

Towards Integrated LLM-Driven and Ensemble Approaches for Early Alzheimer's Disease Detection: Synthesis of Advanced Research and Clinical Roadmaps

I. Executive Summary: LLMs as Strategic Assets in AD Detection

Early, accurate, and interpretable detection of Alzheimer's Disease (AD) remains a central challenge for clinical research and public health, as traditional diagnostic pathways are often reliant on resource-intensive methods like neuroimaging or subjective clinical assessments, frequently leading to diagnoses late in the disease progression.¹ The convergence of specialized Large Language Models (LLMs) and advanced multimodal Artificial Intelligence (AI) systems presents a compelling strategy to overcome these limitations by offering objective, scalable, and computationally efficient diagnostic tools.

This comprehensive synthesis report integrates three complementary lines of research into a cohesive, unified workflow designed to leverage the unique strengths of each component.¹ These pillars include: 1) high-performing stacked fusion ensemble models optimized for screening based on spoken and written language; 2) the CARE AD multi-agent LLM framework tailored for longitudinal risk stratification from Electronic Health Record (EHR) narratives; and 3) a multimodal LLM pipeline that extracts interpretable features from clinical task media, such as those related to the Mini-Mental State Examination (MMSE).¹

The effectiveness of this integrated approach is strongly substantiated by the exceptional performance metrics achieved by its constituent technologies. Specifically, the ensemble stacking methods have yielded remarkable diagnostic acuity, demonstrating AUCs consistently **above 0.98** on language-only screening datasets, with linked experiments

reporting AUCs reaching up to **~99.2%**.¹ The application of LLM frameworks in this domain signals a fundamental evolution: these tools are transitioning beyond simple binary AD/Not AD classification. Instead, they are becoming comprehensive clinical workflow accelerators and risk stratification platforms, encompassing Generative Reporting, Speech Screening, and Granular Phenotyping.¹ The combined system proposes a practical, unified, multi-tier strategy that efficiently moves patients from low-friction frontline screening to transparent in-clinic assessment, culminating in long-term EHR-based risk monitoring.¹

A necessary strategic consideration governing the design of this pipeline is the priority of clinical trust and transparency over raw predictive power. While highly accurate black-box models exist, the integrated pipeline deliberately incorporates features designed for interpretability, such as quantified MMSE metrics and clinician-facing explanations derived from decision trees and counterfactual explanations.¹ This structural decision underscores the recognition that maximum predictive accuracy is clinically insufficient without transparency and justification, making interpretability a non-negotiable prerequisite for integrating these powerful computational tools into regulated clinical workflows.

II. The AI Landscape of Alzheimer's Disease Diagnostics

2.1 The Clinical Imperative: Bridging the Gap in Early Detection

The demand for timely and objective diagnostic methods in AD is intensifying due to the global health crisis posed by the disease.¹ The convergence of deep learning, expansive clinical datasets, and advanced LLM architectures has catalyzed a paradigm shift in AD research, overcoming the limitations of traditional, subjective clinical assessments.¹ LLMs demonstrate high adaptability to complex medical challenges, processing nuanced clinical jargon in EHRs and synthesizing multimodal inputs for diagnosis.¹ Crucially, their power extends beyond analyzing language; LLMs act as central computational hubs capable of integrating diverse patient data.¹

2.2 Defining the Three Domains of LLM Application in AD

The synthesized research efforts align directly with three strategic domains identified in the broader literature regarding LLM applications in AD ¹:

1. **Generative and Multimodal Reporting:** These systems address diagnostic complexity and documentation inefficiency by ingesting heterogeneous data, such as fusing neuroimaging (3D MRI), genetic markers (SNPs), and EHR data to inform a central predictive model.¹ The predictive output is then used to constrain the LLM's output, often via a "Cognition-Guided Prompt," ensuring the generated human-readable clinical report is grounded in the calculated evidence, thereby minimizing the risk of "hallucination".¹
2. **Speech-Based Screening and Diagnostic Equity:** This domain leverages LLMs to analyze acoustic and linguistic markers of AD, providing a low-cost, non-invasive screening method.¹ High-performing systems combine linguistic features extracted from fine-tuned LLMs with acoustic features to create powerful multimodal fusion classifiers. The **Stacked Fusion** model utilizing lexicosyntactic features and n-grams from short speech and writing samples directly addresses this domain.¹
3. **Granular Clinical Phenotyping:** Specialized LLMs are highly effective at extracting fine-grained, quantitative symptom severity from unstructured clinical text, shifting the focus beyond binary diagnosis toward detailed clinical characterization.¹ The **CARE AD** multi-agent framework exemplifies this through the analysis of longitudinal clinical notes to quantify symptom profiles across multiple domains.¹

This current state of research suggests that the convergence of modalities is a crucial hallmark of maturity in AD diagnostics. Earlier research often focused on single modalities (e.g., speech or imaging in isolation). The necessity of sophisticated multimodal fusion (leveraging language, MMSE media, EHR text, genetics, and imaging) indicates that the field is moving past exploratory single-source studies toward mature, integrated systems that effectively mimic the multi-disciplinary nature of a comprehensive clinical diagnosis, where reliance on any single data point is insufficient. Furthermore, the functional versatility of the LLM is evident; it is employed not only for language analysis but also as a feature extractor in the multimodal pipeline and a data synthesizer (for synthetic report generation).¹ This versatility positions the LLM as a core computational utility capable of translating raw, complex clinical inputs (such as video or imaging data) into structured, quantifiable, and, critically, interpretable metrics.¹

III. Data Modalities and Next-Generation Feature Engineering

The power of the integrated pipeline stems from its ability to synthesize data from multiple heterogeneous sources and optimize diverse feature representations, balancing high-performance, low-level characteristics with clinically relevant, high-level indicators.

3.1 Diverse Data Inputs for Comprehensive Analysis

The analysis draws upon comprehensive data inputs used across the constituent studies¹:

- **Spoken and Written Language:** The Cookie Theft Picture Corpus (CTPC), a manually transcribed spoken-description dataset, and the Alzheimer's Disease Blog Corpus (ADBC), written posts from individuals with AD and healthy controls.¹ These sources are used for extracting lexicosyntactic and n-gram features.
- **Longitudinal Clinical Records (EHRs):** VHA longitudinal clinical notes, which span many years and are segmented and filtered to identify clinically relevant encounters for automated extraction and annotation.¹
- **Multimodal MMSE Media:** Media gathered during MMSE-like tasks, including short videos (e.g., paper folding, eye-closing), images (e.g., pentagon drawing, handwriting samples), and recorded speech.¹

3.2 Comprehensive Feature Taxonomy: Analysis and Optimization

The models rely on a strategic combination of feature sets¹:

- **Lexicosyntactic Features:** A broad feature set capturing lexical richness and syntactic complexity, including character, word, and sentence counts, average word and sentence length, type-token ratio (TTR), idea density, content density, and proposition densities (active vs. passive).¹ These features are engineered to be interpretable and are optimized using robust scaling (e.g., RobustScaler) and Correlation-Based Feature Selection (CFS) to drop redundant features and reduce dimensionality.¹
- **Character N-gram Spaces:** Low-level character bigrams and trigrams (after stopword removal) are utilized to capture stylistic and morphological patterns. These proved effective as complementary features, particularly when analyzing short texts and noisy data, such as blog posts.¹
- **EHR-derived Symptom Timelines (CARE AD):** A data-extraction agent processes longitudinal clinical notes to identify AD-relevant sentences, assigning them to

expert-defined categories: cognitive impairment, notice/concern by others, functional decline, physiological changes, and neuropsychiatric symptoms. The outputs are aggregated into age-aware longitudinal symptom profiles.¹

- **Multimodal MMSE-Derived Metrics:** The multimodal LLM pipeline translates raw task media into compact numeric and categorical scores. These scores are explicitly designed for clinician review and include metrics such as *MotorCoordination*, *FollowingInstruction*, *ShapeCompleteness*, *ResponseLatency*, *Legibility*, *Speech Fluency*, *Vocabulary Richness*, and *Coherence*.¹

A strategic observation regarding model performance indicates that combining high-level, interpretable features (lexicosyntactic metrics) with low-level, abstract features (character n-grams) consistently improves classifier performance.¹ This suggests that high diagnostic acuity in language-based screening requires integrating both human-meaningful indicators of linguistic decline

and latent stylistic patterns that are inaccessible to human linguistic analysis, justifying the sophisticated fusion approach. Furthermore, the role of the LLM in the multimodal pipeline is critical; it functions as an interpretability enforcer, translating complex, raw inputs (e.g., a video of a paper folding task) into clean, structured metrics (e.g., a score for *MotorCoordination*).¹ This mechanism is necessary to structure and simplify complex clinical inputs into forms compatible with downstream explainable models, such as decision trees, thereby successfully balancing computational power with the prerequisite for clinical utility.

Table I: Feature Sets and Clinical Interpretability

Feature Family	Source Data	Example Metrics	Primary Utility
Lexicosyntactic & N-grams	Spoken/Written Language (CTPC, ADBC)	Type-Token Ratio (TTR), Proposition Density, Character Bigrams	Capturing sub-clinical decline in lexical richness and stylistic complexity for high-performance classification.
Multimodal MMSE-Derived	MMSE Task Media (Video, Image, Audio)	MotorCoordination, ShapeCompleteness, Speech Fluency, ResponseLatency	Concise, numeric/categorical scores produced by an LLM extractor; designed for immediate clinician interpretability and

			high-resolution assessment.
EHR-Derived Symptom Timelines (CARE AD)	Longitudinal Clinical Notes (VHA)	Age-aware profiles: Cognitive Impairment, Functional Decline, Neuropsychiatric Symptoms	Input for multi-agent reasoning, detecting subtle, progressive patterns and providing long-term context.

IV. Advanced Model Architectures for Diagnosis and Risk Stratification

The integrated pipeline relies on three complementary modeling paradigms, each optimized for a specific role within the diagnostic journey.

4.1 Pillar 1: High-Performance Frontline Triage—Stacked Fusion Ensembles

The language-based classifiers employ stacked generalization, training a meta-classifier on out-of-fold predictions derived from a pool of level0 base learners, including Random Forest (RF), XGBoost (XGB), Linear Discriminant Analysis (LDA), Support Vector Classifier (SVC), Gaussian Naive Bayes (GNB), Logistic Regression (LR), and Multi-Layer Perceptron (MLP).¹ This methodology substantially reduces complexity by purifying the base set and consistently improves the outcome over simpler voting ensembles.¹

The stacked fusion approach achieved state-of-the-art results in language-based screening. Ensemble stacking and feature fusion yielded AUCs consistently **above 0.98** on language-only datasets.¹ Furthermore, the fusion of spoken and written data improved performance, with AUCs reported up to

~99.2% in linked experiments using purified base sets. Accuracy and F1 scores were often

reported in the mid-90s on held-out tests.¹ The robust, high performance, achieved using low-friction data such as short speech or writing samples, makes this method uniquely suitable for initial triage and frontline flagging.¹

4.2 Pillar 2: Transparent In-Clinic Assessment—Multimodal LLM Feature Extraction

The multimodal LLM pipeline processes media gathered during MMSE-like tasks (video, images, audio) to produce concise, structured, and interpretable metrics like *Legibility* and *ShapeCompleteness*.¹ This component is strategically paired with interpretability mechanisms. Decision trees are trained on these engineered features to derive compact rule lists that explicitly expose the prediction pathways.¹ To further support actionable clinical interpretation, counterfactual explanation modules are employed, producing the minimal feature changes required to flip a prediction.¹ This design ensures that the pipeline enhances screening sensitivity while simultaneously preserving transparency, delivering clinician-friendly justifications for predictions.¹

4.3 Pillar 3: Long-Horizon Risk Forecasting—The CARE AD Multi-Agent System

The CARE AD framework is a specialized LLM system designed for analyzing longitudinal EHRs. It employs specialized LLM agents simulating clinical roles (e.g., primary care, neurology, geriatrics, psychiatry) to evaluate symptom trajectories produced by the data-extraction agent.¹ A supervising AD specialist agent then synthesizes these domain-specific assessments to provide individual risk estimates across multiple horizons, ranging from one day to ten years prior to documented diagnosis.¹

CARE AD significantly outperformed both single-model baselines and conversational multi-agent baselines at equivalent LLM call budgets.¹ It demonstrated notably high accuracy at near-term horizons, achieving approximately

0.83 accuracy at -1 day prior to diagnosis. Although performance naturally declined at longer horizons, it showed meaningful discrimination even several years prior to diagnosis, reaching approximately **0.53 accuracy at -10 years**.¹

The observed challenge of predicting AD ten years in advance (0.53 accuracy) is an intrinsic difficulty rooted in sparse, longitudinal data. However, the multi-agent design strategically improves long-horizon sensitivity by combining complementary domain perspectives, thereby maximizing the informational value of subtle, sparse signals that a single, generalized model might overlook. This approach confirms that specialized, role-based reasoning is essential for successful long-range risk stratification.¹

Table II: Comparison of Integrated AD Detection Models

Model Pillar	Data Modality	Primary Output/Role	Key Performance Insight
Stacked Fusion Ensemble (Language)	Spoken/Written Text	Frontline Screening/Initial Triage	High diagnostic acuity: AUCs consistently above 0.98, substantial gain over single models.
Multimodal MMSE Extractor + DT	Video, Image, Audio (MMSE tasks)	Transparent Assessment/Confirmation	Enhanced screening sensitivity; delivers transparent, actionable rules and counterfactual explanations.
CARE AD Multi-Agent System	Longitudinal Clinical Notes (EHRs)	Long-Term Risk Monitoring/Forecasting	Effective long-horizon risk stratification: Accuracy decreases with time (0.83 acc. at -1 day vs. 0.53 at -10 years).

The strategic difference between these models is crucial: the highest diagnostic performance is achieved by the computational complexity of the Stacked Fusion ensemble (AUC up to 0.992).¹ Conversely, the transparent MMSE pipeline is designed for high sensitivity and interpretability. This configuration reinforces the principle of task-appropriate deployment:

the powerful, high-throughput model is utilized for high-sensitivity

flagging (triage), while the transparent model is reserved for *clinical validation* (assessment), ensuring the process adheres to both performance and regulatory requirements.

V. The Strategic Roadmap: Integrating Technologies into Clinical Workflow

The integrated pipeline is designed as a practical, multi-tier screening and risk-stratification strategy that harnesses the complementary strengths of the three modeling pillars to manage clinical resources efficiently and provide comprehensive evidence.¹

5.1 The Proposed Multi-Tier Strategy: Triage, Assessment, and Monitoring

1. **Tier 1: Frontline Screening:** This stage involves deploying the language-based stacked fusion models. Using low-friction data acquisition, such as short speech or brief writing samples, the models quickly flag individuals who require further assessment. The method's high AUC and accuracy make it ideal for resource-efficient initial triage.¹
2. **Tier 2: In-Clinic Multimodal Assessment:** For flagged cases, the workflow escalates to recording brief MMSE-like multimodal tasks. The multimodal LLM extractor is applied, followed by an interpretable decision tree. Clinicians are provided with readable rules and counterfactual explanations, delivering high-resolution screening evidence alongside essential transparency.¹
3. **Tier 3: Long-Term Risk Monitoring:** The CARE AD multi-agent system is tasked with analyzing longitudinal EHR narratives via its specialist agents. This deep analysis aims to detect subtle, progressive patterns that may presage a diagnosis years in advance, enabling earlier interventions and the scheduling of targeted follow-ups.¹

5.2 Synergistic Evidence: How Complementary Strengths Reduce Diagnostic Friction

The integration ensures that outputs from each stage provide complementary evidence necessary for a thorough clinical conclusion¹: the language models deliver sensitive early signals; the multimodal MMSE metrics provide high-resolution screening evidence with transparency; and CARE AD furnishes longitudinal perspective and crucial domain-aware context.¹

This tiered strategy is critical for mitigating physician fatigue and potential overload. By deploying the high-performance, low-cost Tier 1 (language models) as the initial filter, the system effectively ensures that the more complex and resource-intensive analyses (Tier 2 and Tier 3) are reserved only for a smaller cohort of high-risk patients. This strategy addresses the core inefficiencies noted in traditional diagnostics.¹ Furthermore, explainability is operationalized across the workflow: Tier 1 offers diagnostic

confidence (AUC > 0.98), Tier 2 provides *justification* (rules/counterfactuals), and Tier 3 offers *context* (a longitudinal timeline). This cascade of explanatory functions is necessary for any clinically viable AI system, providing clinicians with tailored explanations relevant to their role at each phase of the diagnostic journey.

VI. Critical Limitations and Future Research Imperatives

Despite the significant advancements, several critical challenges must be addressed before widespread clinical adoption, specifically concerning generalizability, interpretability, and ethical governance.

6.1 The Generalizability Challenge: Scaling from Lab to Clinic

The current limitations in model robustness across heterogeneous clinical environments constitute the primary bottleneck for scaling LLM-driven AD tools.¹

- **Diagnostic Equity and Dialectal Variation:** Standard Automatic Speech Recognition (ASR) systems struggle with dialectal variations, resulting in inaccurate transcripts and, consequently, biased diagnostic results in speech-based screening.¹ The necessity of extensive fine-tuning for each new dialect creates a major scalability bottleneck, hindering deployment across diverse global populations.¹ Addressing ASR issues is not merely a technical failure; it is a clinical and regulatory vulnerability that violates the

principle of diagnostic equity.

- **EHR Transferability:** Models developed for granular phenotyping are highly sensitive to variations in physician writing styles, the use of different clinical abbreviations, and local EHR documentation practices.¹ Data preprocessing is often highly tailored for specific EHR providers, severely limiting cross-system generalizability and requiring retraining and calibration for successful transfer.¹

To overcome these issues, future research must prioritize developing models that are inherently robust against variations in dialect, background noise, and heterogeneous clinical documentation styles across different healthcare systems, moving beyond models tailored to single institutions or small, controlled datasets.¹

6.2 The Interpretability Mandate: Addressing Visual Grounding and Agent Auditability

While the integrated pipeline prioritizes interpretability in the MMSE component, crucial gaps remain in complex, large-scale multimodal systems.¹

- **Lack of Integrated Visual Grounding:** A critical limitation for generative multimodal reports, particularly those analyzing complex data like 3D MRI scans, is the inability to provide visual explanations alongside the generated textual findings.¹ Clinical gold standards require reports to include annotations or highlights on the images that directly justify the textual conclusions. The absence of this integrated visual grounding reduces clinical utility and physician trust, ultimately forcing manual verification of the AI's claims and undermining the goal of reducing physician workload.¹ While some researchers use "Cognition-Guided Prompting" to ground textual reports in quantitative evidence, this fails when the evidence is fundamentally visual. The true challenge is integrating these disparate forms of evidence—quantitative probability, textual narrative, and visual annotation—into a single, cohesive, trust-building output.
- **Formalizing Agent Reasoning:** For multi-agent systems like CARE AD, accountability and trustworthiness depend on transparent operations. Future work must formalize **audit trails for agent reasoning** and counterfactual outputs to provide clear pathways for verification and clinical review.¹

A key future imperative is mandating that models include annotated imaging evidence in the LLM output to address the lack of integrated visual grounding and thereby significantly increase physician trust and verification.¹

6.3 Data Integrity and Ethical Governance

The application of LLMs in medicine introduces complex ethical and data governance risks that must be managed.

- **Risks of Synthetic Data:** Researchers are increasingly utilizing LLM-generated synthetic training data to overcome the scarcity of large, labeled multimodal datasets.¹ However, relying on these synthetic corpora introduces fundamental risks, including the possibility of subtle biases or the omission of rare but critical clinical variations, potentially leading to model "hallucinations" or reports lacking necessary clinical precision.¹
- **Data Privacy and Security:** The utilization of EHR data and patient media necessitates strict governance and de-identification procedures. Multi-agent LLMs must be deployed under secure, auditable conditions to maintain data privacy.¹

Given the growing reliance on LLM-generated data, new and rigorous methods for **independent validation** are required to ensure the clinical fidelity of these synthetic corpora and prevent the introduction of systemic biases into diagnostic and reporting pipelines.¹ Furthermore, prospective clinical trials are necessary before these methods can alter established screening guidelines; the current methods are designed to augment, not replace, clinician judgment.¹

VII. Conclusion and Recommendations for Clinical Implementation

The synthesis of stacked fusion language classifiers, the CARE AD multi-agent EHR reasoning framework, and the multimodal MMSE feature extraction pipeline demonstrates a profound shift toward early, objective, and explainable AD detection. These specialized systems showcase excellence across diagnostic acuity (AUCs > 0.98), data synthesis capabilities, and the provision of phenotypic depth.¹

The most critical requirement for realizing the full potential of these advancements is the transition from specialized, siloed tools to a unified, clinically regulated platform. The final step in this strategic roadmap must be the merging of all specialized tools—ranging from the high-performing speech classifier to the transparent multimodal assessment module and the longitudinal risk monitor—into a single, cohesive diagnostic platform. This integration is essential to effectively reduce physician workload, manage computational complexity, and accelerate the rate of early, accurate AD diagnosis.¹

Future work must prioritize the following research imperatives:

1. **Clinical Validation:** Conducting real-world prospective clinical trials to evaluate deployment in existing clinical workflows.¹
2. **Generalizability and Equity:** Expanding multilingual capabilities and developing models robust against heterogeneous documentation styles and dialectal variations.¹
3. **Auditability and Trust:** Formalizing audit trails for agent reasoning and counterfactual outputs to ensure accountability and establish the necessary levels of trust for clinician adoption.¹
4. **Data Integrity:** Establishing rigorous, independent validation frameworks for LLM-generated synthetic data to maintain clinical fidelity.¹