

1 **Agentic AI for automated hypothesis testing in Alzheimer's**  
2 **disease and related dementias**

3

4 Ketan S. Saichandran<sup>a,b</sup>, Karim Elzokm<sup>a,c</sup>, Osman B. Guney<sup>a,c</sup>, Vijaya B. Kolachalama<sup>a,b,d,†</sup>

5

6 <sup>a</sup>Department of Medicine, Boston University Chobanian & Avedisian School of Medicine,  
7 Boston, MA, USA – 02118

8 <sup>b</sup>Department of Computer Science, Boston University, Boston, MA, USA – 02215

9 <sup>c</sup>Department of Electrical & Computer Engineering, Boston University, Boston, MA, USA –  
10 02215

11 <sup>d</sup>Faculty of Computing & Data Sciences, Boston University, Boston, MA, USA– 02215

12

13 †Corresponding author

14 Email: [vkola@bu.edu](mailto:vkola@bu.edu)

15 Phone: (+1)617-358-7253

16 Address: 72 E. Concord St, Evans 636, Boston, MA – 02118

17 ORCID: <https://orcid.org/0000-0002-5312-8644>

18

19 **Abstract**

20 **Background:** Alzheimer’s disease and related dementias (ADRD) are multifactorial disorders  
21 driven by genetic, vascular, inflammatory, lifestyle, and biological factors interacting over  
22 decades, requiring population-level analyses. However, utilizing multimodal datasets is labor-  
23 intensive, impeding discovery and equity. Multi-agentic AI systems promise workflow  
24 automation but lack ADRD-specific adaptations.

25 **Methods:** We created a modular AI system with agents for data preparation, planning, analysis,  
26 critique, and summarization. Using data from the National Alzheimer’s Coordinating Center  
27 (NACC) (n=48,876), we tested 100 literature-derived ADRD hypotheses, subsampled cohorts  
28 (100–10,000 participants), and 100 non-ADRD controls via iterative loops.

29 **Results:** The system verified  $42.0 \pm 2.8$  hypotheses by rejecting the null, found no evidence for  
30  $16.6 \pm 1.1$ , and deemed  $41.4 \pm 2.3$  not testable, with consistency in negated versions. Verification  
31 rates rose from 10.4 (n=100) to 35.4 (n=10,000); controls showed  $85.2 \pm 0.6$  non-testable.

32 **Conclusions:** Our system illustrates a powerful paradigm in which AI can facilitate automated  
33 data analysis, expediting advancements in population-level dementia investigations.

34

## 35 1 | BACKGROUND

36 Alzheimer’s disease (AD) and related dementias (ADRD) arise from multifactorial causes that  
37 include genetic susceptibility, vascular and inflammatory processes, lifestyle influences, and  
38 other biological pathways interacting over years.<sup>1</sup> Understanding how these factors intersect to  
39 shape disease onset and progression demands population-level analyses. Typically, these  
40 analyses are performed on datasets containing tens to thousands of variables ranging from  
41 demographics, comorbidities, and neuropsychological assessments to neuroimaging, fluid  
42 biomarkers, and genetics, reflecting the layered nature of ADRD pathophysiology. A wealth of  
43 population-based cohorts such as the Alzheimer’s Disease Neuroimaging Initiative,<sup>2</sup> and the  
44 National Alzheimer’s Coordinating Center (NACC),<sup>3,4</sup> among others, offer unparalleled access  
45 to rich, multimodal data resources. These datasets enable exploration of complex associations  
46 such as amyloid-vascular interactions,<sup>5,6</sup> or ethnicity-specific trajectories,<sup>7-9</sup> but leveraging them  
47 remains laborious and fragmented.<sup>10</sup> Researchers often spend significant time to identify  
48 variables, align data dictionaries, harmonize features and manage missingness, even for  
49 seemingly simple hypotheses (e.g., whether APOE4 modulates vascular effects on hippocampal  
50 atrophy).<sup>11</sup> This friction slows discovery, increases analytic errors, and widens disparities  
51 between well-resourced and smaller research teams.

52

53 Recent years have witnessed efforts to standardize multimodal ADRD research through  
54 coordinated data infrastructures, analytic tools and AI-driven multimodal analysis.<sup>12-14</sup> Platforms  
55 such as the National Institute on Aging Genetics of Alzheimer’s Disease Data Storage Site Data  
56 Sharing Service,<sup>15,16</sup> and the AD Workbench,<sup>17</sup> have improved accessibility and interoperability  
57 of large-scale cohorts. Technical standards such as the Brain Imaging Data Structure,<sup>18</sup> and  
58 Findable, Accessible, Interoperable, Reusable (FAIR) data principles,<sup>19</sup> have enabled more  
59 consistent data exchange across repositories. Machine learning frameworks like MONAI,<sup>20</sup> have  
60 been developed to analyze neuroimaging data, while federated learning and harmonization  
61 methods have addressed privacy-preserving analysis across sites.<sup>21</sup> Despite these advances, most  
62 systems remain specialized for discrete tasks, such as image harmonization, and do not integrate  
63 the full scientific workflow from hypothesis generation to testing. In parallel, recent  
64 developments in multi-agentic AI systems have begun transforming how scientific discovery can  
65 be automated. These systems, which include specialized agents each responsible for tasks like  
66 hypothesis formulation, experimental design, or critique, can collaborate autonomously to  
67 conduct literature reviews, run simulations, and refine hypotheses based on results.<sup>22-24</sup> They  
68 showcase the feasibility of embedding scientific reasoning within collaborative AI architectures,  
69 reducing manual burden while maintaining interpretability and rigor. Yet, most of them operate  
70 in general-purpose or simulated environments, lacking adaptation to the specific data standards,  
71 ontologies, and multimodal complexities of domains such as ADRD. What remains missing is a  
72 unifying, domain-aware and modular system that can autonomously handle data, translate  
73 unstructured research queries into analytic hypotheses, and test them rigorously. Such a system

74 would bring conceptual automation to discovery itself, bridging data analysis and hypothesis  
75 testing in an iterative, reproducible, and extensible manner.

76

77 Here, we developed a multi-agentic AI system for accelerated hypothesis testing in ADRD  
78 research. This system combines modular reasoning with robust data processing, allowing each  
79 component to be replaced or improved without disrupting the overall workflow. The system  
80 comprises specialized agents, which collaboratively process spreadsheets representing cohort-  
81 level data and their corresponding data dictionaries. The workflow operates through two nested  
82 loops: an outer loop where the data preparation and planning agents filter relevant variables and  
83 outline analysis plans, and an inner loop where the scientist and critic agents execute statistical  
84 analyses and perform bias or reproducibility checks. User inputs, ranging from natural language  
85 hypotheses to folders containing raw data, are converted into structured analytic pipelines. The  
86 critic agent determines outcomes, followed by generation of plots and summaries. We  
87 demonstrated the system's robustness across multiple scenarios: automated testing and  
88 replication of literature-derived ADRD hypotheses on the NACC cohort, verifying known  
89 associations and identifying non-testable claims; stable performance across dataset sizes; and  
90 high specificity in a negative control using randomly selected non-ADRD hypotheses, where  
91 majority were correctly classified as not testable on NACC data. As proof of concept, our system  
92 points toward a new paradigm, where AI systems can automatically analyze data and accelerate  
93 progress in population-level dementia research.

94

## 95 2 | METHODS

### 96 2.1 | Study population

97 The study utilized data from the National Alzheimer’s Coordinating Center (NACC), which  
98 aggregates information from multiple cohorts and is based on the Uniform Data Set (UDS) 3.0  
99 dictionary.<sup>25</sup> A total of 48,876 participants were included, with eligibility based on a diagnosis of  
100 normal cognition (NC), mild cognitive impairment (MCI), or dementia. We did not exclude  
101 cases based on the absence of features or diagnostic labels. For individuals from the NACC  
102 cohort with multiple clinical visits, we prioritized visits at which the participant received a  
103 dementia diagnosis; among these, we selected the visit with the most available data features,  
104 prioritizing neuroimaging information, and chose the most recent visit in cases of ties. This  
105 approach maximized the sample sizes of dementia cases, ensured inclusion of each individual’s  
106 latest record, and optimized the utilization of available data. Various features (n=750) were  
107 included, encompassing subject demographics, medical history, laboratory results, medications,  
108 neuropsychological tests, and functional assessments. All participants or their designated  
109 informants provided written informed consent.

110

### 111 2.2 | Hypotheses curation

112 To curate the literature-derived hypotheses for evaluating agentic AI’s performance, we first  
113 compiled a bag of keywords encompassing key ADRD domains, such as “*Alzheimer’s disease*  
114 *risk factors*,” “*dementia biomarkers*,” “*APOE4 interactions*,” “*vascular contributions to*  
115 *neurodegeneration*,” “*ethnicity-specific AD trajectories*,” and related terms drawn from  
116 established literature in the field. Using a large language model (ChatGPT, based on GPT-4), we  
117 prompted it to generate hundreds of diverse hypotheses inspired by these keywords, along with  
118 suggested peer-reviewed references supporting each one. We then manually reviewed each  
119 generated hypothesis for scientific plausibility and selected 100 hypotheses relevance to  
120 population-level analyses, and alignment with variables available in the NACC cohort; for each,  
121 we accessed the original publications (primarily from 2000 to 2024) to verify accuracy, extract  
122 supporting excerpts or quotes, and resolve any discrepancies or inaccuracies through team  
123 consensus. Hypotheses were refined for clarity, feasibility within our system, and diversity  
124 across pathophysiological pathways, resulting in a final set that represented a broad spectrum of  
125 testable ADRD research questions while excluding those requiring unavailable data modalities.  
126 The list of hypotheses along with corresponding references are provided in **Supplementary**  
127 **Table 1**.

128

### 129 2.3 | Multi agentic system

130 Our agentic system employs four specialized large language model (LLM) agents, data  
131 preparation agent, planning agent, scientist agent, and critic agent, which collaborate to conduct  
132 comprehensive analyses. The pipeline begins with the data preparation agent, where relevant  
133 variables are automatically selected based on hypothesis relevance. Following data preparation,  
134 the planning agent designs the analytical strategy and recommends appropriate statistical tests  
135 and methodologies. Subsequently, an inner loop cycle is initiated between the scientist and critic  
136 agents. The scientist agent, following the plan, writes and executes code to conduct required  
137 analyses. The critic agent evaluates the results, either approving or rejecting the analysis, and  
138 provides feedback to the scientist agent. This inner loop continues until either the hypothesis is  
139 approved (i.e., null rejected), rejected (i.e., null not rejected), or determined to be not testable. If  
140 not approved, the system may return to data preprocessing with alternative variable selections in  
141 an outer loop iteration, allowing for refinement and exploration of different analytical  
142 approaches. The system overview is presented in **Figure 1**. Detailed descriptions of each agent  
143 are provided below and prompts used on the agentic system are provided in **Supplementary**  
144 **Table 2**.

145 **Data preparation agent.** The data preparation agent automatically detects tabular data files  
146 (CSV format) and their corresponding data dictionaries (in JavaScript Object Notation (JSON)  
147 format) within the workspace folder. It loads both into a centralized cache that stores the dataset  
148 along with rich metadata, including variable names, descriptions, inferred data types, and pre-  
149 computed statistical summaries. For each user-provided hypothesis, the agent uses the hypothesis  
150 text together with the full variable dictionary to automatically select the most relevant features.  
151 Selected variables are then systematically preprocessed: cohort-specific missing-value codes are  
152 converted to NaN, categorical variables are one-hot encoded, and continuous/integer variables  
153 are standardized as needed. The cleaned, processed dataset is saved to the working directory, and  
154 both its file path and updated variable descriptions are injected into the conversation context for  
155 downstream agents.

156 **Planning agent.** The planning agent designs multi-step research strategies and analytical  
157 systems, providing methodological recommendations and validation before analysis begins.  
158 After adding the generated plan to the conversation context, the system enters an iterative loop  
159 between the scientist agent and the critic agent, which we term as the “inner loop”.

160 **Scientist agent.** The scientist agent performs comprehensive statistical analysis, including  
161 exploratory data analysis (EDA), hypothesis testing, regression modeling (linear, logistic,  
162 survival), machine learning approaches, mediation and moderation analysis, and subgroup  
163 analyses. Before each analysis, the system automatically provides the scientist agent with  
164 detailed variable descriptions, data types, and statistical summaries, ensuring informed analytical  
165 decisions. The agent then generates Python code for statistical computations and executes  
166 analyses with error handling and missing data management. The code is run in a packaged and  
167 secure environment to ensure safety and prevent unintended system modifications.

168

169 **Critic agent.** The critic agent reviews and validates scientific approaches, providing  
170 methodological feedback, and determining hypothesis verification status. The agent evaluates  
171 code correctness, statistical validity, result interpretation, and hypothesis alignment, ensuring  
172 rigorous scientific standards are maintained. The critic is guided by a comprehensive prompt that  
173 instructs it to examine multiple aspects of the analysis: reasonable statistical analysis and  
174 adequate sample sizes, sound experimental design that properly tests the hypothesis, appropriate  
175 data preprocessing and feature engineering, valid validation methodology, meaningful metrics  
176 that address the hypothesis, clear interpretation of results and acknowledgment of limitations,  
177 and functional, error-free, and well-documented code that produces valid results. The prompt  
178 includes mandatory code checks requiring that code runs without errors, all imports and  
179 dependencies are properly handled, data loading and preprocessing are correct, statistical  
180 calculations are accurate, and output matches the described analysis. Additionally, the prompt  
181 enforces critical statistical terminology verification, ensuring that results are interpreted using  
182 proper language (e.g., “reject the null hypothesis” rather than “accept the null hypothesis”) and  
183 that p-values, effect sizes, and confidence intervals are properly reported. The prompt also  
184 includes a hypothesis alignment check, requiring that results not only show statistical  
185 significance but also that the direction of effect matches the stated hypothesis. The critic must  
186 output its verdict using XML (a markup language for storing and transporting data) tags:  
187 `<verdict>NULL REJECTED</verdict>`, `<verdict>NULL NOT REJECTED</verdict>`, or  
188 `<verdict>NOT TESTABLE</verdict>`, with the latter reserved for cases where dataset  
189 limitations discovered during analysis prevent meaningful hypothesis testing. If the critic agent is  
190 satisfied with the statistical results, it declares the hypothesis as verified. However, if it is not  
191 satisfied, the loop continues for a fixed number of inner loop iterations. After these inner  
192 iterations, the system initiates a new outer iteration, which resets the conversation context but  
193 incorporates summaries from previous attempts. The planning agent is then prompted to generate  
194 a new plan, enabling the system to explore alternative variables, statistical methods, or analytical  
195 systems. We note that even during the inner loop iterations, the system summarizes the  
196 conversation if it grows beyond a certain number of messages, to ensure the context window of  
197 the LLM is not exceeded.

198 The critic agent also includes a specialized debugging mode that activates automatically when  
199 the scientist agent’s code fails to execute. When code execution errors occur, the critic switches  
200 to a debugging query mode that focuses exclusively on identifying and resolving code issues  
201 rather than evaluating the hypothesis. In this mode, the critic analyzes execution failures,  
202 identifies likely causes (such as syntax errors, import issues, data access problems, statistical  
203 analysis errors, or data preprocessing issues), and provides specific, actionable feedback to guide  
204 the scientist toward a working solution. The debugging query emphasizes code robustness, error  
205 handling, and best practices, with particular attention to statistical analysis errors including array  
206 shape compatibility, missing value handling, sample size requirements, and data type

207 conversions. Importantly, when in debugging mode, the critic does not provide a hypothesis  
208 verdict (NULL REJECTED=NULL NOT REJECTED/NOT TESTABLE) since the code has not  
209 successfully executed, allowing the conversation to continue after providing debugging guidance  
210 so the scientist can fix the code and reattempt the analysis.

211

## 212 **2.4 | Technical implementation and reproducibility**

213 The workflow is orchestrated through a conversation management system that maintains the  
214 complete conversation history between agents, tracks turn counts and implements context  
215 window management. The conversation manager stores all messages in a structured format,  
216 including user inputs, agent responses, and system messages containing variable descriptions and  
217 dataset metadata. To prevent exceeding the LLM context window during extended conversations,  
218 the system implements automatic conversation truncation that preserves critical context: previous  
219 iteration summaries (for outer loop iterations), the initial hypothesis and preprocessing messages,  
220 variable description messages, and the most recent conversation exchanges. This selective  
221 preservation ensures that essential context is maintained while removing intermediate messages  
222 that may be redundant.

223 Dataset information is systematically provided to agents at multiple stages of the workflow to  
224 ensure continuous access to relevant context. During the preprocessing phase, the data  
225 preparation agent injects a comprehensive system message containing all available variables with  
226 their descriptions, data types, and dataset statistics (total rows, columns, missing values, numeric  
227 and categorical column counts). Additionally, before each scientist agent round within the inner  
228 loop, variable descriptions are automatically re-injected into the conversation context, providing  
229 fresh access to variable names, descriptions, and metadata. This repeated injection ensures that  
230 agents maintain awareness of available variables and their characteristics throughout extended  
231 conversations, even after conversation truncation. The truncation mechanism specifically  
232 preserves these variable description messages, recognizing their critical importance for informed  
233 analytical decision-making.

234 For outer loop iterations, conversation summaries are generated using the LLM to condense  
235 previous attempts into concise context. These summaries capture the hypothesis being tested,  
236 approaches attempted, variables used, results found, and reasons for non-approval, enabling  
237 subsequent iterations to avoid repeating failed approaches and explore alternative methodologies.  
238 The summary is prepended to the conversation context at the start of each new outer loop  
239 iteration, allowing the planning agent to adapt its strategy based on previous attempts.

240 The system includes a code executor with configurable timeout mechanisms (default: 30  
241 seconds), persistent global namespaces for variable continuity across code blocks, and  
242 comprehensive error handling. All agents use the same LLM instance loaded in memory, each

243 with its own specialized prompt. We primarily used the Qwen3-4B-Instruct model  
244 (Qwen/Qwen3-4B-Instruct-2507) as the client. The system tracks code execution outcomes and  
245 automatically switches the critic agent to debugging mode when code execution fails, allowing  
246 for targeted error resolution before hypothesis evaluation.

247 The system supports configuration management through YAML (a human-readable data  
248 serialization language) files with command-line overrides, enabling reproducible experiments  
249 with documented parameter settings. All runs include configuration logging, random seed  
250 management, and structured logging of experiment parameters for reproducibility. Configuration  
251 details including model name, seed, hypothesis mode, and other experiment parameters are  
252 automatically logged to JSON files for each run. All agent interactions are also logged to  
253 structured JSON files that capture the complete workflow: hypothesis details, preprocessing  
254 steps, conversation turns with agent responses, code execution results, statistical outputs, and  
255 final verdicts. This comprehensive logging enables full reproducibility and post-hoc analysis of  
256 the agentic decision-making process.

257

## 258 **2.5 | Data availability**

259 Data from the NACC cohort can be requested and downloaded at <https://naccdata.org>.

260

### 261 3 | RESULTS

262 We rigorously tested the multi-agentic AI system on the NACC cohort via an exemplar  
263 hypothesis evaluation, batch analyses of 100 curated ADRD hypotheses, scalability trials on  
264 subsampled datasets, and negative controls with non-ADRD queries. These assessments reveal  
265 the system’s proficiency in automated, verifiable hypothesis testing, with outcomes emphasizing  
266 its precision, adaptability, and domain-specific utility. Detailed results from each component are  
267 presented sequentially below.

268

269 In an exemplar application, the multi-agentic AI system processed the hypothesis “*Alcohol*  
270 *consumption is associated with increased dementia risk*” using the NACC cohort (**Figure 2**). The  
271 Data Preparation Agent selected 19 relevant variables, including alcohol measures (ALCFREQ,  
272 ALCOCCAS, ALCOHOL), demographics (NACCAGE, SEX, RACE, EDUC), dementia  
273 outcomes (NACCIDEM, NACCALZD, etc.), and comorbidities (DEP, HYPERT, STROKE); it  
274 replaced unknown values with NaN and expanded categorical values into a dataset with binary  
275 indicators for alcohol frequency. The Planning Agent specified chi-square tests with Cramer’s V  
276 to examine associations between alcohol frequency categories (less than monthly, monthly,  
277 weekly, few times per week, daily/almost daily) and binary dementia progression. Executing via  
278 listwise deletion, the Scientist Agent detected significant associations across all categories ( $p <$   
279  $0.0001$ ), with small effect sizes (Cramer’s V 0.040-0.049, highest at 0.049 for “few times per  
280 week”), implying a dose-response pattern. The Critic Agent affirmed methodological soundness,  
281 noting limitations like cross-sectional causality and recall bias, leading to null hypothesis  
282 rejection in a single iteration. This finding aligns with systematic reviews and cohort studies  
283 indicating that alcohol intake,<sup>26-29</sup> even at moderate levels, elevates dementia risk without  
284 protective effects. Furthermore, analyses using NACC and similar data support heightened  
285 dementia risk from alcohol use disorder, particularly in subgroups like men with  
286 neuropsychiatric conditions.

287

288 We assessed the agentic system’s efficacy by applying it to 100 hypotheses curated from  
289 established ADRD literature (**Supplementary Table 1**), encompassing associations such as  
290 genetic risk factors (e.g., APOE  $\epsilon 4$ ), lifestyle influences (e.g., physical activity), and biomarker  
291 interactions (e.g., amyloid and vascular pathology), using the NACC cohort ( $n=48,876$   
292 participants). As depicted in **Figure 3**, across replicated runs to account for variability in agentic  
293 processing, the system rejected the null hypothesis, indicating statistically significant verification  
294 of the proposed association, for a mean of  $42.0 \pm 2.8$  hypotheses, failed to reject the null (denoting  
295 absence of significant evidence) for  $16.6 \pm 1.1$  hypotheses, and classified  $41.4 \pm 2.3$  as not testable  
296 owing to missing variables, incompatible data formats, or insufficient cohort coverage. This  
297 outcome highlights the system’s capability to autonomously replicate a portion of known

298 literature findings while prudently flagging unfeasible queries, thereby enhancing reliability in  
299 population-level analyses.

300

301 To evaluate the scalability and sensitivity of the multi-agentic AI system to cohort size, we  
302 subsampled the NACC dataset (original  $n=48,876$ ) to create reduced cohorts of 100, 1,000 and  
303 10,000 participants while preserving demographic and clinical distributions, then applied the  
304 system to the same 100 literature-derived ADRD hypotheses as in prior analyses. As shown in  
305 **Figure 4**, performance exhibited a clear trend toward enhanced hypothesis verification with  
306 increasing dataset size: at  $n=100$ , the system rejected the null hypothesis for only 10.4  
307 hypotheses (indicating limited detection of significant associations due to low power), failed to  
308 reject it for 31.6 (suggesting inconclusive evidence), and deemed 58.0 not testable owing to  
309 sparse or missing data in small samples. Scaling to  $n=1,000$  improved resolution, with null  
310 rejections rising to 23.4, non-rejections at 27.4, and non-testable dropping to 49.2, reflecting  
311 better statistical power for identifying associations. At  $n=10,000$ , approaching the full cohort, the  
312 system achieved the highest verification rate, rejecting the null for 35.4 hypotheses, not rejecting  
313 for 21.6, and classifying 43.0 as not testable, demonstrating robustness in larger datasets where  
314 variable coverage and sample diversity enable more comprehensive analyses. Corresponding  
315 statistics between different groups are presented in **Supplementary Table 3**. These results  
316 underscore the system's dependence on adequate data volume for accurate, reproducible  
317 hypothesis testing in ADRD research, with diminishing non-testability and stabilized outcomes  
318 as cohort size grows.

319

320 To assess the robustness of our agentic system in a control condition unrelated to ADRD, we  
321 evaluated its performance on 100 randomly generated hypotheses derived from general medical  
322 literature (**Supplementary Table 4**). These hypotheses spanned diverse topics such as  
323 cardiovascular disease, infectious diseases, endocrinology, and oncology, serving as a  
324 benchmark to test the system's ability to handle queries outside its primary domain while  
325 providing a form of negative control validation. As illustrated in **Figure 5**, our system analyzed  
326 these hypotheses using NACC data. Across multiple runs, the system rejected the null hypothesis  
327 (indicating verification of the stated association) for an average of  $7.6\pm 0.8$  hypotheses, failed to  
328 reject the null (indicating no significant evidence for the association) for  $7.2\pm 1.3$  hypotheses, and  
329 deemed  $85.2\pm 0.6$  hypotheses not testable due to insufficient or incompatible data in the NACC  
330 dataset. This high rate of non-testability validates the system's conservative approach when  
331 applied to extraneous topics, highlighting its specificity for ADRD-related inquiries, minimizing  
332 spurious conclusions in unrelated contexts, and reducing the risk of false positives in out-of-  
333 domain scenarios. Although the hypotheses were randomly generated, few of them were still  
334 testable using the variables available in the NACC dataset. For example, the hypothesis  
335 "*Cholesterol-lowering diets reduce cardiovascular disease risk*" could be evaluated because the

336 dataset included indicators of hypercholesterolemia, cardiac events (such as heart attack, angina,  
337 or bypass surgery), and medication use relevant to cholesterol and blood-pressure management.  
338 The presence of these well-defined clinical and behavioral variables enabled the construction of  
339 measurable exposure and outcome categories, making it possible to statistically assess at least  
340 some of the randomly produced hypotheses. Similarly, the hypothesis “*Obesity increases the risk*  
341 *of obstructive sleep apnea*” was also testable, as the dataset contained BMI category variables  
342 and a clearly defined indicator of diagnosed sleep apnea. These variables allowed for the creation  
343 of categorical exposure and outcome measures, making it feasible to examine the association  
344 statistically despite the hypothesis being randomly generated.

345

## 346 4 | DISCUSSION

347 In this study, we developed and rigorously evaluated a multi-agentic AI system designed to  
348 automate and accelerate hypothesis testing in ADRD research. Using the NACC cohort  
349 (n=48,876 participants with rich demographic, clinical, neuropsychological, and biomarker data),  
350 we systematically tested 100 literature-derived ADRD hypotheses, their negated (opposite)  
351 counterparts, and critically, 100 randomly selected non-ADRD medical hypotheses as a negative  
352 control. The system's specialized agents collaboratively handled the entire scientific workflow,  
353 i.e., from natural-language hypothesis interpretation, ontology-aware variable mapping, and  
354 missing-data handling, through statistical modeling, iterative critique, and reproducible result  
355 generation. Our results demonstrate high fidelity in replicating established ADRD associations  
356 while correctly identifying non-testable claims, maintaining logical consistency when evaluating  
357 negated hypotheses, showing robustness across dataset sizes from 100 to 10,000 participants, and  
358 exhibiting strong specificity by deeming most out-of-domain (non-ADRD) queries untestable.

359  
360 Our multi-agentic system provides a structured approach to ADRD research by automating  
361 routine steps like data harmonization and basic testing, allowing researchers to focus more on  
362 interpretation and follow-up. It could help connect teams with varying resources, encouraging  
363 broader participation in examining disease pathways. For instance, midlife hypertension's  
364 adverse impact on cognitive decline or stroke history's link to dementia. By integrating with  
365 emerging platforms such as the AD Workbench along with Findable, Accessible, Interoperable,  
366 and Reusable (FAIR) principles, it might facilitate blending datasets for cross-study comparisons,  
367 such as metabolic syndrome's role in cognitive decline or sleep disorders' ties to amyloid  
368 buildup, fostering collaborative efforts to address dementia's multifactorial aspects. Over time,  
369 this could contribute to practical applications in biomarker discovery or treatment evaluation,  
370 making analyses more inclusive across settings.

371  
372 Our study has a few limitations. The hypotheses were selected from existing literature, so  
373 outcomes may align closely with those sources rather than highlight novel patterns, and the  
374 NACC dataset's emphasis could underrepresent certain groups, such as specific ethnicities or  
375 regions. Dependence on large language models (e.g., Qwen3-4B-Instruct-2507) might carry over  
376 training biases, influencing interpretations in some areas. The iterative loops addressed many  
377 scenarios but could overlook nuances in complex cases, such as mixed dementia etiologies or  
378 visuospatial issues in Lewy body dementia. Broader challenges include data diversity, where  
379 training biases might constrain applicability, and ethical considerations around AI in health  
380 research call for ongoing human oversight and real-world validation.

381

382 In conclusion, these findings establish the agentic system as a reliable, transparent, and scalable  
383 tool that substantially reduces manual effort and analytic variability in population-level dementia  
384 research. This system also provides a basis for integrating AI into ADRD studies, connecting  
385 datasets with systematic explorations. Future developments could incorporate federated learning  
386 for multi-site privacy or adaptable agents for new biomarkers from wearables. Through such  
387 steps, these tools might support progress in understanding dementia's diverse components, from  
388 genetic and lifestyle elements to therapeutic responses, aiding shared initiatives for improved  
389 outcomes.

390

## 391 REFERENCES

- 392 1. 2025 Alzheimer's disease facts and figures. *Alzheimer's & Dementia*. 2025;21(4). doi:  
393 10.1002/alz.70235.
- 394 2. Petersen RC, Aisen PS, Beckett LA, Donohue MC, Gamst AC, Harvey DJ, Jack CR, Jr.,  
395 Jagust WJ, Shaw LM, Toga AW, Trojanowski JQ, Weiner MW. Alzheimer's Disease  
396 Neuroimaging Initiative (ADNI): clinical characterization. *Neurology*. 2010;74(3):201-9. Epub  
397 20091230. doi: 10.1212/WNL.0b013e3181cb3e25. PubMed PMID: 20042704; PMCID:  
398 PMC2809036.
- 399 3. Beekly DL, Ramos EM, Lee WW, Deitrich WD, Jacka ME, Wu J, Hubbard JL, Koepsell  
400 TD, Morris JC, Kukull WA, Centers NIAAsD. The National Alzheimer's Coordinating Center  
401 (NACC) database: the Uniform Data Set. *Alzheimer Dis Assoc Disord*. 2007;21(3):249-58. doi:  
402 10.1097/WAD.0b013e318142774e. PubMed PMID: 17804958.
- 403 4. Beekly DL, Ramos EM, van Belle G, Deitrich W, Clark AD, Jacka ME, Kukull WA,  
404 Centers NI-AsD. The National Alzheimer's Coordinating Center (NACC) Database: an  
405 Alzheimer disease database. *Alzheimer Dis Assoc Disord*. 2004;18(4):270-7. PubMed PMID:  
406 15592144.
- 407 5. Biesbroek JM, Coenen M, DeCarli C, Fletcher EM, Maillard PM, Barkhof F, Barnes J,  
408 Benke T, Chen CPLH, Dal□Bianco P, Dewenter A, Duering M, Enzinger C, Ewers M, Exalto  
409 LG, Franzmeier N, Hilal S, Hofer E, Koek HL, Maier AB, McCreary CR, Papma JM, Paterson  
410 RW, Pijnenburg YAL, Rubinski A, Schmidt R, Schott JM, Slattery CF, Smith EE, Sudre CH,  
411 Steketee RME, Teunissen CE, van den Berg E, van der Flier WM, Venketasubramanian N,  
412 Venkatraghavan V, Vernooij MW, Wolters FJ, Xin X, Kuijf HJ, Biessels GJ. Amyloid pathology  
413 and vascular risk are associated with distinct patterns of cerebral white matter hyperintensities: A  
414 multicenter study in 3132 memory clinic patients. *Alzheimer's & Dementia*. 2024;20(4):2980-9.  
415 doi: 10.1002/alz.13765.
- 416 6. Chaudhuri S, Dempsey DA, Huang YN, Park T, Cao S, Chumin EJ, Craft H, Crane PK,  
417 Mukherjee S, Choi SE, Scollard P, Lee M, Nakano C, Mez J, Trittschuh EH, Klinedinst BS,  
418 Hohman TJ, Lee JY, Kang KM, Sohn CH, Kim YK, Yi D, Byun MS, Risacher SL, Nho K,  
419 Saykin AJ, Lee DY. Association of amyloid and cardiovascular risk with cognition: Findings  
420 from KBASE. *Alzheimer's & Dementia*. 2024;20(12):8527-40. doi: 10.1002/alz.14290.
- 421 7. Babulal GM, Quiroz YT, Albeni BC, Arenaza□Urquijo E, Astell AJ, Babiloni C,  
422 Bahar□Fuchs A, Bell J, Bowman GL, Brickman AM, Chételat G, Ciro C, Cohen AD,  
423 Dilworth□Anderson P, Dodge HH, Dreux S, Edland S, Esbensen A, Evered L, Ewers M, Fargo  
424 KN, Fortea J, Gonzalez H, Gustafson DR, Head E, Hendrix JA, Hofer SM, Johnson LA, Jutten R,  
425 Kilborn K, Lanctôt KL, Manly JJ, Martins RN, Mielke MM, Morris MC, Murray ME, Oh ES,  
426 Parra MA, Rissman RA, Roe CM, Santos OA, Scarmeas N, Schneider LS, Schupf N, Sikkes S,  
427 Snyder HM, Sohrabi HR, Stern Y, Strydom A, Tang Y, Terrera GM, Teunissen C, Melo van  
428 Lent D, Weinborn M, Wesselman L, Wilcock DM, Zetterberg H, O'Bryant SE. Perspectives on  
429 ethnic and racial disparities in Alzheimer's disease and related dementias: Update and areas of  
430 immediate need. *Alzheimer's & Dementia*. 2018;15(2):292-312. doi: 10.1016/j.jalz.2018.09.009.
- 431 8. Wise EA, Yan H, Oh E, Leoutsakos JM. Racial/ethnic differences in neuropsychiatric  
432 disturbances associated with incident dementia. *Alzheimer's & Dementia: Diagnosis, Assessment  
433 & Disease Monitoring*. 2024;16(3). doi: 10.1002/dad2.12615.
- 434 9. Levine TF, Kaplan NL, Chiao C, Cross CL, Caldwell JZK. Differential effects of  
435 sex/gender on brain structure in diverse Alzheimer's disease cohorts. *Alzheimer's & Dementia:  
436 Behavior & Socioeconomics of Aging*. 2025;1(3). doi: 10.1002/bsa3.70033.

- 437 10. Au R, Bose N, Imam F, Kolachalama VB. Technological advances enabling an enhanced  
438 understanding of early Alzheimer's disease. *Alzheimer's & Dementia*. 2025;21(11). doi:  
439 10.1002/alz.70883.
- 440 11. Shi D, Xie S, Li A, Wang Q, Guo H, Han Y, Xu H, Gan W-B, Zhang L, Guo T. APOE-  
441  $\epsilon 4$  modulates the association among plasma A $\beta$ 42/A $\beta$ 40, vascular diseases, neurodegeneration  
442 and cognitive decline in non-demented elderly adults. *Translational Psychiatry*. 2022;12(1). doi:  
443 10.1038/s41398-022-01899-w.
- 444 12. Jasodanand VH, Kowshik SS, Puducheri S, Romano MF, Xu L, Au R, Kolachalama VB.  
445 AI-driven fusion of multimodal data for Alzheimer's disease biomarker assessment. *Nat*  
446 *Commun*. 2025;16(1):7407. Epub 20250811. doi: 10.1038/s41467-025-62590-4. PubMed PMID:  
447 40789853; PMCID: PMC12339743.
- 448 13. Xue C, Kowshik SS, Lteif D, Puducheri S, Jasodanand VH, Zhou OT, Walia AS, Guney  
449 OB, Zhang JD, Poesy S, Kaliaev A, Andreu-Arasa VC, Dwyer BC, Farris CW, Hao H, Kedar S,  
450 Mian AZ, Murman DL, O'Shea SA, Paul AB, Rohatgi S, Saint-Hilaire MH, Sartor EA, Setty BN,  
451 Small JE, Swaminathan A, Taraschenko O, Yuan J, Zhou Y, Zhu S, Karjadi C, Alvin Ang TF,  
452 Bargal SA, Plummer BA, Poston KL, Ahangaran M, Au R, Kolachalama VB. AI-based  
453 differential diagnosis of dementia etiologies on multimodal data. *Nat Med*. 2024;30(10):2977-89.  
454 Epub 20240704. doi: 10.1038/s41591-024-03118-z. PubMed PMID: 38965435; PMCID:  
455 PMC11485262.
- 456 14. Qiu S, Miller MI, Joshi PS, Lee JC, Xue C, Ni Y, Wang Y, De Anda-Duran I, Hwang PH,  
457 Cramer JA, Dwyer BC, Hao H, Kaku MC, Kedar S, Lee PH, Mian AZ, Murman DL, O'Shea S,  
458 Paul AB, Saint-Hilaire MH, Alton Sartor E, Saxena AR, Shih LC, Small JE, Smith MJ,  
459 Swaminathan A, Takahashi CE, Taraschenko O, You H, Yuan J, Zhou Y, Zhu S, Alosco ML,  
460 Mez J, Stein TD, Poston KL, Au R, Kolachalama VB. Multimodal deep learning for Alzheimer's  
461 disease dementia assessment. *Nat Commun*. 2022;13(1):3404. Epub 20220620. doi:  
462 10.1038/s41467-022-31037-5. PubMed PMID: 35725739; PMCID: PMC9209452.
- 463 15. Kuzma A, Valladares O, Greenfest-Allen E, Nicaretta H, Kirsch M, Ren Y, Katanic Z,  
464 White H, Wilk A, Bass L, Brettschneider J, Carter L, Cifello J, Chuang WH, Clark K,  
465 Gangadharan P, Haut J, Ho PC, Horng W, Iqbal T, Jin Y, Keskinen P, Rose AL, Moon MK,  
466 Manuel J, Qu L, Robbins F, Saravanan N, Sha J, Tate S, Zhao Y, Alzheimer's Disease  
467 Sequencing P, Cantwell L, Gardner J, Chou SY, Tzeng JY, Bush W, Naj A, Kuksa P, Lee WP,  
468 Leung YY, Schellenberg G, Wang LS. NIAGADS: A data repository for Alzheimer's disease  
469 and related dementia genomics. *Alzheimers Dement*. 2025;21(6):e70255. doi: 10.1002/alz.70255.  
470 PubMed PMID: 40545618; PMCID: PMC12183107.
- 471 16. Greenfest-Allen E, Valladares O, Kuksa PP, Gangadharan P, Lee WP, Cifello J, Katanic  
472 Z, Kuzma AB, Wheeler N, Bush WS, Leung YY, Schellenberg G, Stoeckert CJ, Jr., Wang LS.  
473 NIAGADS Alzheimer's GenomicsDB: A resource for exploring Alzheimer's disease genetic and  
474 genomic knowledge. *Alzheimers Dement*. 2024;20(2):1123-36. Epub 20231026. doi:  
475 10.1002/alz.13509. PubMed PMID: 37881831; PMCID: PMC10916966.
- 476 17. McHugh CP, Clement MHS, Phatak M. AD Workbench: Transforming Alzheimer's  
477 research with secure, global, and collaborative data sharing and analysis. *Alzheimers Dement*.  
478 2025;21(5):e70278. doi: 10.1002/alz.70278. PubMed PMID: 40387289; PMCID: PMC12086970.
- 479 18. Gorgolewski KJ, Auer T, Calhoun VD, Craddock RC, Das S, Duff EP, Flandin G, Ghosh  
480 SS, Glatard T, Halchenko YO, Handwerker DA, Hanke M, Keator D, Li X, Michael Z, Maumet  
481 C, Nichols BN, Nichols TE, Pellman J, Poline JB, Rokem A, Schaefer G, Sochat V, Triplett W,  
482 Turner JA, Varoquaux G, Poldrack RA. The brain imaging data structure, a format for

- 483 organizing and describing outputs of neuroimaging experiments. *Sci Data*. 2016;3:160044. Epub  
484 20160621. doi: 10.1038/sdata.2016.44. PubMed PMID: 27326542; PMCID: PMC4978148.
- 485 19. Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, Blomberg  
486 N, Boiten JW, da Silva Santos LB, Bourne PE, Bouwman J, Brookes AJ, Clark T, Crosas M,  
487 Dillo I, Dumon O, Edmunds S, Evelo CT, Finkers R, Gonzalez-Beltran A, Gray AJ, Groth P,  
488 Goble C, Grethe JS, Heringa J, t Hoen PA, Hooft R, Kuhn T, Kok R, Kok J, Lusher SJ, Martone  
489 ME, Mons A, Packer AL, Persson B, Rocca-Serra P, Roos M, van Schaik R, Sansone SA,  
490 Schultes E, Sengstag T, Slater T, Strawn G, Swertz MA, Thompson M, van der Lei J, van  
491 Mulligen E, Velterop J, Waagmeester A, Wittenburg P, Wolstencroft K, Zhao J, Mons B. The  
492 FAIR Guiding Principles for scientific data management and stewardship. *Sci Data*.  
493 2016;3:160018. Epub 20160315. doi: 10.1038/sdata.2016.18. PubMed PMID: 26978244;  
494 PMCID: PMC4792175.
- 495 20. Diaz-Pinto A, Alle S, Nath V, Tang Y, Ihsani A, Asad M, Perez-Garcia F, Mehta P, Li W,  
496 Flores M, Roth HR, Vercauteren T, Xu D, Dogra P, Ourselin S, Feng A, Cardoso MJ. MONAI  
497 Label: A framework for AI-assisted interactive labeling of 3D medical images. *Med Image Anal*.  
498 2024;95:103207. Epub 20240515. doi: 10.1016/j.media.2024.103207. PubMed PMID: 38776843.
- 499 21. Li X, Gu Y, Dvornek N, Staib LH, Ventola P, Duncan JS. Multi-site fMRI analysis using  
500 privacy-preserving federated learning and domain adaptation: ABIDE results. *Med Image Anal*.  
501 2020;65:101765. Epub 20200702. doi: 10.1016/j.media.2020.101765. PubMed PMID:  
502 32679533; PMCID: PMC7569477.
- 503 22. Ghafarollahi A, Buehler MJ. *SciAgents: Automating Scientific Discovery Through*  
504 *Bioinspired Multi-Agent Intelligent Graph Reasoning*. *Adv Mater*. 2025;37(22):e2413523. Epub  
505 20241218. doi: 10.1002/adma.202413523. PubMed PMID: 39696898; PMCID: PMC12138853.
- 506 23. Gao S, Fang A, Huang Y, Giunchiglia V, Noori A, Schwarz JR, Ektefaie Y, Kondic J,  
507 Zitnik M. Empowering biomedical discovery with AI agents. *Cell*. 2024;187(22):6125-51. doi:  
508 10.1016/j.cell.2024.09.022. PubMed PMID: 39486399.
- 509 24. Dai H, Li Y, Liu Z, Zhao L, Wu Z, Song S, Ye S, Zhu D, Li X, Li S, Yao X, Shi L, Peng  
510 TQ, Li Q, Chen Z, Zhang D, Liu T, Mai G. AD-AutoGPT: An autonomous GPT for Alzheimer's  
511 disease infodemiology. *PLOS Glob Public Health*. 2025;5(5):e0004383. Epub 20250507. doi:  
512 10.1371/journal.pgph.0004383. PubMed PMID: 40334220; PMCID: PMC12058166.
- 513 25. Weintraub S, Salmon D, Mercaldo N, Ferris S, Graff-Radford NR, Chui H, Cummings J,  
514 DeCarli C, Foster NL, Galasko D, Peskind E, Dietrich W, Beekly DL, Kukull WA, Morris JC.  
515 The Alzheimer's Disease Centers' Uniform Data Set (UDS). *Alzheimer Disease & Associated*  
516 *Disorders*. 2009;23(2):91-101. doi: 10.1097/WAD.0b013e318191c7dd.
- 517 26. Topiwala A, Levey DF, Zhou H, Deak JD, Adhikari K, Ebmeier KP, Bell S, Burgess S,  
518 Nichols TE, Gaziano M, Stein M, Gelernter J. Alcohol use and risk of dementia in diverse  
519 populations: evidence from cohort, case-control and Mendelian randomisation approaches. *BMJ*  
520 *Evid Based Med*. 2025. Epub 20250923. doi: 10.1136/bmjebm-2025-113913. PubMed PMID:  
521 40987604.
- 522 27. Handing EP, Andel R, Kadlecova P, Gatz M, Pedersen NL. Midlife Alcohol  
523 Consumption and Risk of Dementia Over 43 Years of Follow-Up: A Population-Based Study  
524 From the Swedish Twin Registry. *The Journals of Gerontology Series A: Biological Sciences*  
525 *and Medical Sciences*. 2015;70(10):1248-54. doi: 10.1093/gerona/glv038.
- 526 28. Xu W, Wang H, Wan Y, Tan C, Li J, Tan L, Yu J-T. Alcohol consumption and dementia  
527 risk: a dose-response meta-analysis of prospective studies. *European Journal of Epidemiology*.  
528 2017;32(1):31-42. doi: 10.1007/s10654-017-0225-3.

- 529 29. Wang G, Li DY, Vance DE, Li W. Alcohol Use Disorder as a Risk Factor for Cognitive  
530 Impairment. *Journal of Alzheimer's Disease*. 2023;94(3):899-907. doi: 10.3233/jad-230181.  
531

532 **FUNDING**

533 This project was supported by grants from the National Institute on Aging's Artificial  
534 Intelligence and Technology Collaboratories (P30-AG073104 & P30-AG073105), and the  
535 National Institutes of Health (R01-AG062109, R01-AG083735, R01-HL159620, and R01-  
536 NS142076).

537

538 **ACKNOWLEDGMENTS**

539 The manuscript was originally drafted without the use of AI technologies. Subsequently, our  
540 agentic system was used to edit some portions of the manuscript. Finally, the authors reviewed  
541 and edited the content as needed and take full responsibility for the content of the publication.

542 The NACC database is funded by NIA/NIH Grant U24 AG072122. NACC data are contributed  
543 by the NIA-funded ADRCs: P30 AG062429 (PI James Brewer, MD, PhD), P30 AG066468 (PI  
544 Oscar Lopez, MD), P30 AG062421 (PI Bradley Hyman, MD, PhD), P30 AG066509 (PI Thomas  
545 Grabowski, MD), P30 AG066514 (PI Mary Sano, PhD), P30 AG066530 (PI Helena Chui, MD),  
546 P30 AG066507 (PI Marilyn Albert, PhD), P30 AG066444 (PI David Holtzman, MD), P30  
547 AG066518 (PI Lisa Silbert, MD, MCR), P30 AG066512 (PI Thomas Wisniewski, MD), P30  
548 AG066462 (PI Scott Small, MD), P30 AG072979 (PI David Wolk, MD), P30 AG072972 (PI  
549 Charles DeCarli, MD), P30 AG072976 (PI Andrew Saykin, PsyD), P30 AG072975 (PI Julie A.  
550 Schneider, MD, MS), P30 AG072978 (PI Ann McKee, MD), P30 AG072977 (PI Robert Vassar,  
551 PhD), P30 AG066519 (PI Frank LaFerla, PhD), P30 AG062677 (PI Ronald Petersen, MD, PhD),  
552 P30 AG079280 (PI Jessica Langbaum, PhD), P30 AG062422 (PI Gil Rabinovici, MD), P30  
553 AG066511 (PI Allan Levey, MD, PhD), P30 AG072946 (PI Linda Van Eldik, PhD), P30  
554 AG062715 (PI Sanjay Asthana, MD, FRCP), P30 AG072973 (PI Russell Swerdlow, MD), P30  
555 AG066506 (PI Glenn Smith, PhD, ABPP), P30 AG066508 (PI Stephen Strittmatter, MD, PhD),  
556 P30 AG066515 (PI Victor Henderson, MD, MS), P30 AG072947 (PI Suzanne Craft, PhD), P30  
557 AG072931 (PI Henry Paulson, MD, PhD), P30 AG066546 (PI Sudha Seshadri, MD), P30  
558 AG086401 (PI Erik Roberson, MD, PhD), P30 AG086404 (PI Gary Rosenberg, MD), P20  
559 AG068082 (PI Angela Jefferson, PhD), P30 AG072958 (PI Heather Whitson, MD), P30  
560 AG072959 (PI James Leverenz, MD).

561

562 **DECLARATIONS**

563 V.B.K. is a co-founder and equity holder of deepPath Inc. and Cognimark, Inc. He also serves on  
564 the scientific advisory board of Altoida Inc. The remaining authors declare no competing  
565 interests.

566

567 **KEYWORDS**

568 Artificial intelligence, Agentic AI, Hypothesis testing

569

## 570 **FIGURE CAPTIONS**

571

572 **Figure 1: Agentic AI for automated hypothesis testing.** The diagram depicts a multi-agent  
573 iterative workflow for evaluating user-provided hypotheses (e.g., “*APOE ε4 carriage increases*  
574 *risk for Alzheimer's disease across ethnicities*”). The process begins with the Data Preparation  
575 Agent, which filters relevant variables, prepares datasets for analysis, and saves processed data.  
576 If no relevant variables are found, the hypothesis is deemed “NOT TESTABLE.” Otherwise, if  
577 the null hypothesis is not rejected, the system proceeds to an outer iteration. A core inner loop  
578 involves: the Critic Agent (providing constructive criticism and code revisions), the Scientist  
579 Agent (performing computational/statistical experiments, interpreting results, and assessing  
580 significance), and the Planning Agent (creating execution plans). Additional “NOT TESTABLE”  
581 paths account for dataset limitations. The system supports collaboration with external resources  
582 (e.g., PubMed) and ADRD data inputs (cohort metadata, CSVs), enabling use cases such as  
583 replication of known associations, cross-cohort hypothesis testing, and generation of ranked,  
584 interpretable outputs.

585

586 **Figure 2: Sample execution log for a single hypothesis.** This figure presents a detailed  
587 workflow log from the multi-agentic AI system evaluating the hypothesis “*Alcohol consumption*  
588 *is associated with increased dementia risk*” on the NACC cohort (n=48,876). It showcases  
589 sequential agent outputs: Data Preparation (variable selection, missing data handling, and  
590 transformation); Planning (research strategy with chi-square tests); Scientist (statistical execution  
591 revealing significant associations,  $p < 0.0001$  across categories, Cramer's V 0.040–0.049); Critic  
592 (methodological validation); and Final Outcome (null hypothesis rejection in one iteration),  
593 highlighting automated rigor in hypothesis verification.

594

595 **Figure 3: Performance of the agentic system on NACC data.** Mean frequencies ( $\pm$  SD across  
596 multiple runs) for outcomes of 100 literature-derived ADRD hypotheses: null rejected (verified  
597 by significant association), null not rejected (no significant evidence), and not testable (due to  
598 data limitations).

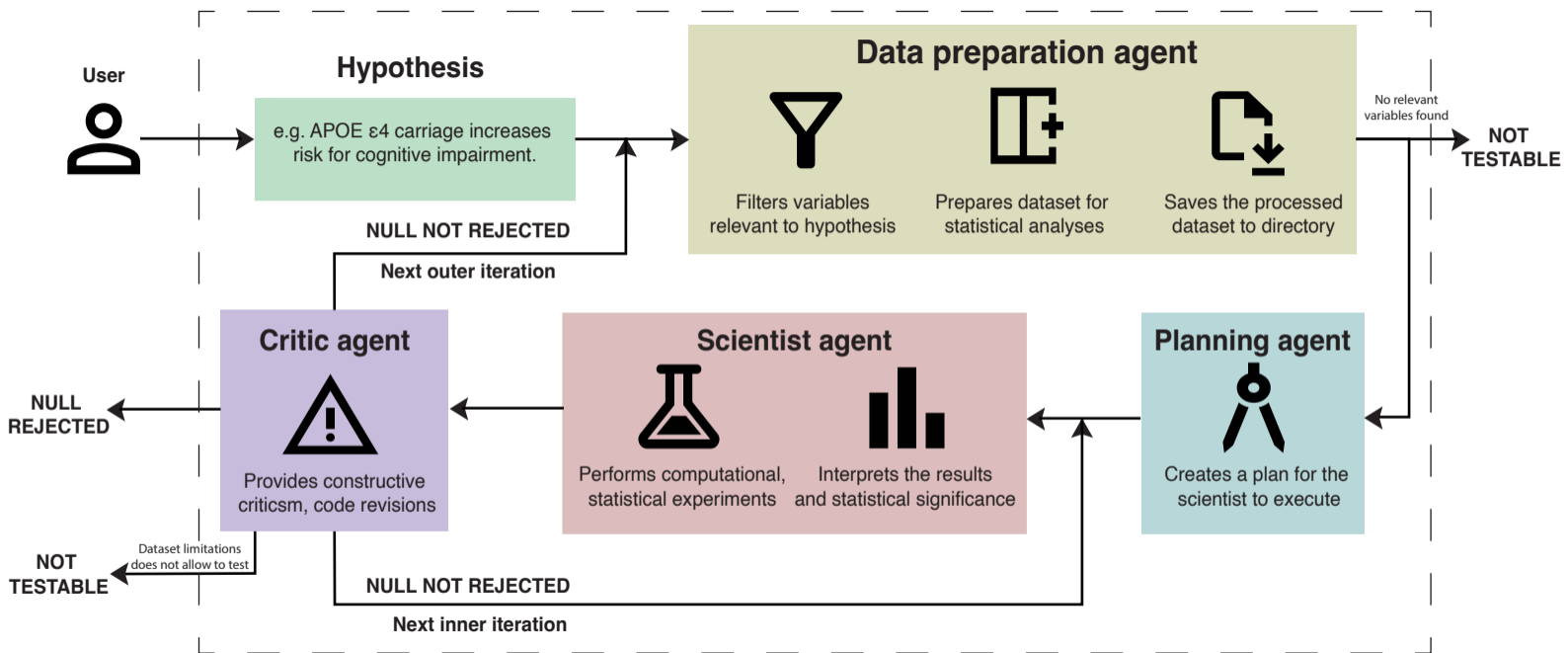
599

600 **Figure 4: Scalability of the agentic system across dataset sizes.** This plot demonstrates the  
601 mean frequencies of outcomes for 100 literature-derived ADRD hypotheses tested on  
602 subsampled NACC data at varying participant counts (100, 1,000, and 10,000 samples),  
603 categorized by null rejected (significant association verified), null not rejected (no significant  
604 evidence), and not testable (data limitations), highlighting improved verification rates with larger  
605 cohorts.

606

607 **Figure 5: Performance of the agentic system in a negative control condition.** This plot  
608 illustrates the performance of the agentic system on 100 randomly generated hypotheses from  
609 general medical literature (Supplementary Table 2) using NACC data. It shows the frequency of  
610 hypotheses where the null was rejected (indicating verification of the association), not rejected  
611 (indicating no significant evidence), and deemed not testable due to insufficient or incompatible  
612 data. The high non-testability rate serves as a negative control, validating the tool's specificity  
613 and conservative approach for out-of-domain queries.

614





**HYPOTHESIS:** Alcohol consumption is associated with increased dementia risk.

## [1] DATA PREPARATION AGENT

Dataset: NACC cohort (n=48,876 participants, 750 variables)

- Variable selection: LLM-guided selection identified 20 candidate variables from 750 available; 19 variables retained after validation
- Selected variables: Alcohol consumption measures (ALCFREQ, ALCOCCAS, ALCOHOL), demographics (NACCAGE, SEX, RACE, EDUC), dementia outcomes (NACCIDEM, NACCALZD, NACCALZP, NACCADMU, NACCFTDM, NACCBLDM), comorbidities (DEP, DEP2YRS, HYPERT, STROKE, CVOTHR, NACCETPR)
- Missing data handling: Systematic replacement of unknown values (codes 8, 9, 88, 99, 0) with NaN; 121,603 total unknown values replaced across variables
- Key patterns: ALCFREQ (2,886 missing), NACCETPR (19,721 missing), NACCADMU (48,756 missing - not assessed)
- Data transformation: Categorical variables converted to binary dummy variables
- Final dataset: 48,876 rows x 75 columns (expanded from 19 categorical variables)
- Alcohol frequency encoded as 5 binary indicators for dose-response analysis



## [2] PLANNING AGENT - Research Strategy

Objective: Test the null hypothesis that alcohol consumption is independent of dementia risk (odds ratio = 1)

- Exposure variable: Alcohol consumption frequency (5 categories: less than monthly, monthly, weekly, few times per week, daily or almost daily)
- Outcome variable: Dementia progression status (NACCIDEM\_PROCESSED\_TO\_DEMENTIA)
- Binary classification: progressed to dementia vs. did not progress
- Statistical approach: Chi-square test of independence with Cramer's V for effect size
- Rationale: Appropriate for categorical exposure x categorical outcome associations
- Alternative considered: Logistic regression (recommended for future adjusted analyses)
- Sample size: All participants with complete data for exposure and outcome variables



## [3] SCIENTIST AGENT - Statistical Analysis

Method: Chi-square tests of association between alcohol frequency and dementia progression

- Alcohol frequency categories: Less than monthly, monthly, weekly, few times per week, daily or almost daily
- Statistical test: Chi-square test of independence with Cramer's V for effect size
- Cramer's V interpretation: 0.01-0.10 = small effect, 0.10-0.30 = medium effect
- Missing data: Participants with missing values, excluded from analysis (listwise deletion)
- Statistical significance threshold:  $\alpha = 0.05$

Results: Significant association found across all frequency categories

- Less than monthly:  $p < 0.0001$ , Cramer's V = 0.040
- Monthly:  $p < 0.0001$ , Cramer's V = 0.040
- Weekly:  $p < 0.0001$ , Cramer's V = 0.044
- Few times per week:  $p < 0.0001$ , Cramer's V = 0.049 (strongest effect)
- Daily or almost daily:  $p < 0.0001$ , Cramer's V = 0.047

Clinical interpretation: Small but consistent effect sizes across all frequency categories suggest dose-response relationship; strongest association observed in "few times per week" category

Conclusion: Null hypothesis rejected across all categories



## [4] CRITIC AGENT - Methodological Review

Assessment:

- Data preprocessing: Appropriate handling of missing values (dropna applied before analysis)
- Statistical test selection: Chi-square test with Cramer's V is valid for nominal x nominal categorical data analysis
- Hypothesis testing framework: Correct use of null hypothesis testing terminology; proper interpretation of p-values and effect sizes
- Hypothesis alignment: Direction of effect supports stated hypothesis; there are consistent significant associations across all alcohol frequency categories
- Code execution: Analysis executed without errors; results reproducible
- Study limitations noted: Cross-sectional design limits causal inference; self-reported alcohol consumption may introduce recall bias; effect sizes are small but statistically significant

VERDICT: Null hypothesis rejected ( $p < 0.0001$  across all frequency categories).

Analysis methodologically sound and statistically valid



## [5] FINAL OUTCOME

Result: Null hypothesis rejected | Analysis completed in first iteration (no revisions needed)

Conclusion: Statistically significant evidence that alcohol consumption frequency is associated with increased dementia risk

- Consistent significant associations ( $p < 0.0001$ ) across all frequency categories
- Small but meaningful effect sizes (Cramer's V: 0.040-0.049) suggest clinically relevant associations
- Strongest effect observed in "few times per week" category, suggesting potential dose-response relationship
- The multi-agent system workflow ensured methodological rigor through independent agent validation



