socioeconomic variation, potentially embedding systematic biases into
 410 downstream models. Individual recordings are often very short, offering only brief snapshots
 411 of patient speech that limit the richness of cognitive assessment. Finally, many datasets reflect
 412 institutional biases in recruitment and annotation practices, which can influence both the linguistic
 413 style of collected samples and the diagnostic framing provided to models.

414 **Privacy and consent issues.** The sensitive nature of speech data also raises privacy and consent
 415 concerns. Audio might contain personal information that, if misused, could expose individuals beyond



Figure 6: **Top:** patient transcript to be analyzed for dementia. **Middle:** zero-shot prompting with a general-purpose MLLM (ChatGPT-5) and a medical MLLM (MedGemma-4B) produces incorrect diagnoses with vague explanations. **Bottom:** structured prompting elicits explicit feature extraction and transcript analysis, which can be used in a manner similar to clinical notes to support more fine-grained and interpretable diagnosis.

416 their cognitive status. Responsible use requires strict supervision for anonymization, controlled
 417 sharing, and adherence to medical data governance frameworks.

418 **LLMs and Hallucinations.** Reliance on LLMs introduces risks of hallucination and misplaced
 419 trust. Current models may generate reasonable but clinically invalid explanations. Without grounding
 420 in clinician-validated features, hallucinated outputs could reinforce overconfidence and inaccurate
 421 information in automated dementia detection models. Our position emphasizes that future work
 422 must treat LLMs as tools for structured analysis, not as substitutes for clinical judgment, ensuring
 423 interpretability and alignment with real-world diagnostic reasoning.