

Improving genetic risk modeling of dementia from real-world data in underrepresented populations

Mingzhou Fu^{1,2}, Leopoldo Valiente-Banuet¹, Satpal S. Wadhwa¹, UCLA Precision Health Data Discovery Repository Working Group, UCLA Precision Health ATLAS Working Group, Bogdan Pasaniuc³, Keith Vossel¹, Timothy S. Chang^{1*}

¹ Department of Neurology, David Geffen School of Medicine, University of California, Los Angeles, Los Angeles, CA, 90095, United States

² Medical Informatics Home Area, Department of Bioinformatics, University of California, Los Angeles, Los Angeles, CA, 90024, United States

³ Department of Computational Medicine, David Geffen School of Medicine at UCLA, Los Angeles, CA, 90095, USA

*** Correspondence:** Timothy S. Chang

timothychang@mednet.ucla.edu

710 Westwood Plaza, Room 3149, Los Angeles, CA 90073

1 **Abstract**

2 **BACKGROUND:** Genetic risk modeling for dementia offers significant benefits, but studies
3 based on real-world data, particularly for underrepresented populations, are limited.

4 **METHODS:** We employed an Elastic Net model for dementia risk prediction using single-
5 nucleotide polymorphisms prioritized by functional genomic data from multiple
6 neurodegenerative disease genome-wide association studies. We compared this model with
7 *APOE* and polygenic risk score models across genetic ancestry groups, using electronic health
8 records from UCLA Health for discovery and All of Us cohort for validation.

9 **RESULTS:** Our model significantly outperforms other models across multiple ancestries,
10 improving the area-under-precision-recall curve by 21-61% and the area-under-the-receiver-
11 operating characteristic by 10-21% compared to the *APOE* and the polygenic risk score models.
12 We identified shared and ancestry-specific risk genes and biological pathways, reinforcing and
13 adding to existing knowledge.

14 **CONCLUSIONS:** Our study highlights benefits of integrating functional mapping, multiple
15 neurodegenerative diseases, and machine learning for genetic risk models in diverse populations.
16 Our findings hold potential for refining precision medicine strategies in dementia diagnosis.

17 18 **Key Words**

19 Dementia, genetic risk prediction, machine learning, electronic health record, non-European
20 population

21 **1 Background**

22 Dementia, a complex and multifaceted syndrome, is characterized by a progressive decline in
23 cognitive function beyond what might be expected from normal aging. Etiologies include
24 Alzheimer's disease (AD), vascular dementia, Lewy body dementia (LBD), Frontotemporal
25 dementia (FTD), and Parkinson's disease dementia (PDD), among others.¹ The prognosis of
26 dementia is generally a gradual and continuous decline in cognitive function, which can
27 significantly impact an individual's ability to perform daily activities.² Dementia represents a
28 significant public health concern, with a global prevalence estimated at around 36 million in
29 2020. Owing to an aging population, this number is projected to triple by 2050.³ The economic
30 burden of dementia is also substantial, with global costs estimated to be around \$594 billion
31 annually.⁴

32 Dementia has a strong genetic predisposition, with numerous significant genetic variants
33 associated with the disease identified through Genome-Wide Association Studies (GWASs). For
34 example, the Apolipoprotein E (*APOE*) gene, which encodes a protein responsible for binding
35 and transporting low-density lipids, significantly influences the risk of late-onset AD, the most
36 prevalent form of dementia.^{5,6} Similarly, the Microtubule-associated protein tau (*MAPT*) is a
37 recognized genetic mutation in FTD,⁷ and Synuclein Alpha (*SNCA*) is associated with PDD.⁸
38 While these studies have deepened our understanding of the genetic architecture of dementia,
39 additional research is necessary to successfully model personal dementia genetic risk and
40 understand the potential limitations.

41 Polygenic risk scores (PRSs), which aggregate the effects of many genetic variants associated
42 with a disease, have recently been used to quantify an individual's genetic predisposition for
43 complex diseases like dementia.⁹ A growing number of studies have underscored the robust links

44 between AD PRS and AD phenotype,¹⁰⁻¹³ declines in memory and executive function,¹⁴⁻¹⁷
45 clinical progression,¹⁵ and amyloid load¹⁸ in the non-Hispanic white population. However, the
46 performance of PRSs in non-European ancestries has been suboptimal. The weights for SNPs in
47 PRSs are predominantly calculated based on European ancestry GWASs, leading to a lack of
48 generalizability in representing genetic risks for non-European individuals.¹⁹⁻²² Using PRSs for
49 245 curated traits from the UK Biobank data, Privé et al.²³ revealed notable disparities in the
50 phenotypic variance explained by PRSs across different populations. Specifically, compared to
51 individuals of Northwestern European ancestry, the PRS-driven phenotypic variance is only
52 64.7% in South Asians, 48.6% in East Asians, and 18% in West Africans. Similarly, using a
53 population from the Health and Retirement Study, Marden et al. demonstrated that the estimated
54 effect of the AD PRS was notably smaller for non-Hispanic black compared to non-Hispanic
55 white in both dementia probability score and memory score.²⁴

56 Another limitation of current genetic risk modeling is differentiating between causal and
57 uninformative variants. Causal variants, such as *APOE* in AD, have been suggested to be
58 included as separate variables in genetic risk modeling due to their independent risk
59 contribution.²⁵ On the other hand, including uninformative, non-causal variants in prediction
60 models may introduce "noise" that obscures the effects of important variants. In a study by
61 Dickson et al.,²⁶ a model incorporating allelic *APOE* terms and just 20 additional Single-
62 Nucleotide Polymorphisms (SNPs) outperformed the model that included thousands of SNPs in
63 AD risk prediction (area under the receiver operating characteristic (AUROC): 0.75 vs. 0.63).
64 Moreover, most current studies used longitudinal cohorts, which perform extensive testing and
65 consensus criteria²⁷ applied by clinicians with expertise in dementias to determine dementia
66 diagnosis. While this approach ensures precision within research cohorts, it does not necessarily

67 mirror the practicalities of real-world community settings. In real-world clinical care, the
68 expertise in dementia may vary, and the criteria used for diagnosis may not always align with the
69 stringent standards of research cohorts. Diagnoses documented in the Electronic Health Records
70 (EHRs) capture these real-world data and, by routinely capturing patient data over extended
71 periods, form an expansive longitudinal cohort ideal for real-world research. Compared to
72 traditional cohorts, EHR cohorts offer additional benefits, such as vast sample sizes, diverse
73 phenotypes, and a more inclusive representation of often underrepresented groups, like
74 minorities and older adults.²⁸ However, only a few genetic studies on dementia have been
75 conducted within the context of EHR, and have predominantly focus on AD^{11,29}
76 Finally, prior studies have primarily focused on the genetic risk prediction of AD. However,
77 while AD accounts for a significant portion of dementia cases, concentrating solely on it risks
78 overlooking the broader scope of cognitive disorders. In real-world scenarios, many dementia
79 cases display mixed pathologies,^{30,31} with mixed dementia being a common occurrence³².
80 Addressing dementia as a whole, rather than exclusively focusing on AD, could better reflect the
81 clinical landscape and lead to interventions and therapies that benefit a larger cohort of affected
82 individuals.³³
83 Unfortunately, dementia remains significantly underdiagnosed in real-world community settings.
84 Research comparing diagnoses from real-world sources like Medicare claims or EHR to the gold
85 standard diagnoses from longitudinal cohort studies reveals a sensitivity range of just 50-65%.³⁴⁻
86 ³⁹ Early detection of all-cause dementia with genetic modeling can empower healthcare providers
87 to pinpoint the appropriate diagnostic processes, streamline care coordination, manage symptoms
88 effectively, and begin suitable treatments. The above-mentioned limitations underscore the need

89 for more refined methodologies to develop genetic risk models across diverse populations
90 accurately.

91 In the present study, we hypothesized that the risk SNPs associated with dementia, and their
92 corresponding weights, may vary across diverse populations, namely Amerindian, African, and
93 East Asian genetic ancestry. We further proposed that the prediction performance of dementia
94 phenotypes in non-European populations could be enhanced by identifying biological-
95 meaningful SNPs followed by sparse machine learning models within each genetic ancestry
96 group. Thus, we present a novel approach for assessing individual dementia genetic risks across
97 diverse populations.

98 Our approach addresses the previously noted limitations through several innovative measures.

99 Firstly, we utilized functional and biological information to prioritize SNPs based on GWAS
100 results, thereby targeting causal SNPs with the highest likelihood of contributing to dementia
101 risk. Secondly, we employed machine learning algorithms to select important genetic variants.

102 Our method allows for the fine-tuning of models across different ancestry groups, offering a
103 significant advantage for non-European populations that are often underrepresented in GWAS
104 studies. Finally, we developed and validated our models within real-world EHR settings,

105 focusing on predicting dementia as an encompassing condition. This innovative approach holds
106 promise for enhancing our understanding of individual dementia genetic risks and promoting
107 health equity in genetic research.

108 **2 Methods**

109 *2.1 Data source*

110 **2.1.1 UCLA ATLAS Community Health Initiative**

111 Our discovery cohort for model development was derived from the biobank-linked EHR of the
112 UCLA Health System.⁴⁰ The UCLA ATLAS Community Health Initiative collects biosamples
113 from participants of a diverse population. Upon obtaining patient consent, these biological
114 samples undergo genotyping using a customized Illumina Global Screening Array.⁴¹ Detailed
115 information regarding the biobanking and consenting procedures can be referenced in our
116 previous publications.^{42,43} After the genotype quality control described below, there were 54,935
117 individuals with genotype and UCLA EHR data. As all genetic data and EHRs utilized in this
118 study were de-identified, the study was deemed exempt from human subject research regulations
119 (UCLA IRB# 21-000435).

120 **2.1.2 All of Us Research Hub**

121 We validated our models and findings using All of Us Research Hub data. As one of the most
122 diverse biomedical data resources in the United States, the All of Us Research Program serves as
123 a centralized data repository, offering secure access to de-identified data from program
124 participants.⁴⁴ For our validation, we utilized data release version 7, encompassing 409,420
125 individuals, of which 245,400 have undergone whole genome sequencing.

126 2.2 *Patient genetic data preprocessing*

127 2.2.1 **Quality control**

128 The quality control process was conducted using PLINK v1.9,⁴⁵ adhering to established
129 guidelines.⁴⁰ We removed samples with a missingness rate exceeding 5%. Low-quality SNPs
130 with >5% missingness and monomorphic and strand-ambiguous SNPs were excluded. Post-
131 quality control, we performed genotype imputation via the Michigan Imputation Server.⁴⁶ This
132 step was crucial to augment the coverage of genetic variants and enable the comparison of results
133 across diverse genotyping platforms. SNPs with imputation $r^2 < 0.90$ or MAF < 1% were pruned
134 from the data. After quality control measures and imputation, there were 21,220,668 genotyped
135 SNPs across a sample of 54,935 individuals. Finally, we restricted our analyses to SNPs that
136 overlapped between UCLA ATLAS and All of Us, amounting to a total of 8,705,988 SNPs. This
137 approach ensured consistency in the genetic variables under consideration across both datasets.

138 2.2.2 **Inferring genetic ancestry**

139 Genetic ancestry refers to the geographic origins of an individual's genome, tracing back to their
140 most recent biological ancestors while largely excluding cultural aspects of their identity.⁴⁷
141 Genetic Inferred Ancestry (GIA) employs genetic data, a reference population, and inferential
142 methodologies to categorize individuals within a group likely to share common geographical
143 ancestors.⁴⁸ In our UCLA ATLAS sample, we used the reference panel from the 1000 Genomes
144 Project⁴⁹ and principal component analysis⁵⁰ to infer a patient's genetic ancestry. GIA groups
145 included European American (EA), African American (AA), Hispanic Latino American (HLA),
146 East Asian American (EAA), and South Asian American (SAA). For instance, we designated
147 individuals within the United States whose recent biological ancestors were inferred to be of

148 Amerindian ancestry as "HLA GIA".⁵¹ In addition, we calculated ancestry-specific principal
149 components within each GIA group using principal component analysis.

150 *2.3 Genetic predictors*

151 **2.3.1 GWAS selection**

152 Our study's initial step is identifying potential risk SNPs as candidate predictors for dementia
153 GWASs. A summary of the GWASs used and steps to select candidate SNPs in our study can be
154 found in **Supplementary Table 1** and **Supplementary Figure 1**.

155 We selected GWASs for AD,^{5,52,53} PDD,⁵⁴ PSP,⁵⁵ LBD,⁵⁶ and stroke⁵⁷ phenotypes. For AD
156 GWASs, we included three different GWASs conducted on diverse populations, including
157 European,⁵ African American,⁵² and multi-ancestries.⁵³ The summary statistics from all these
158 GWAS are publicly available. Detailed information regarding the recruitment procedures and
159 diagnostic criteria can be found in the original publications.

160 **2.3.2 Candidate SNPs identification and annotation**

161 A significant proportion of GWAS hits are found in non-coding or intergenic regions,⁵⁸ and
162 given the correlated nature of genetic variants in Linkage disequilibrium (LD), distinguishing
163 causal from non-causal variants often proves challenging based solely on association P-values
164 from GWASs.⁵⁹ Pinpointing the most likely relevant causal variants typically involves
165 understanding the regional LD patterns and assessing the functional consequences of correlated
166 SNPs, such as protein coding, regulatory, and structural sequences.⁶⁰ Several functionally
167 validated variants have been proved to be clinically relevant to the pathogenesis of diseases, as
168 confirmed through in vitro or in vivo experimental validation.⁶¹ To address this, we utilized the
169 Functional Mapping and Annotation of Genome-Wide Association Studies (FUMA), a tool that

170 leverages information from biological data repositories and other resources to annotate and
171 prioritize SNPs.⁵⁹

172 For each GWAS summary statistic, we first identified genomic risk loci using a P-value
173 threshold ($<5e-8$) and a pre-calculated LD structure ($r^2<0.2$) based on the relevant reference
174 population from the 1000 Genomes.⁴⁹ Subsequently, we identified two distinct sets of SNPs:

175 1. **Independent genome-wide-significant SNPs:** We selected the SNP with the most significant
176 GWAS P-value within each genomic risk locus. This process was iterated until all SNPs were
177 assigned to a risk locus cluster or considered independent.

178 2. **Independent gene-annotated SNPs:** We prioritized SNPs based on their functional
179 consequences on genes. In FUMA, the mapping from SNPs to genes was achieved by performing
180 ANNOVAR⁶² using Ensembl genes (build 85). SNPs were mapped to genes through positional
181 mapping, eQTL associations, and 3D chromatin interactions. The Combined Annotation-
182 Dependent Depletion (CADD) score⁶³ was used to select potential causal SNPs, with the SNP
183 possessing the highest CADD score within each genomic risk locus being chosen, indicating a
184 higher probability of the variant being deleterious.

185 The identified independent genome-wide-significant SNPs and independent gene-annotated
186 SNPs were subsequently used in constructing the disease PRSs and as candidate features in
187 dementia prediction models. To ensure the robustness of our findings, we also adopted a
188 stringent r^2 cut-off (<0.1) to define independent genome-wide-significant SNPs, ensuring the
189 selected SNPs were independent.

190 **2.3.3 Polygenic risk scores and *APOE-ε4***

191 We computed the disease-specific PRS as the sum of an individual's risk allele dosages, each
192 weighted by its corresponding risk allele effect size from the GWAS summary statistics, as

193 shown in the PRS equation $PRS_i = \sum_j^M \hat{\beta}_j \times dosage_{ij}$. All PRSs were then standardized to a
194 mean of 0 and a standard deviation of 1. The standardization process used the 1000 Genome
195 European genetic ancestry as the reference population, ensuring that the scores' range and values
196 are comparable across different GWASs. For each phenotype, we employed two distinct sets of
197 SNPs identified by FUMA, namely the independent genome-wide-significant SNPs and
198 independent gene-annotated SNPs, to calculate two respective PRSs: $PRS_{.psig}$ and $PRS_{.map}$.
199 The *APOE* gene has two variants, rs7412 and rs429358, which determine the three common
200 isoforms of the apoE protein: E2, E3, and E4, encoded by the $\epsilon 2$, $\epsilon 3$, and $\epsilon 4$ alleles.⁶⁴ Previous
201 research has demonstrated that out of the three polymorphic forms of *APOE*, carriers of *APOE-*
202 *e4* are at a higher risk of developing AD, and this association exhibits a dose-dependent effect.⁶⁵
203 Therefore, to quantify the *APOE* genotype in our study, we created a numerical variable,
204 "*APOE-e4count*", with the two variants mentioned above, representing the number of $\epsilon 4$ alleles
205 (0, 1, or 2) carried by each individual.

206 *2.4 Dementia definition and demographic features*

207 The primary outcome of interest was dementia, which we defined using the ICD-10 codes
208 (**Supplementary Table 2**). The demographic variables considered in our study were self-
209 reported sex and age. The age of each participant, measured in years, was calculated based on
210 their self-reported birth date and the dates of their encounters. For individuals diagnosed with
211 dementia, we determined the age at dementia onset.

212 *2.5 Analytical sample selection*

213 To focus on patients with longitudinal records, our analyses included patients with complete
214 demographic data (age and sex) who had at least two medical encounters after age 55. We also

215 applied a restriction of age at the last recorded encounter to be less than 90 as patients in the
216 UCLA EHR dataset are censored when older than 90.

217 We identified eligible dementia cases as patients with at least one encounter with a recorded
218 dementia diagnosis, provided that the initial onset of the condition occurred after age 55. To
219 qualify as an eligible control, subjects were required to meet the following criteria: 1) not have
220 any recorded dementia or related diagnoses, as determined by a set of predefined exclusion
221 phenotypes;⁶⁶ 2) age at the last recorded visit ≥ 70 , to exclude younger patients who may not
222 have manifested signs of dementia; and 3) a minimum of five years' length of records with an
223 average of at least one encounter per year, thereby minimizing the potential for bias associated
224 with misdiagnosis.

225 Upon the application of these selection criteria, the resultant sample served as the pool for
226 permutation resampling and subsequent modeling in our study.

227 *2.6 Prediction of dementia risk with machine learning models*

228 In our discovery study, we developed a series of logistic regression models to predict the binary
229 dementia phenotype in the UCLA ATLAS sample, stratified by GIA groups.

230 **2.6.1 Permutation resampling**

231 In order to fortify the reliability of our findings, we employed the permutation resampling
232 methodology to assess model performance, ascertain feature importance, and evaluate statistical
233 significance. Specifically, we conducted random sampling from the pool of eligible controls,
234 maintaining a case-to-control ratio of 1:3, and utilized the amalgamated case and control samples
235 for the following modeling process. This iterative procedure was repeated 1000 times.

236 2.6.2 Regress out demographic variable effects

237 To distinctly assess genetic influences, our analysis commenced by mitigating the impact of
238 demographic factors, encompassing age, sex, and ancestry-specific principal components (PCs),
239 from the predictive model. We first employed a logistic regression model that exclusively
240 utilized these variables to predict dementia status. Subsequently, we derived the predicted values
241 for each patient through this model. Applying an appropriate inverse link function (e.g., logit),
242 we then subtracted these predicted values from the ultimate outcome (dementia status),
243 generating an "offset" value. These offset values encapsulated the dementia status, after
244 regressing out the effects of demographic variables and genetic population structure.

245 2.6.3 Genetic prediction models

246 Next, we trained genetic risk models to predict the outcome (dementia status) with the offset
247 corrections applied in the linearized space, i.e., $\hat{y}_i = g^{-1}(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + offset_i)$,
248 where \hat{y}_i represents the predicted dementia status, and $g^{-1}(\cdot)$ is the inverse of the link function.⁶⁷
249 We compared four different sets of predictors: 1) *APOE* status, 2) AD PRS, 3) multiple PRSs,
250 and 4) smaller SNP sets with Elastic Net regularization. The latter involved the application of a
251 regularization technique known as Elastic Net to smaller sets of SNPs.⁶⁸ For multiple PRS
252 models, we crafted models utilizing diverse AD PRSs of varying ancestries or PRSs derived
253 from other GWASs focused on neurodegenerative diseases. Across all models, we employed a 5-
254 fold cross-validation methodology to authenticate their predictive efficacy, with the final results
255 reported on the combined hold-out testing set.
256 The primary assessment criterion was the Area Under the Precision-Recall Curve (AUPRC),
257 specifically chosen for its appropriateness in scenarios involving imbalanced datasets where the
258 number of cases is significantly outnumbered by controls.⁶⁹ Additionally, the AUROC was

259 reported as a comprehensive metric for model evaluation. To determine the optimal threshold,
260 we selected the point that maximized the Matthews Correlation Coefficient (MCC).²⁸ Subsequent
261 performance metrics, such as the F1 score, accuracy, precision, recall, and specificity, were
262 computed based on this threshold. The 95% confidence intervals (CIs) and p-values ($P =$
263 $\frac{1}{1000} \{metric_{model1} \geq metric_{model2}\}$) were derived through 1000 permutations as described
264 previously.

265 2.7 *Validations in the All of Us sample*

266 We conducted a validation study using the All of Us cohort to assess the generalizability of our
267 findings derived from the UCLA ATLAS sample. We selected a comparable sample from the All
268 of Us Research Hub, adhering to the same criteria and sampling scheme for the GIA groups in
269 the UCLA ATLAS sample. The same methodologies were employed to define dementia cases
270 and controls. We extracted the same genetic risk loci from the All of Us Whole Genome
271 Sequencing data for PRS construction or those identified through Elastic Net models in the
272 UCLA ATLAS sample. We employed a consistent methodology to regress out demographic
273 variables and genetic population structure (i.e., PCs) as a preliminary step. This approach was
274 undertaken to derive offset corrections, mirroring the procedures employed in our prior research.
275 By regressing out these factors, we aimed to ensure that the statistical models accurately reflect
276 the intrinsic genetic associations, unconfounded by extraneous demographic or population
277 structure influences.

278 We compared three models in the All of Us sample: 1) the *APOE-e4* model; 2) the best-
279 performing PRS model; and 3) the best-performing Elastic Net SNP model. The same evaluation
280 metrics were utilized for model comparisons.

281 2.8 *Gene mapping and gene set analysis*

282 To facilitate biological interpretations, we employed FUMA's positional, eQTL, and chromatin
283 interaction mapping to associate dementia risk SNPs, identified from the top-performing Elastic
284 Net SNP models, with specific genes.⁵⁹ We then tested these mapped genes against gene sets
285 procured from MsigDB, such as positional gene sets and Gene Ontology (GO) gene sets, to
286 assess the enrichment of biological functions through hypergeometric tests. To correct for
287 multiple testing, we implemented the Benjamin-Hochberg adjustment.⁷⁰ Using heatmaps, we
288 reported and visualized gene sets with an adjusted P-value ≤ 0.05 and more than one overlapping
289 gene.

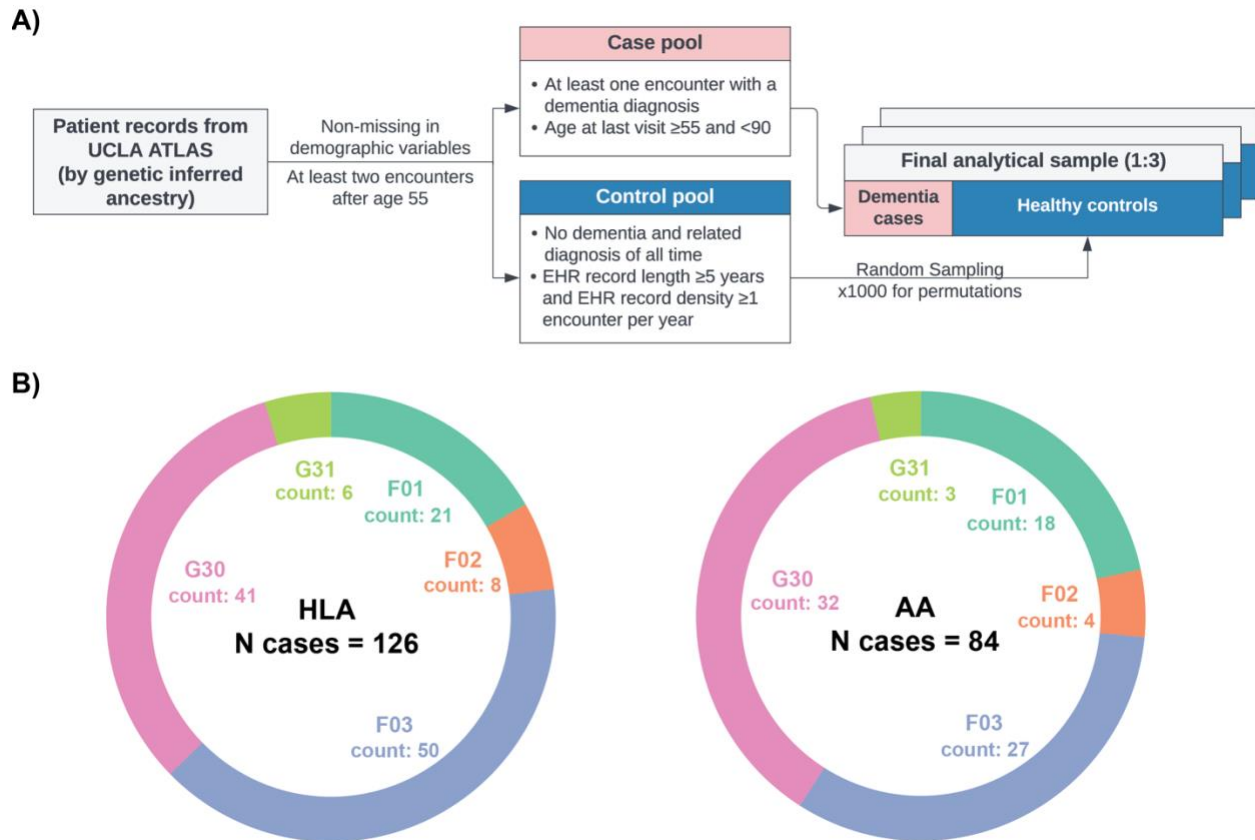
290 3 Results

291 3.1 *Sample description*

292 The study's primary dataset for model development was derived from EHR linked to the biobank
293 of the UCLA Health System.⁴⁰ A detailed depiction of the sample selection steps and resampling
294 scheme is provided in **Figure 1A**.

295 **Figure 1B** illustrates the finalized UCLA ATLAS samples, stratified by GIA groups. Notably,
296 the HLA sample comprised 610 patients, while the AA sample consisted of 440 patients, with
297 126 and 84 dementia cases, respectively, within each group. The distribution of International
298 Classification of Diseases, 10th Revision (ICD-10) diagnosis codes remained relatively
299 consistent across the two GIA samples, with Alzheimer's disease (G30) and unspecified
300 dementia (F03) being the most prevalent diagnoses. However, it is important to highlight that the
301 AA group exhibited a higher proportion of patients diagnosed with vascular dementia (F01)

302 compared to the HLA group. The EAA group, with a limited case count (N = 75), was excluded
 303 from primary analyses but included in sensitivity analyses.



304 **Figure 1. Sample selection steps and dementia patient characteristics by genetic inferred ancestry groups,**
 305 **UCLA ATLAS sample.** A) Inclusion criteria and case-control selection steps. B) Distribution of diagnosis in ICD-
 307 10 codes by genetic inferred ancestry groups. *Abbreviations: AA, African Americans; HLA: Hispanic Latino*
 308 *Americans. ICD-10 codes descriptions: G30, Alzheimer's disease; F03, Unspecified dementia; F02, Dementia in*
 309 *other diseases classified elsewhere; F01, Vascular dementia; G31, Other degenerative diseases of nervous system,*
 310 *not elsewhere classified.*

311 Within each GIA group, we found that eligible controls, due to the more stringent inclusion
 312 criteria, displayed a longer span of records and more encounters. There were no significant
 313 differences in other EHR features between dementia cases and controls (**Table 1**).

Table 1. Descriptive statistics of demographic and electronic health record features by case/control groups, UCLA ATLAS sample, stratified by genetic inferred ancestry group

	Hispanic Latino Americans (N = 610)			African Americans (N = 440)		
	Cases	Controls	P value	Cases	Controls	P value
N	126	484	-	84	356	-
Age	78.4 (71.3, 81.7)	75.3 (72.6, 79.6)	0.2	78.0 (70.1, 82.6)	75.7 (72.7, 79.9)	0.7
Sex (Female)	72 (57%)	300 (62%)	0.30	46 (55%)	218 (61%)	0.30
Span of records (in yrs)	5.9 (2.8, 8.8)	9.6 (7.7, 10.9)	<0.001*	6.2 (3.1, 10.1)	9.9 (8.1, 11.4)	<0.001*
Encounters per year	16 (7, 25)	14 (8, 20)	0.05	14 (6, 28)	13 (9, 21)	0.60

Number of encounters	73 (26, 156)	124 (73, 205)	<0.001*	65 (28, 183)	140 (84, 210)	<0.001*
Number of unique diagnosis	68 (36, 113)	71 (47, 108)	0.40	61 (41, 99)	73 (47, 103)	0.20

Notes: Continuous variables were reported as median (IQR), and categorical variables were reported as n (%). P-values were calculated based on Wilcoxon rank sum test or Pearson's Chi-squared test as appropriate. * Statistically significant at level 0.05.

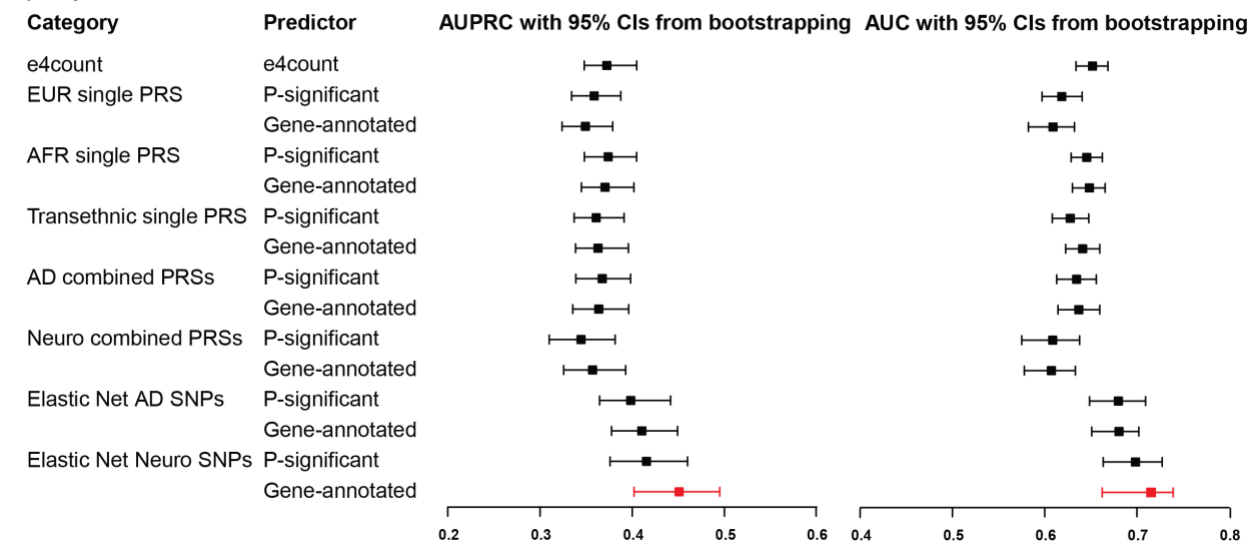
314

315 3.2 Performance comparison for dementia phenotype prediction task

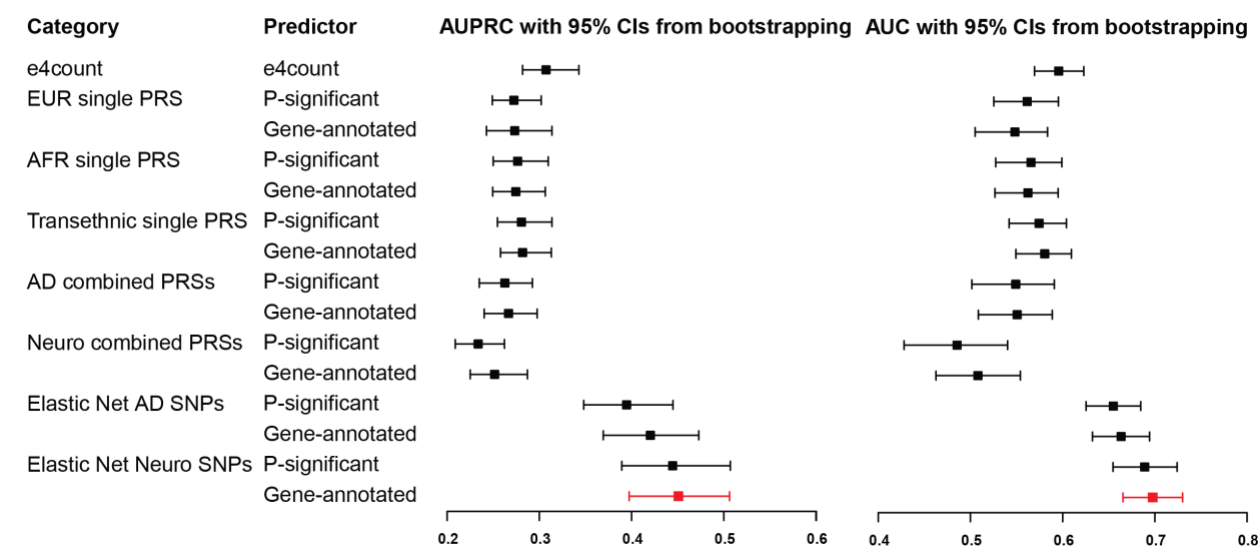
316 We developed and evaluated a series of logistic regression models to predict the binary dementia
317 phenotype within the UCLA ATLAS sample, stratified by GIA groups. After regressing out the
318 effects of age, sex, and ancestry-specific genetic variations as represented by PCs, we
319 constructed genetic risk models for dementia, incorporating offset corrections within a linearized
320 framework. The predictive capabilities of these models were assessed using four distinct sets of
321 genetic markers: 1) *APOE-e4* counts, 2) AD PRS, 3) a composite of multiple PRSs, and 4) select
322 SNPs refined through Elastic Net regularization.⁶⁸ For the selection of SNP sets, we utilized the
323 FUMA tool⁵⁹ to prioritize independent genome-wide-significant SNPs or independent gene-
324 annotated SNPs. We employed the permutation resampling methodology (1000 times) to assess
325 model performance, ascertain feature importance, and evaluate statistical significance (details see
326 **Methods**).

327 The overall performances of models for predicting dementia phenotypes are visually represented
328 in **Figure 2**. No discernible differences were observed among *APOE-e4* and all PRS models,
329 irrespective of the SNP set employed for PRS construction—whether derived from ancestry-
330 specific GWASs, genome-wide-significant SNPs, or gene-annotated SNPs. Notably, the
331 predictive performance of *APOE-e4* and all PRS models within the AA GIA sample exhibited
332 inferior outcomes compared to the HLA GIA sample, particularly evident in the AUPRC.

A) Hispanic Latino Americans



B) African Americans



333 **Figure 2. Overall model performance of *APOE-e4* count, polygenic risk score, and Elastic Net SNP models in**
 334 **dementia genetic prediction, UCLA ATLAS sample, stratified by genetic inferred ancestry group.** All models
 335 (if not other specified) have regressed out age, sex, and ancestry-specific principal components. *Abbreviations: AD,*
 336 *Alzheimer's Disease; AUROC, Area Under the ROC Curve; AUPRC, Area Under the Precision-Recall Curve; EUR,*
 337 *European; PRS, Polygenic Risk Score; SNP, Single-Nucleotide Polymorphism.*

339 Elastic Net SNP models demonstrated an overall improvement in dementia prediction across
 340 both GIA groups. The model incorporating gene-annotated SNPs from AD and other dementia-
 341 related disease GWASs emerged as the most effective, indicating a collective contribution from
 342 SNPs associated with various dementia-related diseases. Specifically, the leading Elastic Net
 343 SNP model for HLA GIA sample significantly enhanced the AUPRC by 22% (0.451 vs. 0.371,
 344 p-value = 0.003), and the AUROC by 11% (0.715 vs. 0.648, p-value = 0.008) compared to the

345 best PRS model. Furthermore, this model outperformed the *APOE-e4* count model, with
 346 increments of 21% in AUPRC (p-value = 0.003) and 10% in AUROC (p-value = 0.007).
 347 This model's efficacy was even more pronounced within the AA GIA sample, with an increase in
 348 AUPRC by 61% (p-value < 0.001) and the AUROC by 21% (p-value < 0.001) in comparison to
 349 the best PRS model. Relative to the *APOE-e4* count model, the improvements were 47% in
 350 AUPRC (p-value < 0.001) and 17% in AUROC (p-value < 0.001).
 351 We also noted a substantial enhancement in the other performance metrics (based on the
 352 threshold that maximized the MCC) of the Elastic Net SNPs models compared to other models
 353 across both GIA samples (**Supplementary Table 3**). This was evidenced by marked
 354 improvements in accuracy, precision, and the F1 score. In our sensitivity analysis, applying a
 355 more stringent r^2 cut-off (<0.1) for defining independent genome-wide-significant SNPs yielded
 356 results consistent with our initial findings, as detailed in **Supplementary Table 4**.
 357 In summary, models leveraging SNPs as features identified through machine learning methods
 358 possess the potential to surpass those relying solely on summary scores such as PRSs.
 359 Furthermore, selecting SNPs mapped to genes using functional genomic data holds promise for
 360 further refining predictive performance.

361 3.3 *Featured risk variants and mapped genes*

362 In our analysis of the best-performing Elastic Net SNPs models, we further examined the
 363 features selected by each model. The HLA and AA models identified 15 and 10 risk SNPs,
 364 respectively. A detailed list of SNPs, including related information, is provided in **Table 2**.

Table 2. Featured risk SNPs from the best-performing Elastic Net SNP model, UCLA ATLAS sample, stratified by genetic ancestry

rsID	CHR	POS	Variable Importance (percentage, 95% CI)	Nearest Gene	AD EUR	AD AFR	AD multi	LBD	PD	PSP	Stroke
Hispanic Latino American ancestry (HLA)											
rs429358	19	44908684	0.088 (0.02, 0.143)	<i>APOE</i>		x					

rs2075650	19	44892362	0.086 (0.02, 0.14)	<i>TOMM40</i>		x	x	x	
rs483082	19	44912921	0.071 (0.019, 0.113)	<i>APOC1</i>		x	x		
rs157581	19	44892457	0.06 (0.015, 0.097)	<i>TOMM40</i>		x			x
rs412776	19	44876259	0.059 (0.019, 0.099)	<i>PVRL2</i>	x		x		
rs62120578	19	44713297	0.049 (0.021, 0.075)	<i>CTB-171A8.1</i>	x				
rs4803765	19	44855191	0.045 (0.015, 0.076)	<i>PVRL2</i>	x				
rs80100206	4	705856	0.044 (0.016, 0.083)	<i>PCGF3</i>					x
rs6857	19	44888997	0.038 (0.011, 0.068)	<i>NECTIN2</i>		x			
rs2276412	11	121590137	0.032 (0.008, 0.062)	<i>SORL1</i>	x				
rs2220427	4	110793733	0.031 (0.007, 0.056)	<i>RP11-777N19.1</i>					x
rs13067212	3	39404095	0.027 (0.004, 0.055)	<i>RPSA</i>					x
rs435380	19	44903861	0.026 (0.003, 0.063)	<i>TOMM40</i>		x	x		
rs10422350	19	44725238	0.025 (0.005, 0.048)	<i>snoZ6</i>	x		x		
rs1551890	19	44829875	0.023 (0.004, 0.046)	<i>BCAM</i>	x		x		
African American ancestry (AA)									
rs2627641	19	45205500	0.092 (0.05, 0.166)	<i>BLOC1S3</i>	x				
rs8073976	17	44955857	0.077 (0.041, 0.128)	<i>CIQL1</i>					x
rs429358	19	44908684	0.065 (0.031, 0.111)	<i>APOE</i>		x			
rs77283277	7	143386852	0.064 (0.03, 0.125)	<i>ZYX</i>	x				
rs2075650	19	44892362	0.06 (0.028, 0.101)	<i>TOMM40</i>		x	x	x	
rs13032148	2	127107524	0.057 (0.02, 0.107)	<i>BINI</i>	x		x		
rs73936967	19	44890485	0.056 (0.022, 0.101)	<i>TOMM40</i>		x			
rs71352239	19	44926286	0.053 (0.023, 0.086)	<i>APOC1P1</i>	x		x	x	
rs11223641	11	133950127	0.04 (0.012, 0.064)	<i>IGSF9B</i>					x
rs435380	19	44903861	0.035 (0.004, 0.073)	<i>TOMM40</i>		x	x		

Abbreviations: AD, Alzheimer's Disease; AFR, African American; CI, confidence interval; EUR, European; LBD, Lewy body dementia; PD, Parkinson's disease; PRS, Polygenic Risk Score; PSP, progressive supranuclear palsy; SNP, Single-Nucleotide Polymorphism. **Note:** SNPs marked in red are overlapped SNPs identified by both samples.

365

366 By assessing the feature importance of the SNPs chosen by the models, we discovered that

367 rs429358 (chr19:44908684, nearest gene: *APOE*), rs2075650 (chr19:44892362, nearest gene:

368 *TOMM40*), and rs483082 (chr19: 44912921, nearest gene: *APOC1*) were selected as the top

369 three important predictor for the HLA GIA group, together accounting for ~25% of the total

370 predictive importance. Conversely, for the AA GIA group, the most influential predictors were

371 identified as rs2627641 (chr19:45205500, nearest gene: *BLOC1S3*), rs8073976

372 (chr17:44955857, nearest gene: *CIQL1*), and rs429358 (chr19:44908684, nearest gene: *APOE*).

373 Two AD-associated risk SNPs, rs429358 and rs2075650, were pinpointed by both GIA Elastic

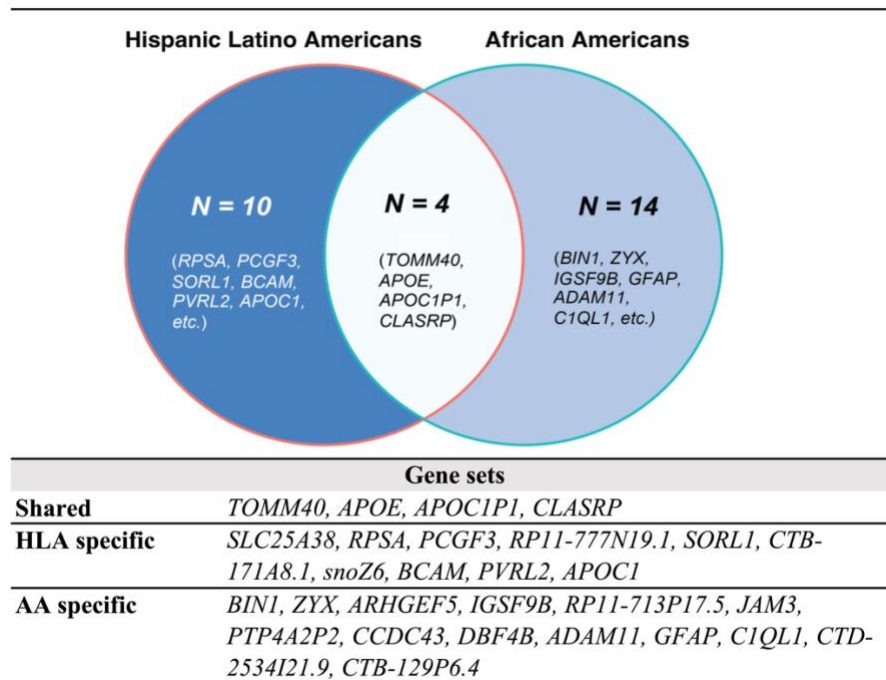
374 Net SNPs models, albeit with slight variations in their relative importance. Moreover, both

375 models identified several risk SNPs of PDD and progressive supranuclear palsy (PSP) as crucial

376 predictors of dementia. However, there were notable differences between the models. For
377 instance, the AA GIA model ascribed significant importance to a PSP-associated risk SNP,
378 rs8073976, located on chromosome 17. Interestingly, stroke-risk SNPs were only identified as
379 important predictors by the HLA GIA model, underscoring the distinct genetic underpinnings
380 influencing these different ancestry groups.

381 To better understand the biological functions and pathways associated with the identified risk
382 variants, we then mapped those featured risk SNPs to genes. This was also achieved using
383 FUMA, which incorporates positional, eQTL, and 3D chromatin mapping.⁵⁹

384 Notably, four genes were identified by both non-European GIA models (**Figure 3 &**
385 **Supplementary Table 5**). All shared genes were located near *chr19q13*, which includes the
386 well-established AD risk gene cluster, *APOE-TOMM40-APOC1*.⁷¹ According to the enrichment
387 analysis results, these shared genes are predominantly involved in biological pathways associated
388 with lipid metabolism. These pathways encompass processes such as the assembly and
389 organization of protein-lipid complexes, as delineated by the GO terms. Additionally, these
390 genes play an essential role in regulating cholesterol, triglyceride, amyloid proteins, and
391 lipoprotein particles, further underscoring the significance of lipid metabolic processes in
392 dementia. In addition, we investigated ancestry-specific genes. For instance, genes near the
393 *chr17q21* (e.g., *CCDC43*, *GFAP*, and *CIQL1*), and the *chr11q25* region (e.g., *GSF9B* and
394 *JAM3*) were uniquely pinpointed by the AA GIA model.



395
396 **Figure 3. Shared and ancestry-specific risk genes identified by the best-performing Elastic Net SNP models,**
397 **UCLA ATLAS sample.**

398 In the sensitivity analyses, we performed dementia risk modeling in the EAA GIA sample (N =
399 673). Similar to other GIA groups, the model incorporating gene-annotated SNPs from AD and
400 other dementia-related disease GWASs performed the best compared to all other models,
401 enhancing the AUPRC by 11% (0.511 vs. 0.459), and the AUC by 7% (0.754 vs. 0.703)
402 compared to the best PRS model. Despite these improvements, the differences in performance
403 between the leading Elastic Net SNP model and other models did not reach statistical
404 significance (AUPRC: p-value = 0.438; AUROC: p-value = 0.376). Among the featured 12 risk
405 SNPs, rs429358 (chr19:44908684, nearest gene: *APOE*), rs35106910 (chr19:44781009, nearest
406 gene: *CBLC*), and rs66626994 (chr19:44924977, nearest gene: *APOC1P1*) were the most
407 significant predictors for the EAA GIA group, collectively accounting for ~32% of the overall
408 predictive importance. After mapping featured SNPs to gene, we also identified the AD-risk
409 gene cluster, *APOE-TOMM40-APOC1*, as well as the gene region near *chr17q21* (e.g., *FMNL1*
410 *and SPPL2C*) (**Supplementary Table 6A-D**).

411 *3.4 Validations in the All of Us sample*

412 We conducted a validation study using the All of Us cohort to evaluate the broad applicability of
 413 our findings obtained from the UCLA ATLAS sample. A comparable sample was selected from
 414 the All of Us Research Hub, employing the same selection scheme to their corresponding GIA
 415 groups in the UCLA ATLAS sample. However, due to the limited number of eligible dementia
 416 cases (N case = 8) in the All of Us EAA GIA sample, we could only validate our models and
 417 findings in the HLA (N_case = 81, N_control = 445) and AA (N_case = 181, N_control = 2,463)
 418 samples. In contrast to the UCLA ATLAS samples, the All of Us cohort samples exhibited a
 419 younger demographic profile, with participants having comparatively shorter durations of EHR
 420 documentation and fewer recorded healthcare visits. Within each GIA sample, we found similar
 421 distributions of demographics and EHR features between dementia cases and eligible controls
 422 (**Supplementary Table 7-8**).

423 We applied the model weights trained from the UCLA ATLAS sample to the All of Us sample,
 424 stratified by GIA groups. In the comparison of three representative models, namely 1) the *APOE-*
 425 *e4* model; 2) the best-performing PRS model; and 3) the best-performing Elastic Net SNP model,
 426 our results mirrored those from the UCLA ATLAS sample, with the Elastic Net SNP model,
 427 which included gene-annotated SNPs from GWASs of AD and other dementia-related diseases,
 428 outperforming all other models in terms of the AUPRC and AUC in both the HLA and AA GIA
 429 samples (**Table 3**).

Table 3. Overall model performance of *APOE-e4* count, polygenic risk score, and Elastic Net SNP models in dementia genetic prediction in validation of All of Us sample, stratified by genetic inferred ancestry

		HLA (N = 526)		AA (N = 2,644)	
N case		Cases	Controls	Cases	Controls
	N	81	445	181	2,463
Model		AUPRC	AUROC	AUPRC	AUROC
<i>APOE</i>	e4 count	0.425 (0.39, 0.468)	0.64 (0.62, 0.67)	0.352 (0.317, 0.39)	0.603 (0.573, 0.632)

Best single AD PRS	AFR gene-annotated	0.395 (0.34, 0.484)	0.62 (0.58, 0.68)	0.347 (0.299, 0.404)	0.599 (0.549, 0.646)
Best SNPs	Gene-annotated Neuro SNPs	0.475 (0.384, 0.533)	0.69 (0.61, 0.73)	0.371 (0.328, 0.414)	0.628 (0.591, 0.66)

Abbreviations: AA, African Americans; AD, Alzheimer's Disease; AFR, African American; APOE, apolipoprotein E; AUROC, Area Under the ROC Curve; AUPRC, Area Under the Precision-Recall Curve; HLA: Hispanic Latino Americans; PRS, Polygenic Risk Score; SNP, Single-Nucleotide Polymorphism.

430

431 In particular, the Elastic Net SNP model demonstrated a substantial improvement in the AUPRC,

432 outperforming the *APOE-e4* model by 12% in AUPRC (p-value = 0.082), and the best AD PRS

433 model (AD AFR PRS.map) by 20% in AUPRC (p-value = 0.034) in the HLA GIA sample.

434 Similarly, in the AA GIA sample, the Elastic Net SNP model showed an enhancement of 5.4%

435 (p-value = 0.083) and 6.9% (p-value = 0.528) in the AUPRC over the *APOE-e4* and best AD

436 PRS model, respectively.

437 **4 Discussion**

438 Traditional genetic risk models have faced limitations in effectively capturing causal disease risk

439 variants and accurately assessing genetic risks across diverse populations. To address these

440 challenges, our present study introduces a novel approach to predicting dementia risks by

441 leveraging functional mapping of genetic data in conjunction with machine learning methods in

442 the real-world EHR setting. Our proposed method shows remarkable improvements in prediction

443 performance compared to well-known approaches like *APOE* gene and PRS models. We

444 successfully identified shared and ancestry-specific risk genes and biological pathways

445 contributing to dementia risks for each non-European GIA group. Finally, we bolstered the

446 reliability and generalizability of our findings by validating our models using a comparable EHR

447 sample from the All of Us cohort.

448 Our study highlights the significance of prioritizing biologically meaningful SNPs in genetic
449 prediction. GWASs often identify genomic regions with multiple correlated SNPs, which may
450 encompass several closely located genes. However, not all of these genes are relevant to the
451 disease.⁷² Functional annotation of genetic variants enabled us to target potential causal SNPs by
452 considering various factors, such as regional LD patterns, functional consequences of variants,
453 their impact on gene expression, and their involvement in chromatin interaction sites.⁵⁹ In our
454 models developed on UCLA ATLAS samples, we achieved significant improvements in model
455 performance by prioritizing biologically meaningful SNPs, ranging from 21-61% in AUPRC and
456 10-21% in AUROC across different GIA groups, compared to the *APOE-e4* count and the best-
457 performing PRS models. These results underscore the critical role of considering functional and
458 biological information in enhancing the performance of genetic prediction models, especially in
459 diverse populations.

460 It is worth highlighting that no discernible performance differences were observed between PRSs
461 constructed using genome-wide-significant and gene-annotated SNPs. This can be attributed to
462 the strong LD between genome-wide-significant and gene-annotated SNPs within the same
463 genomic region. As a result, these SNPs tend to have similar effect estimates in the GWASs.
464 Thus, it is expected that the PRSs built with these two sets of SNPs would exhibit a high
465 correlation (**Supplementary Table 9**), which further supports the notion that the choice of
466 genome-wide-significant or gene-annotated SNPs does not significantly impact the predictive
467 performance of the PRSs in our study.

468 Moreover, our study emphasizes the significance of incorporating risk factors from multiple
469 dementia-related diseases when developing predictive models for complex conditions like
470 dementia. Both ancestry-specific Elastic Net SNP models highlighted several PD and PSP risk

471 variants as significant predictors of dementia. This finding aligns with the well-known
472 complexity of dementia as a multifactorial disorder that shares common features with these
473 related conditions.⁷³ However, it is worth noting that including PRSs of those diseases did not
474 significantly improve the overall performance (**Figure 2**). This result is consistent with research
475 conducted by Clark et al.,⁷⁴ in which they demonstrated that a combined genetic score, which
476 incorporated risk variants for AD and 24 other traits, had an equivalent predictive power as the
477 AD PRS on its own. One possible explanation is that many traits were not dementia etiologies
478 and diluted the effects of the true causal SNPs in the models.

479 Our proposed Elastic Net SNPs models identified several shared risk factors across different
480 ancestries. Notably, a substantial proportion of the identified shared genes were found near the
481 *chr19q13* region, which is well-known for the AD risk gene cluster comprising *APOE-*
482 *TOMM40-APOC1*. These findings align with previous research,^{6,52,64} further supporting the
483 significance of this genomic region in contributing to the genetic risks associated with dementia.

484 At the same time, we have discovered compelling evidence supporting our hypothesis that risk
485 SNPs associated with dementia, along with their corresponding weights, exhibit significant
486 variations across diverse populations. Notably, our analysis of PRS models revealed that the
487 performance of PRS built with the European population GWAS was worse when predicting a
488 non-European GIA group. On the other hand, we also observed that the *APOE-e4* count model
489 performed better than most PRS models in HLA and AA GIA samples. These finding further
490 reinforces the limitations of standard PRS when applied to non-European populations, in which
491 attempting to transfer GWAS effect size from one GIA to another GIA, or when using matched
492 genetic ancestry GWAS with smaller sample size, as demonstrated in several AD and other
493 phenotype studies.⁷⁵⁻⁷⁸

494 In addition, we observed notable differences in the feature importance of various SNPs within
495 the best-performing Elastic Net models across distinct GIA groups. Consequently, this led us to
496 identify ancestry-specific genes and distinct biological pathways implicated in the genetic
497 predisposition to dementia in diverse ancestral samples. These findings highlight the uniqueness
498 of genetic risk factors and functional pathways in diverse population groups.

499 Finally, we validated our models using samples from separate EHR linked with genetic data (All
500 of Us). Our proposed Elastic Net SNP model consistently outperformed the *APOE-e4* and the
501 best PRS models. While the Elastic Net SNP model demonstrated effective performance in both
502 HLA and AA populations, we observed a decrease in the general performance and significance
503 (AUPRC and AUROC) in the All of Us sample compared to the UCLA ATLAS sample,
504 particularly in the AA samples. One potential explanation for this discrepancy is the distinct
505 population structure within each sample, as revealed by comparing patient characteristics
506 (**Supplementary Table 7**). These findings underscore the influence of population-specific
507 factors on the generalizability of genetic risk models, highlighting the critical need to account for
508 population diversity in predictive models for complex diseases.

509 Our study boasts several notable strengths that contribute to its significance and impact. Firstly,
510 machine learning techniques applied in our study allowed us to infer crucial dementia risk factors
511 for underrepresented populations, such as HLA and AA, with GWAS summary statistics from
512 extensively studied populations like Europeans. This approach enabled a deeper understanding of
513 the genetic landscape of dementia in underrepresented populations, particularly valuable given
514 the current limitations in large-sample-size GWASs specific to these groups. Secondly, we
515 fortified the robustness and generalizability of our findings through the validation of our model
516 on an independent dataset from the All of Us cohort. Furthermore, our innovative approach,

517 which incorporated biologically relevant genetic markers and functional annotations,
518 significantly enhanced the accuracy of disease prediction. This approach can be readily adapted
519 to predict other complex diseases, extending the scope of its applications and enriching our
520 understanding of diverse human populations' genetic traits.

521 However, we acknowledge certain limitations. Firstly, we observed variations in the composition
522 of dementia subtypes among different GIA groups' case samples. Consequently, the distinct
523 genes and biological pathways identified by different ancestry models should be interpreted with
524 this consideration. Secondly, although our study identified potential risk SNPs and genes
525 associated with dementia, additional experimentation is necessary to understand the precise
526 mechanisms underlying the association of these factors with dementia. Thirdly, due to the
527 limited number of dementia cases in the All of Us EAA GIA sample after applying our inclusion
528 criteria, we could only validate our models and findings in the HLA and AA samples. As a
529 result, the generalizability of our findings to the EAA ancestry is constrained.

530 In light of these limitations, further research with more extensive and diverse datasets,
531 encompassing a broader range of dementia subtypes and GIA groups is imperative to strengthen
532 the validity and applicability of our study's outcomes. Such efforts will contribute to a more
533 comprehensive understanding of the genetic complexities underlying dementia across diverse
534 populations.

535 **5 Conclusions**

536 Our study introduces a novel and robust approach to assessing individual genetic risks for
537 dementia across diverse populations in a real-world setting. Our study demonstrates the
538 importance of considering functional and biological information and population diversity when
539 developing predictive models for complex diseases like dementia. The findings from our

540 research provide valuable insights into the intricate genetic factors underlying dementia.
541 Moreover, this work opens up promising avenues for developing more accurate and efficient
542 predictive models for complex genetic traits in diverse human populations. Such advancements
543 can potentially be paired with the development of targeted treatments tailored to the specific
544 genetic profiles of individuals affected by dementia and related conditions.

545 **6 List of abbreviations**

Abbr.	Description
AA	African American
AD	Alzheimer's disease
APOE	Apolipoprotein E
AUPRC	Area Under the Precision-Recall Curve
AUROC	area under the receiver operating characteristic
CADD	Combined Annotation-Dependent Depletion
CI	confidence intervals
EA	European American
EAA	East Asian American
EHR	Electronic Health Records
FTD	Frontotemporal dementia
FUMA	Functional Mapping and Annotation of Genome-Wide Association Studies
GIA	Genetic Inferred Ancestry
GO	Gene Ontology
GWAS	Genome-Wide Association Studies
HLA	Hispanic Latino American
LBD	Lewy body dementia
LD	Linkage disequilibrium
MCC	Matthews Correlation Coefficient
PC	principal components
PDD	Parkinson's disease dementia
PRS	Polygenic risk scores
SAA	South Asian American
SNP	Single-Nucleotide Polymorphisms

546

547 **7 Declarations**

548 *7.1 Ethics approval and consent to participate*

549 All human subjects involved in this study provided informed consent, ensuring their
550 understanding and voluntary participation in the research.

551 *7.2 Consent for publication*

552 Not applicable.

553 *7.3 Availability of data and materials*

554 The Genome-Wide Association Study summary statistics data analyzed in this study are publicly
555 available. Individual electronic health record data are not publicly available due to patient
556 confidentiality and security concerns. Collaboration with the study authors who have been
557 approved by UCLA Health for Institutional Review Board-qualified studies are possible and
558 encouraged. Code is available on GitHub: [https://github.com/TSCchang-Lab/Dementia-](https://github.com/TSCchang-Lab/Dementia-prediction)
559 [prediction](https://github.com/TSCchang-Lab/Dementia-prediction). Requests for additional information can be directed to the Lead Contact: Timothy S
560 Chang (timothychang@mednet.ucla.edu).

561 *7.4 Competing interests*

562 The authors declare that the research was conducted in the absence of any commercial or
563 financial relationships that could be construed as a potential conflict of interest.

564 *7.5 Funding*

565 MF, LVB, SSW, and TSC was supported by the National Institutes of Health (NIH) National
566 Institute of Aging (NIA) grant K08AG065519-01A1 and the Fineberg Foundation. KV was

567 supported by NIH grants R01 NS033310, R01 AG058820, R01 AG075955, and R56 AG074473.

568 BP was supported by NIH grants R01HG009120, R01MH115676, and R01HG006399.

569 *7.6 Author Contributions*

570 MF, BP, KV and TSC contributed to conception and design of the study. MF, LVB, and SSW

571 performed the statistical analysis. MF wrote the first draft of the manuscript. All authors

572 contributed to manuscript revision, read, and approved the submitted version.

573 *7.7 Acknowledgments*

574 We gratefully acknowledge the resources provided by the Institute for Precision Health (IPH)

575 and participating UCLA ATLAS Community Health Initiative patients. The UCLA ATLAS

576 Community Health Initiative in collaboration with UCLA ATLAS Precision Health Biobank, is a

577 program of IPH, which directs and supports the biobanking and genotyping of biospecimen

578 samples from participating UCLA patients in collaboration with the David Geffen School of

579 Medicine, UCLA CTSI and UCLA Health. We would also like to acknowledge all participants

580 and researchers at the All of Us program. The All of Us Research Program is supported by the

581 National Institutes of Health, Office of the Director: Regional Medical Centers: 1 OT2

582 OD026549; 1 OT2 OD026554; 1 OT2 OD026557; 1 OT2 OD026556; 1 OT2 OD026550; 1 OT2

583 OD 026552; 1 OT2 OD026553; 1 OT2 OD026548; 1 OT2 OD026551; 1 OT2 OD026555; IAA

584 #: AOD 16037; Federally Qualified Health Centers: HHSN 263201600085U; Data and Research

585 Center: 5 U2C OD023196; Biobank: 1 U24 OD023121; The Participant Center: U24 OD023176;

586 Participant Technology Systems Center: 1 U24 OD023163; Communications and Engagement: 3

587 OT2 OD023205; 3 OT2 OD023206; and Community Partners: 1 OT2 OD025277; 3 OT2

588 OD025315; 1 OT2 OD025337; 1 OT2 OD025276.

589 8 References

- 590 1. 2022 Alzheimer's disease facts and figures. *Alzheimers Dement.* 2022;18(4):700-789.
591 doi:10.1002/alz.12638
- 592 2. Pandey E, Tejan V, Garg S. A novel approach towards behavioral and psychological symptoms of
593 dementia management. *ABP.* 2023;1(1):32-35. doi:10.25259/ABP_7_2023
- 594 3. Aggarwal NT, Tripathi M, Dodge HH, Alladi S, Anstey KJ. Trends in Alzheimer's Disease and
595 Dementia in the Asian-Pacific Region. *International Journal of Alzheimer's Disease.*
596 2012;2012:e171327. doi:10.1155/2012/171327
- 597 4. Pedroza P, Miller-Petrie MK, Chen C, et al. Global and regional spending on dementia care from
598 2000–2019 and expected future health spending scenarios from 2020–2050: An economic modelling
599 exercise. *eClinicalMedicine.* 2022;45. doi:10.1016/j.eclinm.2022.101337
- 600 5. Kunkle BW, Grenier-Boley B, Sims R, et al. Genetic meta-analysis of diagnosed Alzheimer's disease
601 identifies new risk loci and implicates A β , tau, immunity and lipid processing. *Nat Genet.*
602 2019;51(3):414-430. doi:10.1038/s41588-019-0358-2
- 603 6. Kulminski AM, Philipp I, Shu L, Culminskaya I. Definitive roles of TOMM40-APOE-APOC1
604 variants in the Alzheimer's risk. *Neurobiol Aging.* 2022;110:122-131.
605 doi:10.1016/j.neurobiolaging.2021.09.009
- 606 7. Younes K, Miller BL. Frontotemporal Dementia: Neuropathology, Genetics, Neuroimaging, and
607 Treatments. *Psychiatric Clinics of North America.* 2020;43(2):331-344.
608 doi:10.1016/j.psc.2020.02.006
- 609 8. Klein C, Westenberger A. Genetics of Parkinson's Disease. *Cold Spring Harb Perspect Med.*
610 2012;2(1):a008888. doi:10.1101/cshperspect.a008888
- 611 9. Duncan L, Shen H, Gelaye B, et al. Analysis of polygenic risk score usage and performance in
612 diverse human populations. *Nat Commun.* 2019;10(1):3328. doi:10.1038/s41467-019-11112-0
- 613 10. de Rojas I, Moreno-Grau S, Tesi N, et al. Common variants in Alzheimer's disease and risk
614 stratification by polygenic risk scores. *Nat Commun.* 2021;12:3417. doi:10.1038/s41467-021-22491-
615 8
- 616 11. Fu M, Chang TS. Phenome-Wide Association Study of Polygenic Risk Score for Alzheimer's
617 Disease in Electronic Health Records. *Front Aging Neurosci.* 2022;14:800375.
618 doi:10.3389/fnagi.2022.800375
- 619 12. Chaudhury S, Brookes KJ, Patel T, et al. Alzheimer's disease polygenic risk score as a predictor of
620 conversion from mild-cognitive impairment. *Transl Psychiatry.* 2019;9(1):1-7. doi:10.1038/s41398-
621 019-0485-7
- 622 13. Escott-Price V, Myers AJ, Huentelman M, Hardy J. Polygenic risk score analysis of pathologically
623 confirmed Alzheimer disease. *Ann Neurol.* 2017;82(2):311-314. doi:10.1002/ana.24999

- 624 14. Marden JR, Mayeda ER, Walter S, et al. Using an Alzheimer Disease Polygenic Risk Score to Predict
625 Memory Decline in Black and White Americans Over 14 Years of Follow-up. *Alzheimer Dis Assoc*
626 *Disord.* 2016;30(3):195-202. doi:10.1097/WAD.0000000000000137
- 627 15. Mormino EC, Sperling RA, Holmes AJ, et al. Polygenic risk of Alzheimer disease is associated with
628 early- and late-life processes. *Neurology.* 2016;87(5):481-488.
629 doi:10.1212/WNL.0000000000002922
- 630 16. Felsky D, Patrick E, Schneider JA, et al. Polygenic analysis of inflammatory disease variants and
631 effects on microglia in the aging brain. *Molecular Neurodegeneration.* 2018;13(1):38.
632 doi:10.1186/s13024-018-0272-6
- 633 17. Clark K, Leung YY, Lee WP, Voight B, Wang LS. Polygenic Risk Scores in Alzheimer's Disease
634 Genetics: Methodology, Applications, Inclusion, and Diversity. *J Alzheimers Dis.* 89(1):1-12.
635 doi:10.3233/JAD-220025
- 636 18. Tan CH, Fan CC, Mormino EC, et al. Polygenic hazard score: an enrichment marker for Alzheimer's
637 associated amyloid and tau deposition. *Acta Neuropathol.* 2018;135(1):85-93. doi:10.1007/s00401-
638 017-1789-4
- 639 19. Qiao J, Wu Y, Zhang S, et al. Evaluating significance of European-associated index SNPs in the East
640 Asian population for 31 complex phenotypes. *BMC Genomics.* 2023;24:324. doi:10.1186/s12864-
641 023-09425-y
- 642 20. Majara L, Kalungi A, Koen N, et al. Low and differential polygenic score generalizability among
643 African populations due largely to genetic diversity. *HGG Adv.* 2023;4(2):100184.
644 doi:10.1016/j.xhgg.2023.100184
- 645 21. Peterson RE, Kuchenbaecker K, Walters RK, et al. Genome-wide Association Studies in Ancestrally
646 Diverse Populations: Opportunities, Methods, Pitfalls, and Recommendations. *Cell.*
647 2019;179(3):589-603. doi:10.1016/j.cell.2019.08.051
- 648 22. Grinde KE, Qi Q, Thornton TA, et al. Generalizing polygenic risk scores from Europeans to
649 Hispanics/Latinos. *Genet Epidemiol.* 2019;43(1):50-62. doi:10.1002/gepi.22166
- 650 23. Privé F, Aschard H, Carmi S, et al. Portability of 245 polygenic scores when derived from the UK
651 Biobank and applied to 9 ancestry groups from the same cohort. *The American Journal of Human*
652 *Genetics.* 2022;109(1):12-23. doi:10.1016/j.ajhg.2021.11.008
- 653 24. Marden JR, Walter S, Tchetgen Tchetgen EJ, Kawachi I, Glymour MM. Validation of a polygenic
654 risk score for dementia in black and white individuals. *Brain and Behavior.* 2014;4(5):687-697.
655 doi:10.1002/brb3.248
- 656 25. Ware EB, Faul JD, Mitchell CM, Bakulski KM. Considering the APOE locus in Alzheimer's disease
657 polygenic scores in the Health and Retirement Study: a longitudinal panel study. *BMC Medical*
658 *Genomics.* 2020;13(1):164. doi:10.1186/s12920-020-00815-9
- 659 26. Dickson SP, Hendrix SB, Brown BL, et al. GenoRisk: A polygenic risk score for Alzheimer's
660 disease. *Alzheimer's & Dementia: Translational Research & Clinical Interventions.*
661 2021;7(1):e12211. doi:10.1002/trc2.12211

- 662 27. McKhann GM, Knopman DS, Chertkow H, et al. The diagnosis of dementia due to Alzheimer's
663 disease: recommendations from the National Institute on Aging-Alzheimer's Association workgroups
664 on diagnostic guidelines for Alzheimer's disease. *Alzheimers Dement.* 2011;7(3):263-269.
665 doi:10.1016/j.jalz.2011.03.005
- 666 28. Ho Y, Hu F, Lee P. The Advantages and Challenges of Using Real-World Data for Patient Care. *Clin*
667 *Transl Sci.* 2020;13(1):4-7. doi:10.1111/cts.12683
- 668 29. Gao XR, Chiariglione M, Qin K, et al. Explainable machine learning aggregates polygenic risk scores
669 and electronic health records for Alzheimer's disease prediction. *Sci Rep.* 2023;13(1):450.
670 doi:10.1038/s41598-023-27551-1
- 671 30. Robinson JL, Xie SX, Baer DR, et al. Pathological combinations in neurodegenerative disease are
672 heterogeneous and disease-associated. *Brain.* 2023;146(6):2557-2569. doi:10.1093/brain/awad059
- 673 31. Schneider JA, Arvanitakis Z, Bang W, Bennett DA. Mixed brain pathologies account for most
674 dementia cases in community-dwelling older persons. *Neurology.* 2007;69(24):2197-2204.
675 doi:10.1212/01.wnl.0000271090.28148.24
- 676 32. Zekry D, Hauw JJ, Gold G. Mixed Dementia: Epidemiology, Diagnosis, and Treatment. *Journal of*
677 *the American Geriatrics Society.* 2002;50(8):1431-1438. doi:10.1046/j.1532-5415.2002.50367.x
- 678 33. Dubois B, Padovani A, Scheltens P, Rossi A, Dell'Agnello G. Timely Diagnosis for Alzheimer's
679 Disease: A Literature Review on Benefits and Challenges. *J Alzheimers Dis.* 2016;49(3):617-631.
680 doi:10.3233/JAD-150692
- 681 34. Bradford A, Kunik ME, Schulz P, Williams SP, Singh H. Missed and Delayed Diagnosis of
682 Dementia in Primary Care: Prevalence and Contributing Factors. *Alzheimer Dis Assoc Disord.*
683 2009;23(4):306-314. doi:10.1097/WAD.0b013e3181a6bebc
- 684 35. Lang L, Clifford A, Wei L, et al. Prevalence and determinants of undetected dementia in the
685 community: a systematic literature review and a meta-analysis. *BMJ Open.* 2017;7(2):e011146.
686 doi:10.1136/bmjopen-2016-011146
- 687 36. Kotagal V, Langa KM, Plassman BL, et al. Factors associated with cognitive evaluations in the
688 United States. *Neurology.* 2015;84(1):64-71. doi:10.1212/WNL.0000000000001096
- 689 37. Taylor DH, Østbye T, Langa KM, Weir D, Plassman BL. The Accuracy of Medicare Claims as an
690 Epidemiological Tool: The Case of Dementia Revisited. *J Alzheimers Dis.* 2009;17(4):807-815.
691 doi:10.3233/JAD-2009-1099
- 692 38. Amjad H, Roth DL, Sheehan OC, Lyketsos CG, Wolff JL, Samus QM. Underdiagnosis of Dementia:
693 an Observational Study of Patterns in Diagnosis and Awareness in US Older Adults. *J Gen Intern*
694 *Med.* 2018;33(7):1131-1138. doi:10.1007/s11606-018-4377-y
- 695 39. Ponjoan A, Garre-Olmo J, Blanch J, et al. How well can electronic health records from primary care
696 identify Alzheimer's disease cases? *Clin Epidemiol.* 2019;11:509-518. doi:10.2147/CLEP.S206770
- 697 40. Johnson R, Ding Y, Bhattacharya A, et al. The UCLA ATLAS Community Health Initiative:
698 Promoting precision health research in a diverse biobank. *Cell Genomics.* 2023;3(1):100243.
699 doi:10.1016/j.xgen.2022.100243

- 700 41. Illumina. *Infinium Global Diversity Array-8 BeadChip | Array for Human Genotyping Screening*.
- 701 42. Lajonchere C, Naeim A, Dry S, et al. An Integrated, Scalable, Electronic Video Consent Process to
702 Power Precision Health Research: Large, Population-Based, Cohort Implementation and Scalability
703 Study. *Journal of Medical Internet Research*. 2021;23(12):e31121. doi:10.2196/31121
- 704 43. Naeim A, Dry S, Elashoff D, et al. Electronic Video Consent to Power Precision Health Research: A
705 Pilot Cohort Study. *JMIR Formative Research*. 2021;5(9):e29123. doi:10.2196/29123
- 706 44. All of Us Research Program Investigators, Denny JC, Rutter JL, et al. The “All of Us” Research
707 Program. *N Engl J Med*. 2019;381(7):668-676. doi:10.1056/NEJMsr1809937
- 708 45. Shaun Purcell, Christopher Chang. PLINK 1.9. www.cog-genomics.org/plink/1.9/
- 709 46. Das S, Forer L, Schönerr S, et al. Next-generation genotype imputation service and methods. *Nat*
710 *Genet*. 2016;48(10):1284-1287. doi:10.1038/ng.3656
- 711 47. Wagner JK, Yu JH, Ifekwunigwe JO, Harrell TM, Bamshad MJ, Royal CD. Anthropologists’ views
712 on race, ancestry, and genetics. *American Journal of Physical Anthropology*. 2017;162(2):318-327.
713 doi:10.1002/ajpa.23120
- 714 48. Johnson R, Ding Y, Venkateswaran V, et al. *Leveraging Genomic Diversity for Discovery in an*
715 *EHR-Linked Biobank: The UCLA ATLAS Community Health Initiative.*; 2021:2021.09.22.21263987.
716 doi:10.1101/2021.09.22.21263987
- 717 49. 1000 Genomes Project Consortium. 1000 Genomes (20181203_biallelic_SNV). Accessed June 22,
718 2022.
719 [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/release/20181203](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/release/20181203_biallelic_SNV/)
720 [_biallelic_SNV/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/release/20181203_biallelic_SNV/)
- 721 50. Abdi H, Williams LJ. Principal component analysis. *WIREs Computational Statistics*. 2010;2(4):433-
722 459. doi:10.1002/wics.101
- 723 51. Johnson R, Ding Y, Venkateswaran V, et al. Leveraging genomic diversity for discovery in an
724 electronic health record linked biobank: the UCLA ATLAS Community Health Initiative. *Genome*
725 *Med*. 2022;14(1):104. doi:10.1186/s13073-022-01106-x
- 726 52. Kunkle BW, Schmidt M, Klein HU, et al. Novel Alzheimer Disease Risk Loci and Pathways in
727 African American Individuals Using the African Genome Resources Panel: A Meta-analysis. *JAMA*
728 *Neurol*. 2021;78(1):102-113. doi:10.1001/jamaneurol.2020.3536
- 729 53. Jun GR, Chung J, Mez J, et al. Transethnic genome-wide scan identifies novel Alzheimer disease
730 loci. *Alzheimers Dement*. 2017;13(7):727-738. doi:10.1016/j.jalz.2016.12.012
- 731 54. Nalls MA, Blauwendraat C, Vallerga CL, et al. Identification of novel risk loci, causal insights, and
732 heritable risk for Parkinson’s disease: a meta-analysis of genome-wide association studies. *Lancet*
733 *Neurol*. 2019;18(12):1091-1102. doi:10.1016/S1474-4422(19)30320-5
- 734 55. Chen JA, Chen Z, Won H, et al. Joint genome-wide association study of progressive supranuclear
735 palsy identifies novel susceptibility loci and genetic correlation to neurodegenerative diseases.
736 *Molecular Neurodegeneration*. 2018;13(1):41. doi:10.1186/s13024-018-0270-8

- 737 56. Chia R, Sabir MS, Bandres-Ciga S, et al. Genome sequencing analysis identifies new loci associated
738 with Lewy body dementia and provides insights into its genetic architecture. *Nat Genet.*
739 2021;53(3):294-303. doi:10.1038/s41588-021-00785-3
- 740 57. Malik R, Chauhan G, Traylor M, et al. Multiancestry genome-wide association study of 520,000
741 subjects identifies 32 loci associated with stroke and stroke subtypes. *Nat Genet.* 2018;50(4):524-537.
742 doi:10.1038/s41588-018-0058-3
- 743 58. Zhu Y, Tazearslan C, Suh Y. Challenges and progress in interpretation of non-coding genetic variants
744 associated with human disease. *Exp Biol Med (Maywood).* 2017;242(13):1325-1334.
745 doi:10.1177/1535370217713750
- 746 59. Watanabe K, Taskesen E, van Bochoven A, Posthuma D. Functional mapping and annotation of
747 genetic associations with FUMA. *Nat Commun.* 2017;8(1):1826. doi:10.1038/s41467-017-01261-5
- 748 60. Kingsley CB. Identification of Causal Sequence Variants of Disease in the Next Generation
749 Sequencing Era. In: DiStefano JK, ed. *Disease Gene Identification: Methods and Protocols.* Methods
750 in Molecular Biology. Humana Press; 2011:37-46. doi:10.1007/978-1-61737-954-3_3
- 751 61. Lek M, Karczewski KJ, Minikel EV, et al. Analysis of protein-coding genetic variation in 60,706
752 humans. *Nature.* 2016;536(7616):285-291. doi:10.1038/nature19057
- 753 62. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-
754 throughput sequencing data. *Nucleic Acids Res.* 2010;38(16):e164. doi:10.1093/nar/gkq603
- 755 63. Kircher M, Witten DM, Jain P, O’Roak BJ, Cooper GM, Shendure J. A general framework for
756 estimating the relative pathogenicity of human genetic variants. *Nat Genet.* 2014;46(3):310-315.
757 doi:10.1038/ng.2892
- 758 64. Belloy ME, Napolioni V, Greicius MD. A Quarter Century of APOE and Alzheimer’s Disease:
759 Progress to Date and the Path Forward. *Neuron.* 2019;101(5):820-838.
760 doi:10.1016/j.neuron.2019.01.056
- 761 65. Safieh M, Korczyn AD, Michaelson DM. ApoE4: an emerging therapeutic target for Alzheimer’s
762 disease. *BMC Med.* 2019;17(1):64. doi:10.1186/s12916-019-1299-4
- 763 66. Denny JC, Bastarache L, Ritchie MD, et al. Systematic comparison of phenome-wide association
764 study of electronic medical record data and genome-wide association study data. *Nat Biotechnol.*
765 2013;31(12):1102-1110. doi:10.1038/nbt.2749
- 766 67. Generalized Linear Model (GLM) — H2O 3.28.0.2 documentation. Accessed December 28, 2023.
767 <https://h2o-release.s3.amazonaws.com/h2o/rel-yu/2/docs-website/h2o-docs/data-science/glm.html>
- 768 68. Zou H, Hastie T. Regularization and Variable Selection via the Elastic Net. *Journal of the Royal*
769 *Statistical Society Series B (Statistical Methodology).* 2005;67(2):301-320.
- 770 69. Davis J, Goadrich M. The relationship between Precision-Recall and ROC curves. In: *Proceedings of*
771 *the 23rd International Conference on Machine Learning - ICML '06.* ACM Press; 2006:233-240.
772 doi:10.1145/1143844.1143874

- 773 70. Ferreira JA. The Benjamini-Hochberg Method in the Case of Discrete Test Statistics. *The*
774 *International Journal of Biostatistics*. 2007;3(1). doi:10.2202/1557-4679.1065
- 775 71. Kamboh MI, Demirci FY, Wang X, et al. Genome-wide association study of Alzheimer's disease.
776 *Transl Psychiatry*. 2012;2(5):e117-e117. doi:10.1038/tp.2012.45
- 777 72. Bulik-Sullivan BK, Loh PR, Finucane HK, et al. LD Score regression distinguishes confounding
778 from polygenicity in genome-wide association studies. *Nat Genet*. 2015;47(3):291-295.
779 doi:10.1038/ng.3211
- 780 73. Santiago JA, Bottero V, Potashkin JA. Transcriptomic and Network Analysis Identifies Shared and
781 Unique Pathways across Dementia Spectrum Disorders. *International Journal of Molecular Sciences*.
782 2020;21(6):2050. doi:10.3390/ijms21062050
- 783 74. Clark K, Fu W, Liu CL, et al. The prediction of Alzheimer's disease through multi-trait genetic
784 modeling. *Frontiers in Aging Neuroscience*. 2023;15. Accessed August 3, 2023.
785 <https://www.frontiersin.org/articles/10.3389/fnagi.2023.1168638>
- 786 75. Dikilitas O, Schaid DJ, Tcheandjieu C, Clarke SL, Assimes TL, Kullo IJ. Use of Polygenic Risk
787 Scores for Coronary Heart Disease in Ancestrally Diverse Populations. *Curr Cardiol Rep*.
788 2022;24(9):1169-1177. doi:10.1007/s11886-022-01734-0
- 789 76. Sariya S, Felsky D, Reyes-Dumeyer D, et al. Polygenic Risk Score for Alzheimer's Disease in
790 Caribbean Hispanics. *Annals of Neurology*. 2021;90(3):366-376. doi:10.1002/ana.26131
- 791 77. Ruan X, Huang D, Huang J, Xu D, Na R. Application of European-specific polygenic risk scores for
792 predicting prostate cancer risk in different ancestry populations. *The Prostate*. 2023;83(1):30-38.
793 doi:10.1002/pros.24431
- 794 78. Jung SH, Kim HR, Chun MY, et al. Transferability of Alzheimer Disease Polygenic Risk Score
795 Across Populations and Its Association With Alzheimer Disease-Related Phenotypes. *JAMA Network*
796 *Open*. 2022;5(12):e2247162. doi:10.1001/jamanetworkopen.2022.47162
- 797