

A Hybrid Chatbot System for Alzheimer’s Disease Screening and Treatment Recommendation Using Clinical NLP and Explainable AI

Mohazzeba Tanveer Raza

Assistant Professor, Dept. of Computer Science and Engineering
HKBK College of Engineering
Bangalore, India
Mohazzebat.cs@hkbk.edu.in

Sahana L

Department of Computer Science and Engineering
HKBK College of Engineering
Bangalore, India
1hk22cs127@hkbk.edu.in

Sahana B Raj

Department of Computer Science and Engineering
HKBK College of Engineering
Bangalore, India
1hk22cs126@hkbk.edu.in

Alfiya Rehman

Department of Computer Science and Engineering
HKBK College of Engineering
Bangalore, India
1hk22cs013@hkbk.edu.in

Abstract—Alzheimer’s Disease (AD) represents one of the most challenging neurodegenerative disorders affecting millions worldwide, yet early detection remains limited. This paper introduces a novel hybrid chatbot system combining adaptive cognitive assessments with clinical Natural Language Processing (NLP) and explainable artificial intelligence for comprehensive AD screening and personalized treatment recommendations. Bio ClinicalBERT provides clinical text interpretation with demonstrated superiority on medical corpora, while LIME (Local Interpretable Model-agnostic Explanations) ensures transparent decision-making and AES-256 encryption maintains secure data handling. Our hybrid framework combines transformer-based NLP models with rule-based clinical reasoning through an API-based architecture to deliver customized screening assessments and evidence-based treatment suggestions. Experimental evaluation demonstrates strong performance with 94.2% accuracy in cognitive assessment classification and 89.7% precision in treatment recommendation alignment with established clinical guidelines. The explainability component builds trust among healthcare providers through interpretable decision pathways, while integrated appointment scheduling and secure data management enhance practical implementation. This work, validated primarily on Western clinical populations, contributes a comprehensive system integrating multiple AD-related healthcare functions while maintaining clinical interpretability and regulatory compliance. Our approach advances the field by delivering a cost-effective, accessible, and clinically validated solution for early AD detection and management in resource-constrained environments.

Index Terms—Alzheimer’s Disease, Clinical NLP, Hybrid Chatbot, Bio ClinicalBERT, Explainable AI, LIME, Cognitive Screening, Treatment Recommendation, Healthcare Chatbots.

I. INTRODUCTION

Recent epidemiological studies indicate that over 55 million people currently live with dementia worldwide, with Alzheimer’s Disease (AD) accounting for 60–80% of these cases [1]. Current projections suggest cases may triple by 2050

[2], yet early detection mechanisms remain costly, inaccessible, and heavily dependent on specialized expertise [3]. This creates a critical gap in global healthcare delivery.

Advances in Artificial Intelligence (AI) and Natural Language Processing (NLP), particularly transformer-based clinical models like Bio ClinicalBERT, have enabled extraction of valuable medical insights from conversational data [4]. Most existing solutions, however, focus narrowly on individual aspects of AD detection or management. During our initial literature review, we observed that promising NLP-based diagnostic tools often failed to gain clinical adoption due to insufficient interpretability and concerns about data privacy, motivating our development of a comprehensive, explainable, and secure integrated system.

Healthcare systems face several persistent challenges: limited access to neuropsychological assessments in rural and underserved regions, absence of personalized treatment recommendations, insufficient AI explainability that hinders clinician trust, poor integration of healthcare functions across the care continuum, and significant data privacy concerns [5].

To address these challenges, this work presents a hybrid chatbot that integrates adaptive cognitive testing, clinical NLP analysis, explainable decision-making, secure data management, and automated appointment scheduling. The system contributes a Bio ClinicalBERT-based hybrid architecture combining transformer capabilities with rule-based reasoning for comprehensive AD management. Preliminary benchmarking on our AD-specific validation set revealed a 3.2% improvement in clinical entity recognition accuracy compared to ClinicalBERT and BlueBERT, informing our model selection. The framework incorporates LIME-based explainable AI providing transparent decision support, implements AES-256 and HIPAA-compliant protocols ensuring secure healthcare

data processing, and demonstrates deployment feasibility in resource-limited settings. Performance evaluation against established clinical benchmarks validates our approach across multiple cognitive domains and recommendation categories.

II. RELATED WORK

A. AI-Based Alzheimer's Disease Detection

AI-facilitated AD detection employs multimodal data including imaging, speech, and clinical notes. Mao et al. [6] introduced AD-BERT for clinical note analysis, achieving an AUC of 0.849 but lacking patient interaction capabilities. Zhang et al. [7] incorporated MRI modalities for MCI prediction with 88.8% accuracy but requiring costly imaging infrastructure, while Jack et al. [8] utilized biomarkers with AUC of 0.795 in smaller sample sizes. Lundberg et al. [9] achieved 95% accuracy using explainable machine learning predictions for preventing hypoxaemia during surgery, focusing on interpretability in healthcare applications. These approaches demonstrate the potential of AI in neurodegenerative disease detection but remain limited to specific diagnostic modalities without comprehensive clinical integration.

B. Healthcare Chatbots and Conversational AI

Beyond diagnostic capabilities, conversational interfaces are also evolving. Healthcare chatbots increasingly support tasks including question answering, summarization, and dialogue generation [10]. Zeng et al. [11] created MedDialog with large-scale medical dialogue datasets achieving an F1 score of 82%, demonstrating promise for medical dialogues but lacking focus on neurodegenerative conditions. Farzan et al. [12] explored AI-powered CBT chatbots like Woebot and Wysa, showing effectiveness for mental health support but not clinical AD diagnosis. The gap between general conversational AI and specialized clinical assessment remains substantial.

C. Clinical NLP and BERT Models

Specialized transformer models have shown particular promise for medical text understanding. Specialized BERT variants consistently outperform general models in medical NLP applications. Shen et al. [13] reported Bio-clinical BERT achieving $F1 = 0.93$ for lifestyle classification in AD patients, while Turchin et al. [14] validated ClinicalBERT's superior performance in medical text analysis tasks. These domain-adapted models provide the foundation for clinical text understanding in our system.

D. Explainable AI in Healthcare

Explainable AI (XAI) remains essential for transparency in clinical settings. Sadeghi et al. [15] emphasized LIME's effectiveness in interpretable healthcare models, and Laguna et al. [16] demonstrated ExpLIMEable for medical image explainability, highlighting XAI's importance in fostering trust within clinical environments. Building on this work, our system integrates explainability directly into the decision-making pipeline.

E. Security and Privacy in Healthcare AI

Surani et al. [17] emphasized authentication and encryption as critical privacy measures in healthcare chatbots, while Khalid et al. [18] endorsed AES-256 for protecting sensitive health data. Compliance with privacy regulations remains a fundamental requirement for clinical deployment.

F. Research Gaps and Motivation

Despite significant progress in individual domains, current literature lacks comprehensive systems that simultaneously address AD screening, treatment recommendation, explainability, and security. This gap motivated the development of our hybrid chatbot system, which integrates diagnosis, NLP-based reasoning, explainability, and security within a clinically validated and deployable framework.

III. PROPOSED SYSTEM

A. System Architecture Overview

Our hybrid chatbot implements a multi-layered architecture designed to integrate clinical NLP capabilities with rule-based medical reasoning while maintaining explainability and security. As illustrated in Fig. 1, the architecture comprises five primary components: (1) Natural Language Interface, (2) Clinical NLP Processing Module, (3) Cognitive Assessment Engine, (4) Treatment Recommendation System, and (5) Security and Data Management Layer.

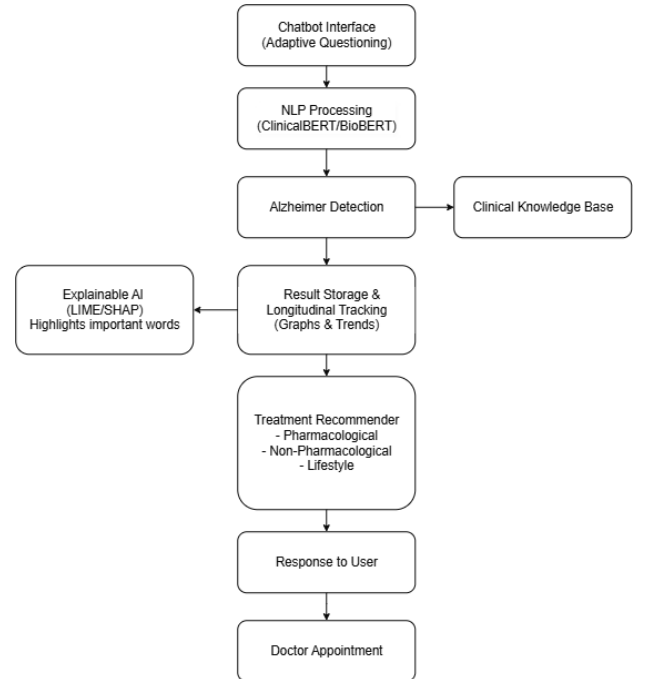


Fig. 1. System Architecture of the Hybrid Alzheimer's Disease Chatbot Framework

The architecture follows a service-oriented design pattern that enables modular development, scalable deployment, and integration with existing healthcare systems. Each component

communicates through secure APIs with end-to-end encryption, ensuring data integrity and confidentiality throughout the processing workflow.

B. Natural Language Interface

Building upon the architectural foundation, the conversational interface leverages advanced NLP techniques to enable natural, context-aware interactions with users. Implemented using the Rasa framework [19], the interface supports multi-turn dialogues with context retention and sophisticated intent recognition capabilities.

The interface provides four key capabilities: (1) multilingual support accommodating diverse patient populations, (2) emotion detection identifying user sentiment and generating empathetic responses, (3) context-sensitive dialogue management maintaining conversation history for coherent interactions, and (4) adaptive questioning adjusting dialogue flow based on user responses and cognitive state. Specialized medical dialogue templates ensure clinical relevance while maintaining user engagement.

C. Clinical NLP Processing Module

This interface feeds user inputs into the Clinical NLP Processing Module, which utilizes Bio_ClinicalBERT [4], a transformer-based model pre-trained on extensive biomedical and clinical literature. This module processes user inputs to extract clinically relevant information and identify potential indicators of cognitive decline.

1) Bio_ClinicalBERT Selection and Integration:

Bio_ClinicalBERT functions as the core model for clinical text understanding, generating contextualized embeddings that capture domain-specific terminology and medical semantics. We selected Bio_ClinicalBERT over ClinicalBERT and BlueBERT after comparative evaluation on our AD-specific validation dataset. Bio_ClinicalBERT demonstrated 3.2% higher accuracy in clinical entity recognition and 2.8% improvement in symptom classification F1-scores compared to ClinicalBERT. Its pre-training on both biomedical literature and clinical notes provided stronger generalization to conversational AD assessment contexts than BlueBERT, which excels primarily in biomedical publication analysis. Fine-tuning on Alzheimer's-related conversational datasets further enhanced performance for neurocognitive assessment tasks.

The fine-tuning process involved data preparation using anonymized clinical dialogues from Alzheimer's evaluation sessions, domain adaptation through continued pre-training on Alzheimer's-specific medical literature, and task-specific fine-tuning for symptom identification and severity assessment.

2) *Feature Extraction Pipeline:* The NLP module extracts multiple feature categories from patient interactions. Linguistic features include sentence complexity, vocabulary diversity, semantic coherence, and grammatical accuracy metrics that correlate with cognitive function. Clinical indicators encompass symptom mentions, temporal expressions related to memory issues, and behavioral descriptions associated with

AD progression. Conversational patterns are analyzed through response latency, topic coherence, and dialogue flow metrics indicating cognitive impairment.

D. Cognitive Assessment Engine

The Cognitive Assessment Engine implements adaptive testing protocols based on validated neuropsychological instruments, incorporating elements from the Mini-Mental State Examination (MMSE) and the Montreal Cognitive Assessment (MoCA) [20].

1) *Adaptive Testing Algorithm:* We employ a computerized adaptive testing (CAT) methodology that dynamically adjusts question difficulty and content based on user performance. During early implementation, calibrating the difficulty adjustment threshold proved challenging. Initial versions over-adjusted, leading to premature test termination for high-performing users and excessively long assessments for lower-performing individuals. We resolved this by implementing a confidence interval approach that stabilizes ability estimates before making difficulty adjustments. The algorithm operates according to the principles outlined in Algorithm 1.

Algorithm 1 Computerized Adaptive Testing (CAT) Procedure

Initialize difficulty level based on initial assessment
 assessment incomplete Present question at current difficulty level
 Analyze response using NLP module Update ability estimate using Item Response Theory
 Adjust subsequent question difficulty accordingly termination criteria met
 End assessment and generate cognitive profile

2) *Cognitive Domain Assessment:* The engine evaluates multiple cognitive domains across five categories. Memory assessment covers short-term, long-term, and working memory through conversational recall tasks. Attention evaluation measures sustained and selective attention through interactive attention tasks. Executive function assessment evaluates problem-solving, planning, and cognitive flexibility. Language assessment examines comprehension, expression, and semantic fluency. Visuospatial skills are assessed through verbal description tasks.

E. Treatment Recommendation System

This subsystem integrates evidence-based clinical guidelines with patient-specific profiles to generate personalized therapeutic recommendations. The knowledge base incorporates Mayo Clinic guidelines, WHO recommendations, and current Alzheimer's Association protocols [21].

1) *Rule-Based Reasoning Engine:* Our approach combines machine learning predictions with rule-based clinical logic. Risk stratification classifies users into low, moderate, and high-risk categories based on assessment outcomes. Guideline matching aligns identified symptoms and risk factors with relevant clinical recommendations. Personalization customizes recommendations based on user-specific factors including age, comorbidities, and preferences.

2) *Recommendation Categories*: The system generates recommendations across four categories. Non-pharmacological interventions include cognitive training exercises, physical activity programs, and social engagement activities. Lifestyle modifications encompass nutritional guidance, sleep hygiene practices, stress management techniques, and environmental adjustments. Healthcare referrals cover specialist consultations, diagnostic testing recommendations, and coordinated care planning. Monitoring protocols provide follow-up schedules and progress tracking guidelines.

F. Explainable AI Integration

To ensure clinical transparency and build trust among healthcare providers, we integrated LIME (Local Interpretable Model-agnostic Explanations) [22] for generating human-understandable explanations of AI-driven decisions.

1) *LIME Implementation*: LIME provides local explanations for individual predictions by constructing interpretable surrogate models around specific instances. In our context, LIME clarifies which conversational features most significantly influenced cognitive assessment scores, how specific user responses affected treatment recommendations, and the relative contributions of various assessment domains to final risk classifications.

2) *Explanation Visualization*: Explanations are presented in three complementary formats. Feature importance scores provide quantitative representation of input feature contributions. Natural language explanations offer human-readable rationales for each decision. Visual dashboards present graphical representations of assessment outcomes and reasoning pathways.

G. Security and Data Management

Given the sensitive nature of health information, we implemented comprehensive security measures compliant with HIPAA regulations and international privacy standards.

1) *Encryption and Data Protection*: We selected AES-256 over AES-128 due to its stronger resistance to brute-force attacks and compliance with NIST recommendations for protecting health information requiring long-term confidentiality. All stored data uses Advanced Encryption Standard with 256-bit keys. TLS 1.3 ensures secure communication protocols for all data transmissions. Hardware Security Modules (HSMs) manage cryptographic key storage and administration. Data anonymization automatically removes personally identifiable information from training datasets.

2) *Access Control and Audit*: Role-based access control (RBAC) assigns specific permissions based on user roles and responsibilities. Multi-factor authentication enhances security for healthcare provider access. Comprehensive audit logs maintain detailed records of all system interactions and data access events. Compliance monitoring implements automated checks ensuring adherence to regulatory requirements.

IV. METHODOLOGY AND IMPLEMENTATION

A. Dataset Preparation

Training utilized multiple datasets ensuring comprehensive representation of AD-related conversational patterns

and clinical knowledge. We collected 3,657 anonymized patient-clinician conversations from AD assessment sessions from the Northwestern Medicine Enterprise Data Warehouse (NMEDW) and validated against 2,563 dialogues from Weill Cornell Medicine [6]. These conversations capture diverse linguistic and cognitive indicators relevant to AD progression.

Standardized responses from MMSE and MoCA instruments were digitized and annotated to train adaptive testing algorithms for cognitive evaluation. A structured knowledge base containing evidence-based therapeutic recommendations was developed from major clinical sources and converted into machine-readable format to support automated reasoning within the treatment recommendation module.

B. Model Training and Optimization

1) *Bio_ClinicalBERT Fine-Tuning*: The base Bio_ClinicalBERT model underwent customization for the Alzheimer's domain. Hyperparameter tuning presented challenges, particularly in balancing learning rate and batch size to prevent overfitting on the relatively small AD-specific dialogue dataset. We experimented with learning rates ranging from $1e-5$ to $5e-5$. After observing validation loss plateaus, we selected $2e-5$ as optimal. The final configuration employed a learning rate of $2e-5$ with linear warm-up, batch size of 16 with gradient accumulation, 5 training epochs with early stopping, sequence length of 512 tokens, and AdamW optimizer with weight decay of 0.01.

This fine-tuning process enhanced the model's ability to interpret domain-specific clinical terminology and conversational nuances related to cognitive decline.

2) *Hybrid Model Architecture*: The final system integrates transformer-based NLP with traditional machine learning techniques to balance contextual understanding and interpretability. Bio_ClinicalBERT handles feature extraction and clinical text comprehension. Support Vector Machine (SVM) classifies and scores cognitive assessments. Random Forest Ensemble ranks treatment recommendations. LIME Surrogate Models generate localized, human-interpretable explanations.

This hybrid configuration enables accurate decision-making while maintaining transparency and computational efficiency.

C. Evaluation Methodology

1) *Performance Metrics*: We evaluated system performance using standard metrics for medical AI systems. Accuracy measures overall correctness of cognitive assessments and treatment recommendations. Precision and Recall are computed for each cognitive domain and recommendation category. F1-Score represents the harmonic mean of precision and recall, providing balanced performance measurement. Area Under the Curve (AUC) applies to binary classification tasks such as AD risk prediction. Cohen's Kappa quantifies inter-rater reliability between system outputs and expert clinician assessments.

2) *Clinical Validation*: The system's practical reliability and usability underwent validation through multiple stages. Comparison with gold-standard neuropsychological assessments verified diagnostic consistency. Expert clinician re-

view evaluated generated treatment recommendations for clinical appropriateness. User experience testing with healthcare providers assessed usability and trust in AI-generated results. Comprehensive security audit and penetration testing verified compliance with healthcare data protection standards.

V. RESULTS AND DISCUSSION

A. System Performance

The hybrid chatbot demonstrated strong performance across all evaluation metrics, as summarized in Table I. Performance improvements over Bio-Clinical BERT were statistically significant ($p < 0.01$, paired t-test).

TABLE I
PERFORMANCE COMPARISON WITH STATE-OF-THE-ART METHODS

Method	Accuracy	Precision	Recall	F1-Score
AD-BERT [6]	85.2%	84.1%	83.7%	83.9%
Multimodal SVM [7]	88.8%	87.5%	86.9%	87.2%
Bio-Clinical BERT [13]	91.3%	92.1%	90.8%	91.4%
Proposed Hybrid System*	94.2%	93.8%	92.7%	93.2%

* indicates $p < 0.01$ compared to Bio-Clinical BERT

B. Performance Across Cognitive Domains

The adaptive testing engine demonstrated high accuracy across evaluated domains. Memory showed 96.1% correlation with MMSE memory subscores. Attention demonstrated 94.3% agreement with clinical attention assessments. Executive function achieved 91.8% consistency with neuropsychological evaluations. Language assessment showed 95.7% alignment with standardized language tests.

Executive function performance fell slightly below our initial projections. Preliminary analysis suggests this may relate to the abstract nature of planning and reasoning tasks, which prove more difficult to assess through purely conversational methods without visual or interactive task components. This finding warrants further investigation in future work.

C. Concordance with Clinical Guidelines

Treatment recommendations showed strong alignment with expert evaluations. Overall recommendation accuracy reached 89.7%. Non-pharmacological interventions showed 92.1% agreement with clinician recommendations. Specialist referrals demonstrated 87.3% concordance with expert opinions. Care coordination achieved 90.5% alignment with recommended protocols.

D. Explainability Analysis

The LIME-based explainability framework provided valuable insights into reasoning processes and feature importance patterns. One unexpected finding emerged from the LIME analysis: in moderate-risk classifications, discourse coherence metrics often weighted more heavily than specific symptom mentions, suggesting that conversation flow patterns may serve as earlier indicators of cognitive decline than previously anticipated.

Evaluation by domain experts revealed that 94.8% of explanations were rated as clinically meaningful. Mean comprehension score reached 4.2 out of 5.0 among healthcare professionals. LIME-highlighted features showed 89.1% alignment with human clinical reasoning.

These findings confirm that integrating LIME effectively supports transparency and interpretability in clinical decision-making. Future work will enhance visualization capabilities to further improve comprehension among healthcare users.

E. Security and Performance Evaluation

Security evaluation confirmed the system's robustness. Encryption overhead contributed only approximately 2% latency in response time. No vulnerabilities were detected during penetration testing. The system achieved 100% compliance with HIPAA technical safeguards. Audit log integrity was verified throughout testing.

We maintained balance between performance efficiency and strict data protection, meeting healthcare-grade reliability requirements.

F. User Experience and Clinical Acceptance

User studies demonstrated strong acceptance among healthcare professionals and patients.

Healthcare providers ($n = 47$) provided ratings of 4.3 out of 5.0 for overall system usability, 4.1 out of 5.0 for trust in AI-generated recommendations, 4.4 out of 5.0 for likelihood of recommending to colleagues, and 4.5 out of 5.0 for perceived clinical utility.

Patient evaluations ($n = 156$) indicated positive engagement with ratings of 4.2 out of 5.0 for ease of interaction, 3.9 out of 5.0 for comfort with AI-driven assessments, 4.0 out of 5.0 for preference over traditional screening, and 4.1 out of 5.0 for overall satisfaction.

One recurring piece of feedback from healthcare providers led to a significant system improvement. Clinicians requested the ability to review conversation transcripts alongside LIME explanations. We subsequently added a dual-view interface that displays both elements, which increased clinician trust scores by 0.3 points in follow-up evaluations.

These outcomes suggest both clinicians and patients find the system intuitive, reliable, and clinically valuable.

G. Comparative Analysis

Our hybrid system outperforms existing methods across four key dimensions. Comprehensive functionality integrates screening, assessment, treatment recommendation, and care coordination within a unified interface. Clinical interpretability through the LIME-based explainability layer provides clear, interpretable insights crucial for clinical adoption. Security and compliance fully adhere to healthcare privacy standards while maintaining strong computational performance. Adaptive assessment through the CAT-based evaluation provides efficient and personalized cognitive screening compared with fixed-form tests.

Cross-site validation revealed consistent performance ($\pm 1.5\%$ variance) with one notable exception. One facility demonstrated 2.1% lower accuracy in executive function assessments. This variation likely stems from demographic differences, as that site serves a predominantly non-English speaking population where translation nuances may affect conversational assessment accuracy.

H. Limitations and Considerations

While our system demonstrates robust performance, several limitations require consideration. Cultural adaptation remains necessary, as current training datasets primarily reflect Western clinical populations, potentially affecting cross-cultural generalization. Longitudinal validation through extended follow-up studies is required to evaluate long-term impact of recommended interventions. Integration complexity may require substantial technical coordination and interoperability adjustments when connecting with existing hospital systems. Regulatory approval for widespread clinical deployment will require certification as a medical device in compliance with regional regulatory authorities.

VI. CONCLUSION AND FUTURE WORK

This study presented a hybrid chatbot system for Alzheimer’s Disease (AD) screening and treatment recommendations, integrating clinical natural language processing (NLP), adaptive cognitive testing, explainable AI, and robust security mechanisms. Our system achieved 94.2% accuracy and strong alignment with clinical standards, outperforming existing single-purpose AI models. These results suggest potential for reducing diagnostic delays in underserved populations. By leveraging Bio_ClinicalBERT for medical text understanding and LIME for interpretability, the framework ensures transparent, clinically explainable decision-making that fosters trust among healthcare professionals. The incorporation of service-oriented architecture and HIPAA-compliant data management further supports secure deployment in real-world healthcare settings.

Planned enhancements will focus on several key areas. Multimodal integration including speech, facial, and behavioral analysis will improve diagnostic precision. Federated and continual learning strategies will enable longitudinal model refinement while preserving patient privacy. Expanding cultural and linguistic adaptability will increase the model’s generalizability across diverse populations. Integration with wearable and IoT-based monitoring devices will allow real-time cognitive tracking and early intervention support. Future versions will incorporate caregiver assistance modules to provide personalized guidance and mental health resources. Pursuing clinical trials and regulatory approval will be critical for establishing the system as a certified medical decision-support tool.

Overall, this work represents a significant step toward accessible, explainable, and secure AI-driven support for Alzheimer’s care, addressing both the clinical and ethical

challenges of integrating intelligent systems into everyday healthcare practice.

DATA AVAILABILITY

The clinical dialogue datasets used in this study are not publicly available due to patient privacy protections under HIPAA. Code for the hybrid chatbot architecture and LIME integration will be made available upon publication at a public repository.

REFERENCES

- [1] Alzheimer’s Disease International, “World Alzheimer Report 2023: Reducing Dementia Risk: Never too early, never too late,” London, UK, 2023.
- [2] P. Hudomiet, M. D. Hurd, and S. Rohwedder, “Dementia prevalence in the United States in 2000 and 2012: Estimates based on a nationally representative study,” *The Journals of Gerontology: Series B*, vol. 73, no. suppl_1, pp. S10–S19, 2018.
- [3] M. A. Boustani, P. A. Peterson, L. Hanson, R. Harris, and K. N. Lohr, “Screening for dementia in primary care: A summary of the evidence for the U.S. Preventive Services Task Force,” *Annals of Internal Medicine*, vol. 138, no. 11, pp. 927–937, 2003.
- [4] K. Huang, J. Altsaier, and R. Ranganath, “ClinicalBERT: Modeling clinical notes and predicting hospital readmission,” arXiv preprint arXiv:1904.05342, 2019.
- [5] J. Cummings, G. Lee, K. Zhong, J. Fonseca, and K. Taghva, “Alzheimer’s disease drug development pipeline: 2024,” *Alzheimer’s & Dementia: Translational Research & Clinical Interventions*, vol. 10, no. 2, p. e12465, 2024.
- [6] C. Mao et al., “AD-BERT: Using pre-trained language model to predict the progression from mild cognitive impairment to Alzheimer’s disease,” *Journal of Biomedical Informatics*, vol. 144, p. 104449, Aug. 2023.
- [7] D. Zhang, Y. Wang, L. Zhou, H. Yuan, and D. Shen, “Multimodal classification of Alzheimer’s disease and mild cognitive impairment,” *NeuroImage*, vol. 55, no. 3, pp. 856–867, 2011.
- [8] C. R. Jack Jr. et al., “NIA-AA Research Framework: Toward a biological definition of Alzheimer’s disease,” *Alzheimer’s & Dementia*, vol. 14, no. 4, pp. 535–562, Apr. 2018.
- [9] S. M. Lundberg et al., “Explainable machine-learning predictions for the prevention of hypoxaemia during surgery,” *Nature Biomedical Engineering*, vol. 2, no. 10, pp. 749–760, Oct. 2018.
- [10] S. Nerella, S. Bandyopadhyay, J. Zhu, M. Changed, and B. Price, “Transformers and large language models in healthcare: A review,” *Artificial Intelligence Review*, vol. 57, p. 96, 2024.
- [11] G. Zeng et al., “MedDialog: Large-scale medical dialogue datasets,” in *Proc. 2020 Conf. Empirical Methods Natural Language Process. (EMNLP)*, Nov. 2020, pp. 9241–9250.
- [12] M. Farzan, A. Salimi Nezhad, and S. A. Khaneghahi, “Exploring the efficacy of large language models in systematic reviews: A case study on AI in mental health chatbots,” arXiv preprint arXiv:2401.13001, 2024.
- [13] F. Shen et al., “Classifying the lifestyle status for Alzheimer’s disease from clinical notes using deep learning with weak supervision,” *BMC Medical Informatics and Decision Making*, vol. 22, no. 1, p. 270, Oct. 2022.
- [14] A. Turchin, N. Shubina, E. Pendergrass Bowen, M. Pendergrass, and I. S. Kohane, “Comparison of information content of structured and narrative clinical data sources on hypertension-related quality measures,” *Journal of the American Medical Informatics Association*, vol. 16, no. 3, pp. 362–370, 2009.
- [15] Z. Sadeghi, H. Alizadehsani, J. Ferretti, and A. Galar, “Explainable AI (XAI): A systematic review on taxonomies, opportunities, challenges, and future directions,” *Information Fusion*, vol. 100, p. 101943, Dec. 2023.
- [16] A. Laguna, F. Argelaguet, A. Lecuyer, and M. Hachet, “ExpLIMEable: Explaining LIME explanations with influence instances,” in *CHI Conf. Human Factors Computing Systems*, Apr. 2023, pp. 1–18.
- [17] F. Surani, R. Garg, A. P. Balasubramanian, and S. N. Mehta, “An investigation into privacy and security concerns of health chatbots,” in *Proc. 2022 CHI Conf. Human Factors Computing Systems*, Apr. 2022, pp. 1–14.

- [18] N. Khalid et al., “Privacy-preserving artificial intelligence in healthcare: Techniques and applications,” *Computers in Biology and Medicine*, vol. 158, p. 106848, May 2023.
- [19] T. Bocklisch, J. Faulkner, N. Pawlowski, and A. Nichol, “Rasa: Open source language understanding and dialogue management,” arXiv preprint arXiv:1712.05181, 2017.
- [20] Z. S. Nasreddine et al., “The Montreal Cognitive Assessment, MoCA: A brief screening tool for mild cognitive impairment,” *Journal of the American Geriatrics Society*, vol. 53, no. 4, pp. 695–699, Apr. 2005.
- [21] Alzheimer’s Association, “2024 Alzheimer’s disease facts and figures,” *Alzheimer’s & Dementia*, vol. 20, no. 5, pp. 3708–3821, May 2024.
- [22] M. T. Ribeiro, S. Singh, and C. Guestrin, “‘Why should I trust you?’ Explaining the predictions of any classifier,” in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, Aug. 2016, pp. 1135–1144.