

Integrating Genomic and Contextual Determinants to Investigate Disparities in  
Alzheimer's Disease and Dementia

Diane Xue

A dissertation  
submitted in partial fulfillment of the  
requirements for the degree of

Doctor of Philosophy

University of Washington

2024

Reading Committee:

Alison Fohner, Chair

Elizabeth Blue

Joel D Kaufman

Program Authorized to Offer Degree:

Public Health Genetics

©Copyright 2024

Diane Xue

University of Washington

**Abstract**

Integrating Genomic and Contextual Data to Investigate Disparities in  
Alzheimer's Disease and Dementia

Diane Xue

Chair of the Supervisory Committee:

Alison Fohner

Department of Epidemiology

One in three older adults dies with dementia. Alzheimer's disease (AD), a progressive neurodegenerative disorder influenced by genetic and environmental risk factors, is the most common cause of dementia. By 2050, >12 million people in the United States will have AD, and the risk for AD is not evenly distributed across the population. To predict and prevent AD and related dementia, identify precision treatments, and reduce disparities, it is crucial to understand the influence of genetic and environmental risk factors on disease susceptibility. In the following papers, we investigated the racial/ethnic representation of participants in United States-based AD genetic studies, characterized the predictive accuracy of various polygenic risk scores for AD in a multi-ethnic cohort, and identified social, built, and physical environment determinants associated with dementia and cognition independent of and modified by genetic risk. We demonstrate that the lack of diversity in current genetic datasets results in insufficient statistical power to detect genetic variants associated with AD in non-European ancestry populations,

particularly for variants with small to moderate effect sizes. Beyond the potential for the lack of diversity to exacerbate inequalities in AD outcomes, it also leads to an incomplete understanding of the genetic architecture of AD – an argument supported by our polygenic risk score comparisons. We show that while variants identified by genome-wide associated studies are meaningful for capturing genetic risk in some populations, the overall predictive accuracy of current polygenic risk score models are limited. Finally, we show that after controlling for individual-level genetic and demographic risk factors, contextual determinants actionable at the population-level are associated with cognition and dementia risk. There is evidence that neighborhood socioeconomic status is differentially associated with dementia across groups with different genetic risk. In sum, this body of work contributes toward a better understanding of the etiology of AD in diverse populations, with an emphasis on how systemic population-level solutions can help reduce disparities in AD outcomes from both genetic and environmental perspectives.

## **Dedication**

*This dissertation is dedicated to my parents, Joel Xue (薛求真) and Sandy Xue (万分)*

## **Acknowledgements**

First, I must thank my dissertation committee members: Alison Fohner, Elizabeth Blue, Joel Kaufman, and Suman Jayadev. Alie, I feel so grateful to have found a mentor who champions every aspect of public health genetics and who has helped me weave through each of these projects. Liz, I am a genetic epidemiologist because of you. Thank you for setting my standards of mentorship, rigor, and kindness. Thank you both for being constant sources of support and inspiration from my first day to my last here at the University of Washington.

Second, thank you to Annique Atwater for helping me navigate graduate school, dealing with all my emergency emails, and supporting me through difficult days.

Third, thank you to my IPHG cohort who have become family. You have been the source of some of my happiest moments over the last five years. I am so grateful we got to do this together.

Fourth, thank you to my grant funders. Receiving the T32 and F99/K00 were pivotal moments in my career and provided me with invaluable training and research freedom. Additionally, thank you to Matthew Hawkins and the Department of Epidemiology grants team for managing my grant and helping me navigate the NIH funding process.

Finally, thank you to my partner (Adam), lifelong friends, and my family (Baba, Mama, and Ben). Your love, understanding, and endless support made this, and everything else in my life, possible.

## Table of Contents

Primary Manuscripts .....	10
<b>Paper 1: The power of representation: Statistical analysis of diversity in US Alzheimer's disease genetics data .....</b>	<b>10</b>
Abstract .....	10
Abbreviations .....	12
Introduction .....	13
Methods .....	15
Results .....	19
Discussion .....	22
References .....	28
Tables .....	34
Figures .....	36
<b>Paper 2: Polygenic risk scores for incident dementia in the Multi-Ethnic Study of Atherosclerosis .....</b>	<b>39</b>
Abstract .....	39
Abbreviations .....	41
Introduction .....	42
Methods .....	44
Results .....	44
Discussion .....	48
References .....	55
Tables .....	55
Figures .....	63
<b>Paper 3: Integrating Contextual Determinants and Polygenic Risk to Examine Dementia and Cognition in the Multi-Ethnic Study of Atherosclerosis .....</b>	<b>68</b>
Abstract .....	68
Abbreviations .....	70
Introduction .....	71
Methods .....	73
Results .....	77
Discussion .....	80

References .....	83
Tables .....	87
Figures.....	91
Supplementary Materials .....	95
<b>Paper 1 Supplement: The power of representation: Statistical analysis of diversity in US Alzheimer's disease genetics data .....</b>	<b>95</b>
Supplementary Table 1. Study-specific sample sizes for array datasets.....	95
Supplementary Table 2. Study-specific sample sizes for sequencing datasets .....	97
Supplementary Table 3. Case Rates for Power Simulations.....	98
Supplementary Table 4. Age-adjusted mortality rates for AD per racial/ethnic group.....	99
Supplementary Table 5. Minimum sample size to detect significant loci .....	100
Supplementary Table 6. AD loci global allele frequencies.....	101
Supplementary Figure 1. Simulated power curves for suggestive hits by effect size and allele frequency .....	104
Supplementary Figure 2. Simulated power curves for whole exome sequencing data. ....	105
<b>Paper 2 Supplement: Polygenic risk scores for incident dementia in the Multi-Ethnic Study of Atherosclerosis .....</b>	<b>106</b>
Supplementary Figure 1. Global Ancestry Proportions for MESA participants.....	106
Supplementary Figure 2. Clumping and Thresholding Polygenic Hazard Scores and Polygenic Risk Scores by Race.....	107
Supplementary Figure 3. Association between adjusted PRS and incident dementia stratified by proportion of NFE ancestry .....	109
Supplementary Figure 4. Association between adjusted PRS and dementia case-control status .....	110
Supplementary Figure 5. PRS predictive performance measured by AUC .....	111
Supplementary Figure 6. PRS predictive performance comparisons stratified by proportion of European ancestry.....	112
Supplementary Figure 7. Predictive performance of C+T Models, including PHS .....	113
Supplementary Figure 8. Inter-model reliability among self-reported white MESA participants.....	114
<b>Paper 3 Supplement: Integrating Contextual Determinants and Polygenic Risk to Examine Dementia and Cognition in the Multi-Ethnic Study of Atherosclerosis.....</b>	<b>116</b>
Supplementary Table 1. Sample size and dementia case rate in groups stratified by genetic risk.....	116

Supplementary Table 2. Effect sizes and confidence intervals of contextual determinants on dementia status in genetically stratified groups .....	117
Supplementary Table 3. Posterior Inclusion Probabilities .....	118
Supplementary Figure 1. Distribution of contextual determinants .....	119
Supplementary Figure 2. Distribution of built environment variables after log2 transformation .....	120
Supplementary Figure 3. Nonlinear interactions between Neighborhood SES and genetic risk on incident dementia risk .....	121
Supplementary Figure 4. The effect of single exposures on the overall mixture .....	123
Supplementary Figures 5-10. Nonlinear interaction plots .....	124

## Primary Manuscripts

Paper 1: The power of representation: Statistical analysis of diversity in US Alzheimer's disease genetics data

### **Abstract**

**Background:** Alzheimer's disease (AD) is a complex disease influenced by genetics and the environment. More than 75 susceptibility loci have been linked to late-onset AD, but most of these loci were discovered in genome-wide association studies (GWAS) exclusive to non-Hispanic White individuals. There are wide disparities in AD risk across racially stratified groups, and while these disparities are not due to genetic differences, underrepresentation in genetic research can further exacerbate and contribute to their persistence. We investigated the racial/ethnic representation of participants in United States-based AD genetics and the statistical implications of current representation.

**Methods:** We compared racial/ethnic data of participants from array and sequencing studies in US AD genetics databases, including National Institute on Aging Genetics of Alzheimer's Disease Data Storage Site (NIAGADS) and NIAGADS Data Sharing Service (dssNIAGADS), to AD and related dementia (ADRD) prevalence and mortality. We then simulated the statistical power of these datasets to identify risk variants from non-White populations.

**Results:** There is insufficient statistical power (probability <80%) to detect single nucleotide polymorphisms (SNPs) with low to moderate effect sizes (odds ratio [OR]<1.5) using array data from Black and Hispanic participants; studies of Asian participants are not powered to detect variants OR  $\leq$  2. Using available and projected sequencing data from Black and Hispanic participants, risk variants with OR = 1.2 are detectable at high allele frequencies. Sample sizes remain insufficiently powered to detect these variants in Asian populations.

**Discussion:** AD genetics datasets are largely representative of US ADRD burden.

However, there is a wide discrepancy between proportional representation and statistically meaningful representation. Most variation identified in GWAS of non-Hispanic White individuals have low to moderate effects. Comparable risk variants in non-White populations are not detectable given current sample sizes, which could lead to disparities in future studies and drug development. We urge AD genetics researchers and institutions to continue investing in recruiting diverse participants and use community-based participatory research practices.

## **Abbreviations**

Africans and African Americans that cluster with the African Genome Resources Reference (AA),

American Indian/Alaska Native (AI/AN),

Alzheimer's disease (AD),

Alzheimer's Disease Sequencing Project (ADSP),

genome-wide association study (GWAS),

community-based participatory research (CBPR),

non-Hispanic White individuals who cluster with 1KG European ancestry groups (EUR),

odds ratio (OR),

Whole exome sequencing (WES),

Whole genome sequencing (WGS).

## Introduction

By 2050, >12 million people in the US will have Alzheimer's disease (AD), and risk for AD is not evenly distributed across the population. A recent study of age-adjusted incidence among 1.8 million veterans found substantially higher rates of dementia among self-reported Hispanic and Black participants compared to American Indian or Alaska Native (AI/AN) participants, with the lowest incidence among Asian and White participants: estimates ranged from 20.7, 19.4, 14.2, 12.4, and 11.5 per 1000 person-years, respectively. While AD and dementia are not synonymous, AD is the most common cause of dementia, accounting for 60-80% of dementia cases[1]. These results are consistent with previous studies that have shown disparities in AD outcomes across racialized groups[2]. Much of these differences arise from inequalities in social determinants, such as those influencing education and risk for cardiovascular disease and hypertension, which are known risk factors for AD and dementia[3,4].

AD risk is also strongly influenced by genetics. AD is a complex, highly heritable ( $h^2 = 58-79\%$ ) disease[5]. To date, >75 susceptibility loci have been implicated in late-onset AD[6,7]. Due to differences in linkage disequilibrium and both genetic and non-genetic modifiers, genetic architecture of AD differs across ancestry groups in terms of associated variants and effect sizes of commonly implicated variants[8–11]. Because demographic histories create structure in human genetics[12,13], differences in allele frequency and linkage disequilibrium across global populations correlate with racialized groups[10]. While wide disparities in AD risk across racially stratified groups are not caused by genetic differences, inequality in genetic research can further exacerbate health disparities and contribute to their persistence. Most known risk variants were discovered in genome-wide association studies (GWAS), which now include >1 million participants, primarily focused on self-described non-Hispanic White individuals who cluster

with 1KG European ancestry groups (EUR)[6,14]. Meanwhile, the largest GWAS of Africans and African Americans that cluster with the African Genome Resources Reference (AA) included a mere 2,784 cases and 5,222 controls[15] – and studies include even fewer participants for other populations. This disparity translates to an understanding of AD genetic architecture that is both incomplete and inequitable[15–17].

Because this paper focuses on representation in genetic studies, it is important to distinguish between biological and social population descriptors[18]. Race/ethnicity are socially constructed without biological meaning while genetic ancestry refers to the continental or geographic origins of biological ancestors[10]. Our study relies on population descriptors aligned with social categorizations for both practicality and future use of findings. Our study is based on previous data collection efforts, and the demographics reported in previous studies are typically socio-political categorizations in adherence with U.S. Office of Management and Budget standards[19]. Furthermore, because genetic ancestry is not known at the time of recruitment, and barriers and willingness to participate in genetic research are more closely related to social and environmental differences, using social categorization when describing participants is more relevant for future applications[20]. When the studies contributing to our work describe procedures of using genetic ancestry information to filter participants, such as using principal components, we will describe the genetic reference used for filtering in addition to the self-reported or ascribed racial/ethnic categorization used for recruitment (ie. EUR = non-Hispanic White individuals who cluster with 1KG European ancestry groups).

Over a decade since the launch of the National Alzheimer’s Project Act, we are on the cusp of its initial goal to prevent and effectively treat AD by 2025. Now is a critical time to assess the state of representation and diversity in AD genetics research. The NIH allocates >\$3

billion annually to deepen our understanding of AD and facilitate the development of effective treatments. It is crucial, however, to ensure equity in who is benefiting from this extensive investment. The NIH has devoted resources to this effort, including the launch of Outreach Pro (<https://outreachpro.nia.nih.gov/>), which provides study recruitment materials in multiple languages and funding the Alzheimer’s Disease Sequencing Project (ADSP) Follow-Up Study 2.0 Diversity Initiative Phase, which is committed to identifying therapeutic targets benefitting a diverse population[21]. Here, we investigate how well US-based AD genetic datasets represent the racial and ethnic demographic characteristics of those living with AD in the US, and whether current and planned AD genetics studies are adequately powered to advance racial/ethnic equity in our understanding of the genetic architecture of AD. We conclude by offering suggestions for future recruitment priorities for AD genetics studies.

## **Methods**

### *Quantifying AD Burden in the US*

We aimed to quantify the demographics of AD burden in the U.S. by estimating disease prevalence by race and ethnicity. We categorized individuals by self-reported or ascribed race/ethnicity into five groups defined by the U.S. Office of Management and Budget[19], which guides how the federal government collects ethno-racial data: AI/AN, Asian, Black, Hispanic or Latino, and White. Participants who identified as “other” were excluded from analysis.

The most widely reported AD prevalence estimates are based on forward projections derived from the Chicago Health and Aging Project (CHAP)[22]. This study estimated prevalence for non-Hispanic White, Black, and Hispanic individuals but did not estimate prevalence for Asian or AI/AN groups. To approximate AD burden in Asian and AI/AN peoples,

we used estimates of dementia prevalence for AI/AN, Asian, Black, Hispanic, and White individuals based on Medicare Fee-for-Service beneficiaries and the U.S. Census data[23].

Because dementia prevalence includes non-AD dementia, we analyzed AD mortality as a supplemental measure of public health burden. We obtained de-identified age-adjusted mortality data for AD in the US from the Centers for Disease Control Wide-Ranging Online Data for Epidemiologic Research (CDC WONDER) Underlying Cause of Death database[24]. CDC WONDER data are based on death certificates for US residents, collected from 1999-2020. This dataset considers one underlying cause of death per person. Deaths for 1999 and beyond are classified using the Tenth Revision of the International Classification of Disease (ICD). Race and ethnicity are obtained either from self-report prior to death or reported by surviving next of kin, an informant, or by observation. We queried crude and age-adjusted death rates due to AD by race, Hispanic ethnicity, and year for the most recent five years of data availability (2016-2020) using the same five racial/ethnic categories defined by the OMB guidelines. Crude proportions can be compared to the other sources of AD population demographic data, but because racial and ethnic groups represent different proportions of the US population and have different average age-at-death, we also evaluated age-adjusted mortality rates.

Comparison between proportions of disease burden from the three sources was performed using a Chi-square test for proportions.

#### *Quantifying Racial/Ethnic Representation in Genetic Datasets*

We obtained demographic data for participants in US AD genetic studies within the National Institute on Aging Genetics of Alzheimer's Disease Data Storage Site (NIAGADS, <https://www.niagads.org/>) and the NIAGADS Data Sharing Service (dssNIAGADS, <https://dss.niagads.org/>). NIAGADS is responsible for harmonizing and sharing AD genetics,

genomics, and phenotypic data derived from NIA-funded AD genetics studies. Access to this publicly available data was approved by NIAGADS. We reviewed all genotype array datasets within NIAGADS that met the following criteria: Disease = “AD,” Molecular Data = “Genotype,” Type = “GWAS.” Array datasets from dssNIAGADS were selected by filtering for Disease = “AD” and Data Type = “GWAS”. Some AD GWAS data predate NIAGADS and are stored elsewhere. We therefore also include the following US-based AD GWAS data sets with clinical phenotyping and race/ethnicity data not captured in NIAGADS: African American Alzheimer’s Disease Genetics Study, Alzheimer’s Disease Neuroimaging Initiative, BIOCARD, Cohorts for Heart and Aging Research in Genomic Epidemiology consortium (CHARGE; includes Atherosclerosis Risk in Communities Study, Cardiovascular Health Study, and Framingham Heart Study), and the Genetic and Environmental Risk Factors for Alzheimer Disease Among African Americans Study. All participant demographic data for AD sequencing studies are from the Alzheimer’s Disease Sequencing Project (ADSP) Umbrella Study, obtained from dssNIAGADS, representing 25 datasets. In addition to sequencing data that is currently available, we analyzed the reported demographics of whole genome sequencing data that ADSP has planned for release through 2027.

Race/ethnicity data were extracted directly from study-specific covariate files where possible, or either approximated from study-specific publications or obtained directly through correspondence with study coordinators. Ancestry, race, and ethnicity labels have been inconsistently used across AD GWAS. The largest GWAS with individuals who identify as White uses “European ancestry” as inclusion criteria[14,16,17,25], while the largest GWAS with individuals who identify as Black is referred to as a study of “African American” individuals[15]. We therefore grouped labels when necessary; ex., “Caucasian” was grouped

with White, “African American” into Black. All Hispanic participants were evaluated exclusively as Hispanic and were not included in a racial category (ie. White = non-Hispanic White).

We evaluated participant demographics separately for array, whole-exome (WES), and whole-genome sequencing (WGS) data, as these could be considered different types of biological data. Array-based data is restricted to a pre-selected array of single nucleotide polymorphisms (SNPs). Sequencing data is a read-out of every base-pair in one’s exome or genome. Array data was more popular historically due to the relative ease and lower cost of conducting the assays, but more recent studies have favored WES and WGS data as costs have decreased and technology has improved. Some individuals represented in array data are represented in sequencing data.

Chi-square tests for proportions were used to compare disease burden in the population and racial/ethnic representation in AD genetics datasets.

#### *Determining Statistical Power of Existing and Planned Data*

We conducted power analyses for hypothetical GWAS of AD case-control status stratified by race based on demographics of participants across all available datasets. GWAS continue to be the dominant method used to identify risk alleles in populations. Power was simulated separately for array and sequencing data using the R function *genpwr::genpwr.calc* (version 1.0.4), available on CRAN. We assumed an additive model and simulated case rates based on population-specific case proportions of each dataset. More information on study case proportions and *genpwr* case rate selection are included in the supplement (**Supplementary Tables 1-3**). We simulated power to identify variants with odds ratios (OR) equal to 1.1, 1.2, 1.5, and 2 given significance levels of  $p < 5e-08$  (genome-wide) and  $p < 2.5e-06$  (exome-wide), and a

continuous range of minor allele frequencies from 0 to 0.5. We define “low” effect size as OR = 1.1 or less (ex., *ACE*[16]), “modest” effect size as OR = 1.2 (ex., *BINI*[26]), “intermediate” as OR = 1.5 (ex., *NCK2*[27]), and “high” effect size as OR = 2 (ex., *TREM2*[28,29]) or more. These designations follow the most recent comprehensive review of the genetic architecture of AD[6].

## Results

### *Quantifying AD Burden in the US*

Approximately 6.7 million adults aged 65 and older are currently living with AD in the US with the following distribution across racial and ethnic groups: 70.8% White, 17.4% Black, and 11.7% Hispanic (**Table 1**). Because Asian and AI/AN individuals were not represented in CHAP, we extended our analyses to dementia prevalence values[23], the majority of which represent AD (60-80%[1]). Dementia prevalence estimates were consistent with the AD prevalence estimates, where those categorized as White made up most of projected dementia cases (72.7%), followed by Black (12.6%), Hispanic (10.3%), Asian (3.7%), and AI/AN (0.6%) (**Table 1**). US cause-of-death estimates from the CDC WONDER database indicate 83.4% of AD deaths were among individuals identified as White, a slightly higher proportion than our AD and dementia prevalence estimates, followed by Black (7.6%), Hispanic (6.4%), Asian (2.4%), and American Indian/Alaskan Native (0.3%).

Racial and ethnic groups represent different proportions of the US population and have different average age-at-death, which can be accounted for in age-adjusted mortality rates. AD mortality rates per 100,000 individuals were as follows: White (254.2), Black (223.9), Hispanic (213.7), AI/AN (151.8), and Asian (125.6) (**Supplementary Table 4**). The proportional representation of racial and ethnic groups across AD prevalence, dementia prevalence, and AD mortality did not significantly differ ( $X^2 = 10.711$ ,  $p = 0.2186$ , Bonferroni-corrected  $\alpha = 0.0167$ ).

### *Quantifying Proportional Representation of Existing Genetic Data*

AD GWAS studies using array data are proportionally representative of AD (**Figure 1, Table 1**). We identified 36 genotype array datasets encompassing 65,733 individuals (**Supplementary Table 1**); among them, 77% of participants are classified as White, 14.4% Black, 6.8% Hispanic, 1.8% Asian, and 0.02% AI/AN. These proportions are similar to those in our AD and dementia prevalence estimates above, and do not differ significantly (**Table 2**).

The currently available sequencing data are more diverse than the array data, mostly due to better representation of Hispanic populations. Figure 1 displays 17 available WGS sample sets that are part of the ADSP Umbrella encompassing 36,336 individuals (**Supplementary Table 2**); among these participants, 45.0% are classified as White, 15.7% Black, and 31.1% Hispanic, while Asian and AI/AN participation remains low (7.8% and 0.4%, respectively.) Almost all Asian participants with genetic sequencing are from the Harmonized Diagnostic Assessment of Dementia for the Longitudinal Aging Study of India (LASI-DAD) study, a subset of the Longitudinal Aging Study in India. Unlike other included studies, LASI-DAD participants are not from the US, but the study is funded and administered by US institutions and investigators and data is stored in US repositories. While the inclusion of LASI-DAD is a significant improvement for South Asian representation compared to the array data, there is little improvement in representation for East and Southeast Asians. The proportion of each racial/ethnic group in the WGS studies significantly differed from the racial/ethnic proportions of AD prevalence, dementia prevalence, and AD mortality (**Table 2**).

### *Determining Statistical Power of Existing Genetic Data*

We conducted power calculations using the sample sizes derived from existing array and sequencing data as well as planned WGS data releases to ascertain the ability to identify

association signals in GWAS stratified by race/ethnicity. Power calculations simulated genotype array data using the following sample sizes: 50,000 non-Hispanic White, 8,600 Black, 1,800 Hispanic, and 1,200 Asian participants, while calculations for currently available WGS sequencing data simulated data for 16,300 non-Hispanic White, 5,700 Black, 11,300 Hispanic, and 2,800 Asian participants. WGS available thru 2027 is projected to be approximately 27,000 non-Hispanic White, 18,400 Black, 29,800 Hispanic, and 7,600 Asian participants. These estimated sample sizes represent a best-case-scenario assuming all samples meet quality control standards and are not duplicated across sample sets within array and sequencing data. Power simulations of WES data used the following sample sizes: 13,500 non-Hispanic White, 4,400 Black, and 2,200 Hispanic participants. Available genetic data for AI/AN were too small to identify any genome-wide significant hits using either genotyping array or sequencing data of any frequency or effect size ( $p < 5E-08$ ). Similarly, we were unable to model power to detect exome-wide significant hits ( $p < 2.5E-06$ ) with current sample sizes of Asian participants.

Based on sample sizes of existing genotyping array data, only studies of non-Hispanic White individuals have adequate sample sizes to detect variants with low effect sizes at genome-wide significant or suggestive thresholds (**Figure 2, Supplementary Figure 1, Supplementary Table 5**). Sample sizes comparable to existing array data alone from Black and Hispanic individuals have insufficient statistical power ( $\text{Pr} [p < 5E-08] < 80\%$ ) to detect variants with low effect size ( $\text{OR} = 1.1$ ), even when these variants are very common (frequency  $\sim 0.5$ ). In the case of current array sample sizes of Hispanic participants, statistical power is only adequate to identify common variants with high effect sizes ( $\text{OR} \geq 2$ ).

WGS and WES samples remain smaller than available array data, leading to studies that remain underpowered to detect variants of low or moderate effect sizes in studies of Black and

Hispanic participants (**Figure 2, Supplementary Figure 2, Supplementary Table 5**). However, the statistical power will certainly improve as sequencing data from Black and Hispanic participants are projected to dramatically increase in the next five years. While sequencing data from Asian individuals will more than double in the next five years, studies will remain underpowered to detect common variants of low or even moderate effects.

## **Discussion**

There is a wide discrepancy between proportional representation and statistically meaningful representation in AD genetic datasets. While racial/ethnic representation in older array datasets are largely comparable to proportions of AD burden in the US, proportional sampling results in inherently unequal understanding of genetic architecture across populations – evident in the striking lack of statistical power to find genetic variants with modest effects on AD risk using all available data from non-White populations. Participation in AD sequencing studies is poised to be enriched for individuals from historically underrepresented groups relative to their proportions in AD epidemiological data. The “oversampling” is justified and necessary – these recruitment efforts have substantially increased the power of GWAS to identify AD variants with modest to intermediate effect sizes in Black and Hispanic populations. Similar population-specific breakthroughs in Asian and AI/AN populations will lag as sample sizes remain insufficient for comparable discoveries. Notably, most variants identified thus far in GWAS of EUR populations have low effect sizes, and comparable discoveries in other populations continue to be unidentifiable with current sample sizes.

While most risk variants are not exclusive to any one ancestry background – AD associated SNPs first discovered in GWAS of EUR populations have been identified in other populations and vice versa (**Supplementary Table 6**)[30,31] – gaps in statistical power continue

to undermine our overall understanding of disease. Studying genomes with diverse ancestry is necessary for the discovery of novel risk variants; genomes with AA ancestry capture much more genetic diversity, with significant variation that is not present in the EUR genomes[32]. Indeed, association studies conducted in Caribbean Hispanics and African Americans with 1KG-YRI-like variation have identified common variants in *FBXL7* and *ABCA7* not replicated in EUR due to differences in allele frequency[33,34].

Disparities in genetic knowledge have implications for downstream applications including risk prediction and understanding underlying disease biology, drug development, and elucidating causal relationships between non-genetic risk factors and AD risk. Numerous papers have described disparities in predictive performance across diverse populations when using genetic risk prediction models developed using summary statistics from GWAS of EUR individuals[35,36]. This can result in inequalities in the ability to accurately identify individuals at high risk of disease for risk stratification in clinical trials or interventions.

There is not a one-size fits all approach for recruiting diverse participants. For example, while mistrust of biomedical research resulting from historical events (e.g., Tuskegee Syphilis Study, HeLa cells)[37] are often cited for low participation among Black and AI/AN people, recent studies have shown that low invitation rates may be to blame for low participation among Black individuals[38]. Meanwhile, American Indian/Alaska Native communities report a lack of involvement in study planning and use of research methods that do not respect community traditions – leading to hesitancy about participating in genomics research[39,40]. Furthermore, Asian and Hispanic participants have identified language and cultural barriers in study materials and communication as hindering their participation in genetic studies. For example, the use of the Spanish word *demenxia* in study materials can dissuade participation of Hispanic participants

because the meaning of *demencia* is close to “crazy”[41]. Thus, efforts to increase enrollment of participants must be tailored to the target populations and their specific concerns.

All efforts to classify participants by race and/or ethnicity create large, heterogeneous, and imprecise groups. Racial/ethnic categorizations are poor proxies of environmental factors, and there are myriad socio-cultural and environmental differences within racial/ethnic groups that impact recruitment and participation[20]. One way to better address these diverse concerns is through community-based participatory research (CBPR). CBPR engages community stakeholders as peers in all stages of the study from design to dissemination of results. For example, hiring research specialists from a community to translate study materials increases the chances of using appropriate, non-stigmatizing language[42]. CBPR, though underutilized, has successfully led to increased participation of non-White populations in genetic research and could be a useful approach for increasing recruitment across many diverse populations[43–45]. Efforts for recruiting historically underserved participants into AD genetics studies using community-based approaches are underway. The Asian Cohort for Alzheimer’s Disease (ACAD), for example, is currently recruiting Asian American and Asian Canadian participants using CBPR approaches including partnership with clinics and senior homes that serve Asian communities and translation of materials into Mandarin, Cantonese, Vietnamese, and Korean.

ACAD and other efforts to recruit Asian Americans are critically needed. To increase knowledge of AD genetics in Asians, the ADSP is primarily relying on partnerships with foreign-based studies in India and Korea[46–48]. While this strategy may help overcome potential cultural barriers to obtaining genetic data that represents individuals with genetic ancestry similar to those currently residing in India and Korea, there are limitations to interpretation and generalizability of findings conducted in these studies. First, individuals with Indian or Korean

ancestry make up only a quarter of those who identify as Asian American[49]. Perhaps, more importantly to AD genetic research, GWAS associations are influenced by context[50], and there are vast differences in environmental factors across countries that could modify genetic effects. There may also be barriers in future efforts to include social determinants and electronic health record phenotypes across AD sequencing participants that could lead to further exacerbating the disparity in AD knowledge for Asians in the US.

Advances in statistical methods offer additional tools for increasing genetic discoveries in diverse populations. Specific ancestry backgrounds enable alternative or complementary gene-discovery approaches. For example, admixture mapping, which leverages the mix of pre-diaspora ancestry in contemporary populations, can have more statistical power than GWAS to discover genomic regions associated with traits or diseases[51]. Admixture mapping has already implicated novel AD loci in studies of African Americans with HGDP-African/European-like ancestry and Caribbean Hispanics with 1KG-CEU/YRI-like and HGDP-Pima/Maya/Colombia ancestry despite samples sizes that are relatively small (~ 10,000)[52–54]. Methods have been developed that allow meta-analysis of GWAS across ancestries or inclusion of participants with diverse ancestries in the same GWAS[55]. These ‘transethnic’ GWAS may be a superior alternative to stratified studies because of the boost in statistical power from variants that are found in many populations and the reinforcement of the fact that global genetic ancestry cannot be accurately stratified into biologically meaningful ‘racial’ groups. Indeed, the association of the *SHARPIN* genetic region with AD risk was observed in a recent multi-ancestry GWAS meta-analysis including 56,241 individuals; previously this region had only been detected in an AD GWAS of exclusively non-Hispanic White individuals with sample size 13X greater[17,56].

While statistical advancements offer significant benefits, they do not alleviate the burden of recruiting larger sample sizes of diverse participants.

Our study has several limitations. Population-based prevalence estimates may be biased – likely underestimating the burden on Indigenous individuals and those racialized as Black who are less likely to be formally diagnosed with AD due to inequitable treatment in healthcare settings and diagnostic thresholds primarily based on White individuals[57–59]. Furthermore, the CDC WONDER dataset provides statistics for a single underlying cause of death for each person. It’s likely that many people were not identified as dying with AD if another condition was a more immediate cause of their death (ex., someone living with AD who died from heart failure, in which case AD would not be listed as the underlying cause of death on the death certificate). Disparities driven by structural racism exist for chronic diseases including cardiovascular disease and cancer[60,61], which can lead to biased underreporting of AD mortality as these causes of death may be disproportionately masking AD-related deaths. Differences in survival rates after dementia diagnosis could also contribute to the differences in proportions between prevalence and mortality rates[62]. There may also be some unreliable reporting of racial/ethnic classification of mortality data due to reporting by an observer rather than self-report. Our summaries of existing AD genetic datasets do not account for the possibility of unaccounted sample overlap within array datasets or within sequencing datasets, which would cause us to have over-estimated existing sample sizes. In this case, the problems we described would only be more important to address. Lastly, power calculations did not specifically model rare variant aggregate tests. Despite these constraints, we describe the ‘best case’ of genetic data representation, which indicated that while participation in AD genetics datasets is approximately

proportional to AD burden, studies remain underpowered to elucidate the genetic architecture of AD in diverse populations.

In conclusion, we must recognize that non-White populations are simultaneously overexposed to AD risk and underrepresented in AD genetics research. Substantial effort must continue to be made to build trust, foster engagement, and actively involve historically underrepresented groups in AD genetics research to ensure that research outcomes and resulting therapies are effective for individuals of all backgrounds.

## References

- [1] Sosa-Ortiz AL, Acosta-Castillo I, Prince MJ. Epidemiology of dementias and Alzheimer's disease. *Arch Med Res* 2012;43:600–8. <https://doi.org/10.1097/00019442-199821001-00002>.
- [2] Alzheimer's Association. 2023 Alzheimer's disease facts and figures. *Alzheimer's & Dementia* 2023;19:1598–695. <https://doi.org/10.1002/alz.13016>.
- [3] Krishnamurthy S, Rollin FG. We must be clear that the root cause of racial disparities in Alzheimer's disease is racism. *Alzheimer's & Dementia* 2023;n/a. <https://doi.org/https://doi.org/10.1002/alz.13389>.
- [4] Adkins-Jackson PB, George KM, Besser LM, Hyun J, Lamar M, Hill-Jarrett TG, et al. The structural and social determinants of Alzheimer's disease related dementias. *Alzheimer's & Dementia* 2023;19:3171–85. <https://doi.org/https://doi.org/10.1002/alz.13027>.
- [5] Gatz M, Reynolds CA, Fratiglioni L, Johansson B, Mortimer JA, Berg S, et al. Role of genes and environments for explaining Alzheimer disease. *Arch Gen Psychiatry* 2006;63:168–74. <https://doi.org/10.1001/archpsyc.63.2.168>.
- [6] Andrews SJ, Renton AE, Fulton-Howard B, Podlesny-Drabiniok A, Marcora E, Goate AM. The complex genetic architecture of Alzheimer's disease: novel insights and future directions. *EBioMedicine* 2023;90. <https://doi.org/10.1016/j.ebiom.2023.104511>.
- [7] Kamboh MI. Genomics and functional genomics of Alzheimer's disease. *Neurotherapeutics* 2022;19:152–72. <https://doi.org/10.1007/s13311-021-01152-0>.
- [8] Blue EE, Horimoto ARVR, Mukherjee S, Wijsman EM, Thornton TA. Local ancestry at APOE modifies Alzheimer's disease risk in Caribbean Hispanics. *Alzheimers Dement* 2019;15:1524–32. <https://doi.org/10.1016/j.jalz.2019.07.016>.
- [9] Reitz C, Mayeux R. Genetics of Alzheimer's disease in Caribbean Hispanic and African American populations. *Biol Psychiatry* 2014;75:534–41. <https://doi.org/10.1016/J.BIOPSYCH.2013.06.003>.
- [10] Khan AT, Gogarten SM, McHugh CP, Stilp AM, Sofer T, Bowers ML, et al. Recommendations on the use and reporting of race, ethnicity, and ancestry in genetic research: Experiences from the NHLBI TOPMed program. *Cell Genomics* 2022;2:100155. <https://doi.org/10.1016/j.xgen.2022.100155>.
- [11] Tishkoff SA, Verrelli BC. Patterns of Human Genetic Diversity: Implications for Human Evolutionary History and Disease. *Annu Rev Genomics Hum Genet* 2003;4:293–340. <https://doi.org/10.1146/annurev.genom.4.070802.110226>.

- [12] Novembre J, Stephens M. Interpreting principal component analyses of spatial population genetic variation. *Nat Genet* 2008;40:646–9. <https://doi.org/10.1038/ng.139>.
- [13] Marchani EE, Watkins WS, Bulayeva K, Harpending HC, Jorde LB. Culture creates genetic structure in the Caucasus: Autosomal, mitochondrial, and Y-chromosomal variation in Daghestan. *BMC Genet* 2008;9:47. <https://doi.org/10.1186/1471-2156-9-47>.
- [14] Wightman DP, Jansen IE, Savage JE, Shadrin AA, Bahrami S, Holland D, et al. A genome-wide association study with 1,126,563 individuals identifies new risk loci for Alzheimer’s disease. *Nat Genet* 2021;53:1276–82. <https://doi.org/10.1038/s41588-021-00921-z>.
- [15] Kunkle BW, Schmidt M, Klein H-U, Naj AC, Hamilton-Nelson KL, Larson EB, et al. Novel Alzheimer Disease Risk Loci and Pathways in African American Individuals Using the African Genome Resources Panel: A Meta-analysis. *JAMA Neurol* 2021;78:102–13. <https://doi.org/10.1001/jamaneurol.2020.3536>.
- [16] Kunkle BW, Grenier-Boley B, Sims R, Bis JC, Damotte V, Naj AC, et al. Genetic meta-analysis of diagnosed Alzheimer’s disease identifies new risk loci and implicates A $\beta$ , tau, immunity and lipid processing. *Nat Genet* 2019;51:414–30. <https://doi.org/10.1038/s41588-019-0358-2>.
- [17] Bellenguez C, Küçükali F, Jansen IE, Kleindam L, Moreno-Grau S, Amin N, et al. New insights into the genetic etiology of Alzheimer’s disease and related dementias. *Nat Genet* 2022;54:412–36. <https://doi.org/10.1038/s41588-022-01024-z>.
- [18] National Academies of Sciences and Medicine E. Using Population Descriptors in Genetics and Genomics Research: A New Framework for an Evolving Field. Washington, DC: The National Academies Press; 2023. <https://doi.org/10.17226/26902>.
- [19] Office Of Management and Budget (OMB) Standards. National Institutes of Health n.d. <https://orwh.od.nih.gov/toolkit/other-relevant-federal-policies/OMB-standards>.
- [20] McConkie-Rosell A, Spillmann RC, Schoch K, Sullivan JA, Walley N, McDonald M, et al. Unraveling non-participation in genomic research: A complex interplay of barriers, facilitators, and sociocultural factors. *J Genet Couns* 2023. <https://doi.org/10.1002/jgc4.1707>.
- [21] Mena PR, Kunkle BW, Faber KM, Adams LD, Inciute JD, Whitehead PL, et al. The Alzheimer’s Disease Sequencing Project Follow Up Study (ADSP-FUS): increasing ethnic diversity in Alzheimer’s disease (AD) genetics research. *Alzheimer’s & Dementia* 2022;18:e068083. <https://doi.org/10.1002/alz.068083>.
- [22] Rajan KB, Weuve J, Barnes LL, McAninch EA, Wilson RS, Evans DA. Population estimate of people with clinical Alzheimer’s disease and mild cognitive impairment in the

- United States (2020–2060). *Alzheimer's & Dementia* 2021.  
<https://doi.org/10.1002/alz.12362>.
- [23] Matthews KA, Xu W, Gaglioti AH, Holt JB, Croft JB, Mack D, et al. Racial and ethnic estimates of Alzheimer's disease and related dementias in the United States (2015–2060) in adults aged  $\geq 65$  years. *Alzheimer's & Dementia* 2019;15:17–24.  
<https://doi.org/10.1016/j.jalz.2018.06.3063>.
- [24] Centers for Disease Control and Prevention NC for HStatistics. National Vital Statistics System, Mortality 1999-2020 on CDC WONDER Online Database, released in 2021. Data are from the Multiple Cause of Death Files, 1999-2020, as compiled from data provided by the 57 vital statistics jurisdictions through the Vital Statistics Cooperative Program. n.d. <http://wonder.cdc.gov/ucd-icd10.html>.
- [25] Jansen IE, Savage JE, Watanabe K, Bryois J, Williams DM, Steinberg S, et al. Genome-wide meta-analysis identifies new loci and functional pathways influencing Alzheimer's disease risk. *Nat Genet* 2019;51:404–13. <https://doi.org/10.1038/s41588-018-0311-9>.
- [26] Harold D, Abraham R, Hollingworth P, Sims R, Gerrish A, Hamshere ML, et al. Genome-wide association study identifies variants at *CLU* and *PICALM* associated with Alzheimer's disease. *Nat Genet* 2009;41:1088–93. <https://doi.org/10.1038/ng.440>.
- [27] Schwartzenuber J, Cooper S, Liu JZ, Barrio-Hernandez I, Bello E, Kumasaka N, et al. Genome-wide meta-analysis, fine-mapping and integrative prioritization implicate new Alzheimer's disease risk genes. *Nat Genet* 2021;53:392–402.  
<https://doi.org/10.1038/s41588-020-00776-w>.
- [28] Jonsson T, Stefansson H, Steinberg S, Jonsdottir I, Jonsson P V, Snaedal J, et al. Variant of *TREM2* Associated with the Risk of Alzheimer's Disease. *New England Journal of Medicine* 2012;368:107–16. <https://doi.org/10.1056/NEJMoa1211103>.
- [29] Guerreiro R, Wojtas A, Bras J, Carrasquillo M, Rogaeva E, Majounie E, et al. *TREM2* Variants in Alzheimer's Disease. *New England Journal of Medicine* 2012;368:117–27.  
<https://doi.org/10.1056/NEJMoa1211851>.
- [30] Chen H, Wu G, Jiang Y, Feng R, Liao M, Zhang L, et al. Analyzing 54,936 samples supports the association between *CD2AP* rs9349407 polymorphism and Alzheimer's disease susceptibility. *Mol Neurobiol* 2015;52:1–7. <https://doi.org/10.1007/s12035-014-8834-2>.
- [31] Miyashita A, Koike A, Jun G, Wang L-S, Takahashi S, Matsubara E, et al. *SORL1* is genetically associated with late-onset Alzheimer's disease in Japanese, Koreans and Caucasians. *PLoS One* 2013;8:e58618. <https://doi.org/10.1371/journal.pone.0058618>.

- [32] Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti A, et al. A global reference for human genetic variation. *Nature* 2015;526:68–74. <https://doi.org/10.1038/nature15393>.
- [33] Tosto G, Fu H, Vardarajan BN, Lee JH, Cheng R, Reyes-Dumeyer D, et al. F-box/LRR-repeat protein 7 is genetically associated with Alzheimer’s disease. *Ann Clin Transl Neurol* 2015;2:810–20. <https://doi.org/10.1002/acn3.223>.
- [34] Cukier HN, Kunkle BW, Vardarajan BN, Rolati S, Hamilton-Nelson KL, Kohli MA, et al. ABCA7 frameshift deletion associated with Alzheimer disease in African Americans. *Neurol Genet* 2016;2:e79. <https://doi.org/10.1212/NXG.0000000000000079>.
- [35] Privé F, Aschard H, Carmi S, Folkersen L, Hoggart C, O’Reilly PF, et al. Portability of 245 polygenic scores when derived from the UK Biobank and applied to 9 ancestry groups from the same cohort. *The American Journal of Human Genetics* 2022;109:12–23. <https://doi.org/https://doi.org/10.1016/j.ajhg.2021.11.008>.
- [36] Martin AR, Kanai M, Kamatani Y, Okada Y, Neale BM, Daly MJ. Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat Genet* 2019;51:584–91. <https://doi.org/10.1038/s41588-019-0379-x>.
- [37] Vernon LF. Tuskegee syphilis study not America’s only medical scandal: Chester M. Southam, MD, Henrietta Lacks, and the Sloan-Kettering research scandal. *Journal of Health Ethics* 2020;16. <https://doi.org/10.18785/ojhe.1602.03>.
- [38] Jones BL, Vyhldal CA, Bradley-Ewing A, Sherman A, Goggin K. If We Would Only Ask: How Henrietta Lacks Continues to Teach Us About Perceptions of Research and Genetic Research Among African Americans Today. *J Racial Ethn Health Disparities* 2017;4:735–45. <https://doi.org/10.1007/s40615-016-0277-1>.
- [39] Claw KG, Anderson MZ, Begay RL, Tsosie KS, Fox K, Garrison NA, et al. A framework for enhancing ethical genomic research with Indigenous communities. *Nat Commun* 2018;9:2957. <https://doi.org/10.1038/s41467-018-05188-3>.
- [40] Ewing A, Thompson N, Ricks-Santi L. Strategies for Enrollment of African Americans into Cancer Genetic Studies. *Journal of Cancer Education* 2015;30:108–15. <https://doi.org/10.1007/s13187-014-0669-z>.
- [41] Dilworth-Anderson P, Hendrie HC, Manly JJ, Khachaturian AS, Fazio S. Diagnosis and assessment of Alzheimer’s disease in diverse populations. *Alzheimer’s & Dementia* 2008;4:305–9. <https://doi.org/10.1016/j.jalz.2008.03.001>.
- [42] Gonzalez S, Strizich G, Isasi CR, Hua S, Comas B, Sofer T, et al. Consent for use of genetic data among US Hispanics/Latinos: results from the Hispanic Community Health Study/Study of Latinos. *Ethn Dis* 2021;31:547. <https://doi.org/10.18865/ed.31.4.547>.

- [43] Brown KE, Fohner AE, Woodahl EL. Beyond the Individual: Community-Centric Approaches to Increase Diversity in Biomedical Research. *Clin Pharmacol Ther* 2023;113:509–17. <https://doi.org/10.1002/cpt.2808>.
- [44] Boyer BB, Mohatt G V, Pasker RL, Drew EM, McGlone KK. Sharing results from complex disease genetics studies: a community based participatory research approach. *Int J Circumpolar Health* 2007;66:19–30. <https://doi.org/10.3402/ijch.v66i1.18221>.
- [45] Ochs-Balcom HM, Jandorf L, Wang Y, Johnson D, Meadows Ray V, Willis MJ, et al. “It takes a village”: multilevel approaches to recruit African Americans and their families for genetic research. *J Community Genet* 2015;6:39–45. <https://doi.org/10.1007/s12687-014-0199-8>.
- [46] Kang S, Gim J, Lee J, Gunasekaran TI, Choi KY, Lee JJ, et al. Potential novel genes for late-onset Alzheimer’s disease in East-Asian descent identified by APOE-stratified genome-wide association study. *Journal of Alzheimer’s Disease* 2021;82:1451–60. <https://doi.org/10.3233/JAD-210145>.
- [47] Yi D, Byun MS, Risacher SL, Craft H, Crane PK, Trittschuh EH, et al. The Korean brain aging study for the early diagnosis and prediction of Alzheimer’s disease (KBASE): Cognitive data harmonization. *Alzheimer’s & Dementia* 2023;19:e064533. <https://doi.org/10.4306/pi.2017.14.6.851>.
- [48] Lee J, Dey AB. Introduction to LASI-DAD: The longitudinal aging study in India-diagnostic assessment of dementia. *J Am Geriatr Soc* 2020;68:S3. <https://doi.org/10.1111/jgs.16740>.
- [49] Ruiz NG, Noe-Bustamante L, Shah S. Appendix: Demographic profile of Asian American adults. Pew Research Center 2023. <https://www.pewresearch.org/race-and-ethnicity/2023/05/08/asian-american-identity-appendix-demographic-profile-of-asian-american-adults/> (accessed November 26, 2024).
- [50] Riehm KE, Keyes KM, Susser ES. Social determinants of health and selection bias in genome-wide association studies. *World Psychiatry* 2023;22:160. <https://doi.org/10.1002/wps.21047>.
- [51] Shriner D. Overview of admixture mapping. *Curr Protoc Hum Genet* 2013;Chapter 1. <https://doi.org/10.1002/0471142905.hg0123s76>.
- [52] Horimoto ARVR, Xue D, Thornton TA, Blue EE. Admixture mapping reveals the association between Native American ancestry at 3q13.11 and reduced risk of Alzheimer’s disease in Caribbean Hispanics. *Alzheimers Res Ther* 2021;13:122. <https://doi.org/10.1186/s13195-021-00866-9>.

- [53] Horimoto ARVR, Boyken LA, Blue EE, Grinde KE, Nafikov RA, Sohi HK, et al. Admixture mapping implicates 13q33. 3 as ancestry-of-origin locus for Alzheimer disease in Hispanic and Latino populations. *Human Genetics and Genomics Advances* 2023;4. <https://doi.org/10.1016/j.xhgg.2023.100207>.
- [54] Kizil C, Sariya S, Kim YA, Rajabli F, Martin E, Reyes-Dumeyer D, et al. Admixture Mapping of Alzheimer’s disease in Caribbean Hispanics identifies a new locus on 22q13.1. *Mol Psychiatry* 2022;27:2813–20. <https://doi.org/10.1038/s41380-022-01526-6>.
- [55] Atkinson EG, Maihofer AX, Kanai M, Martin AR, Karczewski KJ, Santoro ML, et al. Tractor uses local ancestry to enable the inclusion of admixed individuals in GWAS and to boost power. *Nat Genet* 2021;53:195–204. <https://doi.org/10.1038/s41588-020-00766-y>.
- [56] Rajabli F, Benchek P, Tosto G, Kushch N, Sha J, Bazemore K, et al. Multi-ancestry genome-wide meta-analysis of 56,241 individuals identifies LRRC4C, LHX5-AS1 and nominates ancestry-specific loci PTPRK, GRB14, and KIAA0825 as novel risk loci for Alzheimer disease: the Alzheimer Disease Genetics Consortium. *MedRxiv* 2023:2023–7. <https://doi.org/10.1101/2023.07.06.23292311>.
- [57] Clark PC, Kutner NG, Goldstein FC, Peterson-Hazen S, Garner V, Zhang R, et al. Impediments to timely diagnosis of Alzheimer’s disease in African Americans. *J Am Geriatr Soc* 2005;53:2012–7. <https://doi.org/10.1111/j.1532-5415.2005.53569.x>.
- [58] Griffin-Pierce T, Silverberg N, Connor D, Jim M, Peters J, Kaszniak A, et al. Challenges to the recognition and assessment of Alzheimer’s disease in American Indians of the southwestern United States. *Alzheimer’s & Dementia* 2008;4:291–9. <https://doi.org/10.1016/j.jalz.2007.10.012>.
- [59] Barnes LL. Alzheimer disease in African American individuals: increased incidence or not enough data? *Nat Rev Neurol* 2022;18:56–62. <https://doi.org/10.1038/s41582-021-00589-3>.
- [60] Ski CF, King-Shier KM, Thompson DR. Gender, socioeconomic and ethnic/racial disparities in cardiovascular disease: A time for change. *Int J Cardiol* 2014;170:255–7. <https://doi.org/10.1016/j.ijcard.2013.10.082>.
- [61] O’Keefe EB, Meltzer JP, Bethea TN. Health disparities and cancer: racial disparities in cancer mortality in the United States, 2000–2010. *Front Public Health* 2015;3:51. <https://doi.org/10.3389/fpubh.2015.00051>.
- [62] Mayeda ER, Glymour MM, Quesenberry CP, Johnson JK, Pérez-Stable EJ, Whitmer RA. Survival after dementia diagnosis in five racial/ethnic groups. *Alzheimer’s & Dementia* 2017;13:761–9. <https://doi.org/10.1016/j.jalz.2016.12.008>.

## Tables

**Table 1. Population-specific AD burden and representation in AD genetics data in the United States.**

<b>Race/ Ethnicity</b>	<b>Projected AD Prevalence 2020 in 1000s (%)</b>	<b>Projected Dementia Prevalence 2020 in 1000s (%)</b>	<b>Deaths 2016-2020 (%)</b>	<b>US population (2000 standard)</b>	<b>Array Data Sample Size (%)</b>	<b>WGS Data Sample Size (%)</b>
<b>AI/AN</b>	--	38 (0.6)	1,865 (0.3)	1,584,958	14 (0.02)	152 (0.42)
<b>Asian</b>	--	212 (3.7)	14,272 (2.4)	12,526,017	1170 (1.78)	2820 (7.76)
<b>Black</b>	1060 (17.4)	726 (12.6)	45,946 (7.6)	24,282,298	9439 (14.36)	5695 (15.67)
<b>Hispanic</b>	710 (11.7)	594 (10.3)	38,960 (6.4)	20,361,950	4491 (6.83)	11329 (31.18)
<b>White</b>	4300 (70.8)	4186 (72.7)	505,889 (83.4)	201,746,665	50619 (77.01)	16340 (44.97)

*AI/AN: American Indian/Alaska Native*

**Table 2. Statistical comparison of observed race/ethnicity representation in AD genetics studies versus expected values based on AD and dementia prevalence and mortality.**

Table 2 Legend. Chi-squared tests for proportions were used to compare disease burden in the population (AD prevalence, Dementia prevalence, AD mortality) and racial/ethnic representation in AD genetics datasets.

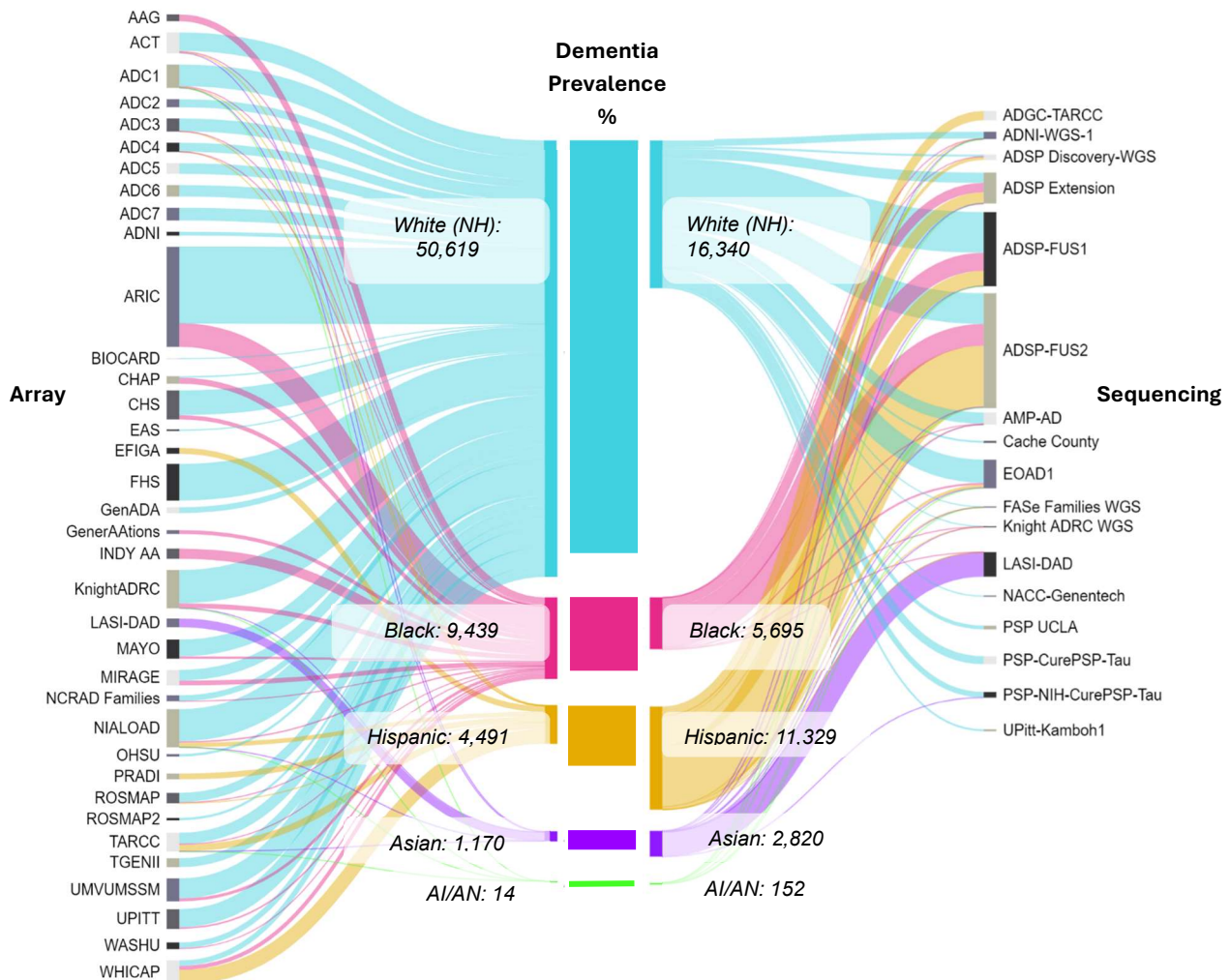
\*Bonferroni adjusted  $\alpha = 0.00833$

	<b>Array</b>	<b>Sequencing</b>
<b>AD prevalence</b>	$X^2 = 3.63, p = 0.458$	$X^2 = 22.88, p = 0.0001^*$
<b>Dementia prevalence</b>	$X^2 = 2.16, p = 0.707$	$X^2 = 18.849, p = 0.0008^*$
<b>AD mortality</b>	$X^2 = 2.69, p = 0.476$	$X^2 = 33.49, p = 9.5e-07^*$

## Figures

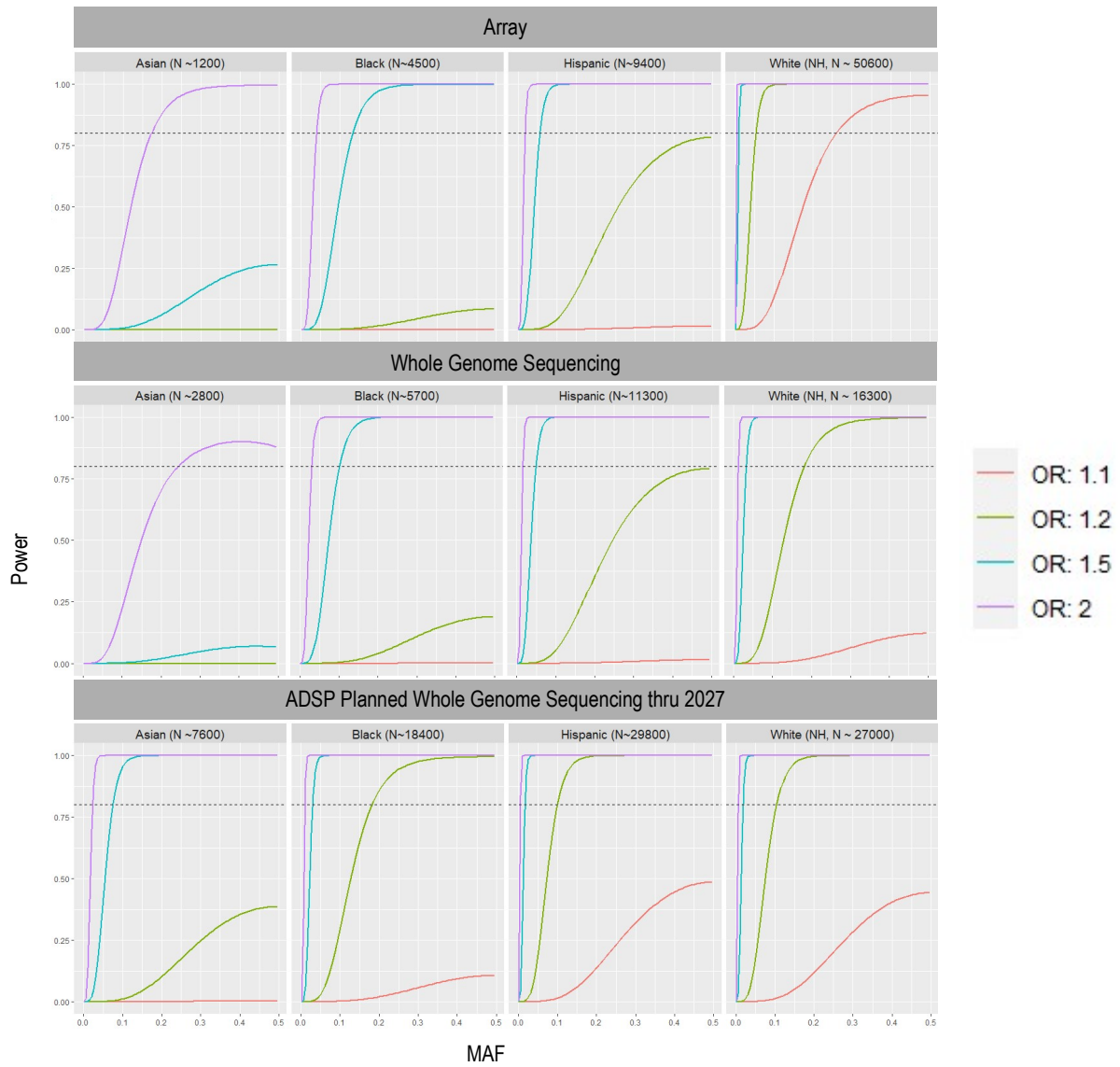
### Figure 1. Racial/ethnic profile of AD genetic data.

Figure 1 Legend. The left-hand side depicts a Sankey plot showing racial/ethnic representation in each array dataset, flowing from left to right. The right-hand side depicts a Sankey plot of racial/ethnic representation in the whole genome sequencing (WGS) data, flowing right to left. On the outside edges are the individual cohorts. The participants are grouped by five broad racial/ethnic categories: White (non-Hispanic), Black, Hispanic, Asian, and American Indian/Alaska Native (AI/AN). In the center of the figure is the relative race/ethnicity specific burden of dementia. Exact proportions of dementia burden are found in [Table 1](#).



**Figure 2. Statistical power to detect loci representing AD genetic architecture across populations.**

Figure 2 Legend. Each panel shows power to detect a significantly associated single nucleotide polymorphism (SNP) ( $p = 5e-08$ ) for a different racial/ethnic group using genome-wide association studies (GWAS). Power simulations are based on current or projected sample sizes and case proportions. Power was simulated based on a set of effect sizes (odds ratio = 1.1, 1.2, 1.5, 2) and minor allele frequencies (MAF) ranging from 0.001 to 0.5. Separate simulations were conducted for array, current whole genome sequencing (WGS), and projected WGS data. The dashed line represents power = 0.80. American Indian/Alaska Native (AI/AN) are not included because current sample sizes are too small to detect any SNPs regardless of MAF or effect size.



## Paper 2: Polygenic risk scores for incident dementia in the Multi-Ethnic Study of Atherosclerosis

### Abstract

**Background:** Over 75 Alzheimer's disease (AD) and dementia-associated variants have been identified through genome-wide association studies (GWAS). Like many other traits, genomic studies of AD have been biased towards those with European ancestry, and the utility of polygenic risk scores (PRS) for predicting AD and dementia in diverse and admixed populations remains unclear. We compared how PRS approaches differing in  $p$ -value thresholds, variant weights, and source ancestry perform in predicting late-onset dementia in a multi-ethnic cohort.

**Methods:** We compared clumping and thresholding (C+T) methods with varying parameters against Bayesian approaches (PRS-CS, PRS-CSx) in the Multi-Ethnic Study of Atherosclerosis (MESA), a longitudinal cohort of 6,814 men and women aged 45-84 at baseline who identified as Black/African American, Chinese, Hispanic, or White. We compared the ability of each method to predict incident dementia for MESA participants overall and stratified by self-reported race/ethnicity. We additionally analyzed performance across groups stratified by the estimated proportion of non-Finnish European (NFE) ancestry.

**Results:** The median follow-up time was 16.8 years, and 569 (8.8%) participants developed dementia according to hospital and death certificate ICD codes. The C+T method with  $p < 5e-08$  had the strongest predictive performance overall ( $C = 0.55$ ,  $SE = 0.01$ ) and among the Black ( $C = 0.57$ ,  $SE = 0.02$ ) and White ( $C = 0.56$ ,  $SE = 0.02$ ) participants. PRS performance did not vary significantly with NFE proportion. Comparisons across models revealed that the inclusion of more SNPs in either C+T models using Bayesian approaches does not improve predictive accuracy.

**Conclusion:** The PRS based on C+T method with only 15 SNPs with strong evidence of association with AD is more strongly associated with dementia than PRS derived from Bayesian models that include >800,000 SNPs, even in target populations genetically dissimilar from that of the source data. When it comes to PRS for dementia, it pays to be picky: including more variants does not improve PRS performance.

## **Abbreviations**

age-at-onset (AAO),

African (AMR, referring to 1000 Genomes superpopulation)

African Genome Resources Panel (AGRP),

American (AMR, referring to 1000 Genomes superpopulation)

Alzheimer's disease (AD),

clumping and thresholding (C+T)

genome-wide association study (GWAS),

East Asian (EAS, referring to 1000 Genomes superpopulation)

hazard ratio (HR)

International Classification of Diseases (ICD)

International Genomics of Alzheimer's Disease Project (IGAP)

likelihood ratio test (LRT)

linkage disequilibrium (LD)

Multi-Ethnic Study of Atherosclerosis (MESA),

National Institute on Aging Genetics of Alzheimer's Disease Data Storage Site (NIAGADS),

Non-Finnish European (NFE, referring to 1000 Genomes superpopulation)

odds ratio (OR)

polygenic risk score (PRS),

single nucleotide polymorphism (SNP)

## Introduction

Dementia is a growing global health challenge, projected to affect over 150 million people worldwide by 2050[1]. Populations are aging around the world, and with nearly one third of adults over 65 dying with Alzheimer's disease (AD) or other dementias, there is an urgent need for effective predictive tools that can aid in risk stratification and lead to more precise treatment and prevention[2]. AD is the most common cause of dementia and is strongly influenced by genetic variation. Less than one percent of AD cases are early-onset and caused by the Mendelian inheritance of variants in *APP*, *PSEN1*, or *PSEN2*[3]. The remaining cases have far more complex etiology. The apolipoprotein E (*APOE*) *e4* allele is the strongest genetic risk factor for late-onset AD, but less than half of AD patients carry an *e4* allele[4,5]. Aside from *APOE*, dozens of loci have been found to be significantly associated with AD through genome-wide association studies (GWAS, [6]). While these loci typically have small effect sizes, they are more common in the population, and their joint effects can place individuals at elevated genetic risk for disease.

Polygenic risk scores (PRS) based on the effects of common genetic variants have been shown to be predictive of disease[7–9]. However, there is no clear consensus on the optimal approach for constructing PRS for AD or dementia, particularly in diverse and admixed populations. The most common approach for constructing PRS is clumping and thresholding (C+T), which involves initially including all SNPs tested in a GWAS and then filtering them based on *p*-value threshold and linkage disequilibrium[10]. Some previous studies have found that less stringent *p*-value thresholds, which allow for the inclusion of a larger number of SNPs in the PRS, lead to better predictive performance of AD. Escott-Price *et al.* (2015) reported that the PRS with *p*-value threshold  $<0.5$  was most strongly associated with AD[11]. Another study suggests that a threshold of  $p < 0.10$  had optimal performance [12]. In contrast, Zhang and

colleagues found that restricting PRS to SNPs that are significantly or suggestively associated with AD in GWAS had better performance, implying PRS constructed using fewer than 100 SNPs can achieve superior prediction[13]. Notably, the comparisons discussed thus far are limited to populations with European ancestry.

The performance of PRS for AD in groups with diverse ancestry remains underexplored. Studies of other conditions have shown that PRS performance deteriorates as the genetic distance between the target and GWAS training populations increases[14–16]. Most of the risk loci discovered to be significantly associated with AD have been found in large GWAS studies of self-reported non-Hispanic white individuals who cluster with 1000 Genome (1KG) European references[17,18]. Only a fraction of these loci have been replicated in populations with different genetic ancestral backgrounds including *APOE*, *ABCA7*, *TREM2*, *SORL1*, and *CLU*[19]. GWAS of non-European ancestry populations remain underpowered to discover risk loci with low to moderate effects, which comprise most of the risk loci identified thus in the large European ancestry GWAS[20]. To determine PRS models for AD that are accurate and actionable, it is crucial to understand how various PRS models perform in diverse groups.

In this study, we assess the performance of various PRS methodologies in predicting late-onset dementia. The Multi-Ethnic Study of Atherosclerosis (MESA), a longitudinal cohort study, includes participants who self-identify as Black/African American, Chinese, Hispanic, or White, providing an opportunity to examine PRS performance in a diverse population. Based on GWAS of clinically ascertained AD, we compared the performance of traditional C+T methods at a range of  $p$ -value thresholds against Bayesian approaches (PRS-CS, PRS-CSx) using multiple GWAS summary statistics with differing ancestral backgrounds.

## **Methods**

### *Study Population*

### *Study Population*

MESA has been previously described[21]. Briefly, MESA is a prospective cohort study originally designed to study cardiovascular disease. Between 2000 and 2002, MESA recruited 6,814 Black/African American, Chinese, Hispanic/Latino, and White participants aged 45-84 from six sites in the United States: Baltimore, Maryland; Chicago, Illinois; Forsyth County, North Carolina; Los Angeles County, California; Northern Manhattan and the Bronx, New York; and Saint Paul, Minnesota. All participants were free from clinical cardiovascular disease and dementia at baseline. All participants provided written informed consent at baseline and all exams following. Institutional Review Board approval was received from each of the six sites. MESA participants with imputed genotypes and information on dementia status were included.

### *Inferring Global Ancestry Proportions*

We estimated global ancestry proportions for all genotyped MESA participants. Global ancestry proportions are based on local ancestry estimates from RFMix2 using HGDP+1KG samples accessed through gnomAD v3.1 as references[22–24] Samples were randomly selected from the following superpopulation groups to construct balanced sample maps: American (AMR), African (AFR), East Asian (EAS), and Non-Finnish European (NFE[24]). The genetic map with coordinates from the latest build of the human reference genome GRCh38 was downloaded from the Eagle v2.4.1 package (<http://data.broadinstitute.org/alkesgroup/Eagle/downloads/>). Using the global ancestry proportions, participants were assigned to low NFE, intermediate NFE, and high NFE groups

where low NFE was considered less than 33% NFE ancestry, intermediate NFE captured those who had between 33% and 67% NFE. High NFE included those with greater than 67% NFE.

### *Incident Dementia Outcome*

Participants were followed up with via telephone interview every 9 to 12 months for updates on hospital admissions or deaths. Incident dementia was identified based on a set of ICD codes at either hospitalization or death. The candidate dementia cases were identified using the following diagnosis codes: ICD-9: 290, 294, 331.0, 331.1, 331.2, 331.82, 331.83, 331.9, 438.0, and 780.93; ICD-10: F00, F01, F03, F04, G30, G31 (excluding G31.2), I69.91, and R41. The ICD code-based identification has been validated against medical record text that indicates significant decline in cognitive function compared with a previous level[25]

### *Genotyping*

SNPs for all MESA participants were genotyped using the Affymetrix 6.0 SNP array. SNPs were imputed using IMPUTE version 2.2.2 and 1000 Genomes cosmopolitan phase 3 version 5 reference haplotypes. KING was used to infer relatedness and an unrelated subset of individuals was selected by randomly choosing one individual from each first-degree related group[26].

### *Calculating Polygenic Risk Scores*

We compared six polygenic risk score models, all of which excluded the *APOE* region (GRCh38 chr19: 44408822 - 45408822). The C+T methods were used to estimate PRS for each target sample using the following *p*-value thresholds: 0.01, 1e-05, and 5e-08. SNPs were filtered based on LD  $<r^2 = 0.01$ . C+T PRS were calculated using PLINK v1.90[27]. After filtering based on *p*-value and LD, PRS were calculated based on the dosage of the SNP effect allele multiplied by the effect sizes. The SNPs and effect sizes for the C+T models were derived from the 2019

International Genomics of Alzheimer's Project genetic meta-analysis of clinically diagnosed late-onset Alzheimer's disease among those of European descent, which includes 21,982 cases and 41,944 controls across 46 studies[17]. All summary statistics were obtained from the National Institute on Aging Genetics of Alzheimer's Disease Data Storage Site (NIAGADS).

In addition to the C+T models with varying  $p$ -value stringency, two Bayesian models were compared: PRS-CS and PRS-CSx[28,29]. Both PRS-CS and PRS-CSx use a continuous shrinkage model that accounts for linkage disequilibrium by tuning or shrinking the effect sizes. We did not use a separate validation set to tune parameters, instead using the -auto option for both PRS-CS and PRS-CSx. Two GWAS summary statistics were used for the PRS-CS models: the European ancestry IGAP study and a cross-population GWAS of 15,579 cases and 17,690 controls that included self-reported Whites, African Americans, Japanese, and Israeli-Arabs (NG00056, [30]).

PRS-CSx allows for multiple summary statistics with differing ancestral backgrounds to be used concurrently. In addition to the European ancestry IGAP study, we also included summary statistics from the African Genome Resources Panel GWAS of 2,748 cases and 5,222 controls in the same model (NG00100, [31]).

To adjust the polygenic risk scores for population structure, we used a previously described procedure to calculate residualized scores based on the principal components[32]. We fit the raw PRS as a function of the first three principal components in non-affected individuals. We used the linear model to calculate a predicted PRS for all individuals. We then computed the residualized, population-structure adjusted PRS by calculating the difference between the raw and predicted PRS. The residualized score was then standardized based on the mean and standard deviation.

## *Statistical Analysis*

Cox proportional hazards models were used to examine the association between the PRS scores and incident dementia in all MESA participants, groups stratified by self-reported race/ethnicity, and groups stratified by quantile of NFE ancestry. Univariate models were computed separately for each PRS method.

Harrell's concordance (C-Index) was used to compare predictive performance. Comparisons were conducted in all participants and groups stratified by self-reported race/ethnicity. Additional comparisons were made across groups in different tertiles of European ancestry to examine if model performance was biased for those with greater proportion of European ancestry.

We computed the added predictive value of each PRS method compared to baseline models that included sex, age, and *APOEε4* carrier status and used likelihood ratio tests to test the significance of improved model fit when including the PRS.

To assess the consistency of individual-level rankings across different PRS models, we calculated Spearman's rank correlation applied quantile mapping across models. Furthermore, we conducted pairwise comparisons between models that used the same GWAS summary statistics (IGAP) to evaluate the consistency of genetic risk classification. For example, we identified participants in the top decile of genetic risk according to the C+T,  $p < 5e-08$  model and calculated the distribution of these participants across all deciles in the other models. If two models had perfectly consistent rankings, 100% of the participants in the top decile of the first model would be in the top decile according to the second model. Pairwise comparisons were made among all MESA participants and additionally among restricted to those who self-reported as white and most closely match the ancestry of the IGAP GWAS.

### *Sensitivity Analysis:*

In addition to polygenic risk scores derived primarily from GWAS of Alzheimer's case-control status in European ancestry samples, we also compared polygenic hazard scores (PHS) based on a multi-ancestry GWAS of age-at-onset (AAO) of Alzheimer's disease that included participants with European, African American, Hispanic, and East Asian ancestry[33]. Like the PRS models, PHS models were constructed using  $p$ -value thresholds of  $p < 0.01$ ,  $1e-05$ , and  $5e-08$ .

While time-to-event data was available, time to hospitalization or death due to dementia may not accurately capture dementia symptom onset or diagnosis. In addition to the comparisons of association and predictive performance based on hazard models, we also fit univariate logistic regression models and calculated the area under the curve (AUC) for each PRS method.

## **Results**

### *Study population and baseline characteristics*

We calculated PRS and inferred global ancestry proportions for all 8,224 participants in the MESA cohort who consented to genetic analyses as part of the SNP Health Association Resource (SHARe). Of these individuals, 6,338 participants had dementia follow-up data and were included in this study. At enrollment, the mean age of the participants was 62. After a median follow up of 16.8 years, 560 (8.8%) incident all-cause dementia events were observed. Complete demographic characteristics are provided in **Table 1**. We estimated global ancestry proportions for all genotyped MESA participants. The African Americans and Hispanic/Latino groups have high amounts of admixture of NFE and AFR ancestry and NFE, AMR, and AFR ancestry, respectively. (**Supplementary Figure 1**). The low NFE subgroup included 2,431 participants, intermediate NFE included 1,256, and high NFE included 2,651.

### *PRS distributions*

**Table 2** outlines the number of SNPs included in each PRS model. The PRS-CSx model incorporated the largest number of SNPs (968,595) that intersect across the 1KG linkage disequilibrium reference maps, IGAP European ancestry GWAS summary statistics, African Genome Resources Panel GWAS summary statistics, and MESA genotyped + imputation data. The PRS-CS models with the European IGAP GWAS and cross-population GWAS included 862,647 SNPs and 851,128 SNPs, respectively. In contrast, the C+T model with a stringent genome-wide significant  $p$ -value cutoff ( $p < 5e-08$  C+T) included only 15 SNPs after filtering for linkage disequilibrium.

For all models, marked differences in PRS distributions were observed across self-reported racial groups using raw scores, although the differences are less pronounced in the conservative C+T models. After adjusting for population structure, the variation was attenuated (**Figure 1, Supplementary Figure 2**)

### *Association between PRS and Incident Dementia*

Univariate Cox proportional hazards models were fit to test the association of each PRS model with incident dementia. In the full sample, PRS derived from all models except for the C+T model with  $p$ -value threshold  $< 0.1$  were associated with incident dementia (**Figure 2**). The hazard ratio for the PRS constructed from the most stringent  $p$ -value threshold ( $p < 5e-08$  C+T) was slightly higher than the other models, but this difference was not statistically significant ( $HR_{total\_5e-08} = 1.21$ , 95% CI: 1.11-1.31, **Figure 2**).

Among the race/ethnicity-stratified models, the cross-population PRS-CS model was associated with dementia in Black/African American ( $HR_{AA\_cs-cp} = 1.24$ , 95% CI: 1.03-1.49), Hispanic/Latino ( $HR_{HIS\_cs-cp} = 1.26$ , 95% CI: 1.03-1.54), and White groups ( $HR_{WHI\_cs-cp} = 1.18$ ,

95% CI: 1.06-1.32). The  $p < 5e-08$  C+T model was associated with incident dementia among the African American/Black ( $HR_{AA\_5e-08} = 1.33$ , 95% CI: 1.11-1.59) and White participants ( $HR_{WHI\_5e-08} = 1.24$ , 95% CI: 1.11-1.40). No PRS were significantly associated with dementia in the Chinese participants, likely due to the small sample size and limited number of dementia cases.

Under the assumption that model performance may be more dependent on genetic similarity to the GWAS sample ancestry than self-reported race/ethnicity, we also stratified MESA participants into tertiles of NFE ancestry. Among the low, intermediate, and high NFE groups, we found that in the high NFE group ( $p_{NFE} > 0.67$ ), all methods aside from the  $p < 0.1$  C+T model resulted in PRS associated with dementia hazard (**Supplementary Figure 3**). In contrast, none of the methods resulted in PRS associated with dementia in the intermediate NFE group and only the PRS calculated from the  $p < 5e-08$  C+T model was associated with dementia in the low NFE group ( $HR_{lowNFE\_5e-08} = 1.26$ , 95% CI: 1.08-1.47).

As a sensitivity analysis, we also fit univariate logistic regression models to assess the association between the PRS models and dementia case-control status. The results from the logistic regression model mirror findings from the Cox proportional hazards models. Among all MESA participants, the  $p < 5e-08$  C+T model had stronger estimated association with dementia status compared to other models, but the difference was not statistically significant (**Supplementary Figure 4**).

#### *Assessing Model Performance*

Model performance was evaluated using Harrell's concordance index (C-index), with the highest value observed for the  $p < 5e-08$  C+T model ( $C_{5e-08} = 0.55$ , standard deviation (SD) = 0.01). Comparisons across models revealed that the inclusion of more SNPs in either C+T

models using Bayesian approaches does not improve predictive accuracy (**Figure 3**). These findings were also supported by comparisons of model AUC ( $AUC_{5e-08} = 0.56$ ,  $SD = 0.01$ , **Supplementary Figure 5**).

In groups with low and high proportions of NFE, the  $p < 5e-08$  C+T model had the highest C-index ( $C_{lowNFE_{5e-08}} = 0.57$ ,  $SD = 0.02$ ;  $C_{highNFE_{5e-08}} = 0.56$ ,  $SD = 0.02$ , **Supplementary Figure 6**). In the group with intermediate NFE proportion, the PRS-CSx model had the best performance ( $C_{midNFE_{csx}} = 0.54$ ,  $SD = 0.03$ ) while the  $p < 5e-08$  C+T model had the worst performance ( $C_{midNFE_{5e-08}} = 0.51$ ,  $SD = 0.02$ ).

Despite the greater similarity in ancestry background, the PHS models did not outperform the  $p < 5e-08$  C+T model derived from the European ancestry GWAS when analyzing all MESA participants as a whole or in groups stratified by self-reported race/ethnicity. (**Supplementary Figure 7**).

#### *Value added from prediction using PRS*

Compared to a baseline model that included age, sex, and *APOE* genotype, the addition of PRS derived from  $p < 5e-08$  and  $p < 1e-05$  C+T models and the PRS-CSx model led to a marginal increase in C-index. The baseline model had a C-index of 0.84. In the  $p < 5e-08$  C+T model,  $p < 1e05$  C+T model, and PRS-CSx model, inclusion of the PRS significantly improved model fit ( $p_{LRT_{5e-08}} = 5e-05$ ,  $p_{LRT_{1e05}} = 0.008$ ,  $p_{LRT_{CSx}} = 0.001$ , **Table 3**). The remaining PRS models did not significantly change model fit.

#### *Correlation of PRS across Models*

The Spearman rank correlation analysis showed varying degrees of correlation between the different PRS models. The C+T models with  $p < 5e-08$  and  $p < 1e-05$  were strongly correlated ( $R = 0.64$ ). Correlations among the Bayesian models were also strong ( $R = 0.45$  to  $0.70$ , **Table**

4). However, correlations between the  $p < 5e-08$  C+T model and the Bayesian models were low ( $R = 0.11$  to  $0.25$ , **Table 4**), suggesting that being classified as high risk using the  $p < 5e-08$  C+T model is not predictive of being high risk based on the Bayesian models.

To directly examine the consistency in classification of high or low risk across models, we conducted pairwise comparisons of inter-PRS model reliability. Of those who are in the top risk decile of PRS in the  $p < 5e-08$  C+T model, 43% are considered in the top decile in the  $p < 1e-05$  C+T model, while only 16% are considered in the top risk decile in the PRS-CS model derived from the same summary statistics (**Figure 4, Supplementary Figure 8**).

## Discussion

PRS are increasingly being used for assessing genetic susceptibility for a wide spectrum of diseases, allowing for earlier identification of individuals at higher risk. This could have a great impact on dementia, a disease that is difficult to treat, let alone predict or prevent, but few studies have assessed the association of AD PRS for predicting dementia in diverse populations. Our study demonstrates that PRS models, even when excluding the *APOE* region, remain significantly associated with incident dementia in a multi-ancestry sample. We also found that including more SNPs does not improve predictive performance. A smaller panel of variants is more robustly associated with incident dementia across diverse groups than the larger models.

However, even the best performing PRS models in our study have low predictive power (C-index  $< 0.6$ ) and add only slight improvements to predictive models that include age, sex, and *APOE*. Early prediction of dementia will likely require integrating demographic and environmental information in addition to genetics. The low correlation across PRS models is also a major concern that must be addressed before implementation, as our results suggest that differing models can lead to significantly different conclusions of an individual's risk percentile.

Our comparisons also demonstrate the benefit of diversifying genomic studies of AD, adding to the growing calls for diversity across genetic research. The PRS-CS model using summary statistics from cross-population GWAS of AD consistently performed better than the PRS-CS model using summary statistics from the European-ancestry GWAS, despite the smaller sample size of the cross-population GWAS. This finding is consistent with previous work examining other non-dementia traits that demonstrated the superior performance of PRS derived from multi-ancestry GWAS meta-analyses compared to single-ancestry GWAS[34]. Of note, neither of the PRS-CS models outperformed the restrictive C+T model with 15 SNPs derived from the European ancestry GWAS. It's likely that the restricted set of SNPs are more likely to tag regions that have true biological impact on risk of disease development.

PRS are least predictive in individuals with high amounts of genetic admixture – those in the intermediate NFE proportion group. We observed that PRS-CSx performed best in the group with intermediate proportions of NFE ancestry. This aligns with previous findings that have observed increased predictive performance in admixed groups when using a linear combination of summary statistics that combine ancestry-specific effect sizes[35]. Nevertheless, the PRS-CSx performance in this group remained lower than the best performing models in the high-NFE group and the PRS derived from PRS-CSx was not associated with incident dementia (**Supplementary Figure 3**) in this group. The limited association could be due to the relatively small sample size of the African Genome Resources Panel GWAS and lack of GWAS information from studies with substantial AMR ancestry.

Our study has several limitations. Incident dementia cases were ascertained from hospitalization and death records using ICD codes, and this method likely underestimates the true incidence because a portion of individuals living with dementia will not be hospitalized or

have an alternative listed cause of death. However, the external validation by physician review of electronic health records and the association of *APOE* genotype and our polygenic risk score with incident dementia suggests that dementia cases are true positives. Furthermore, there are larger GWAS of dementia-by-proxy phenotypes conducted in European ancestry samples that provide a larger pool of SNPs considered to be significantly associated with parental history of dementia. Due to the sub-optimal dementia adjudication of our target data, we chose to prioritize depth of phenotyping over sample size. We also limited our summary statistics to variants identified in GWAS and, therefore, did not consider rare variants that may have large effects on Alzheimer's disease. In addition, numerous new polygenic risk approaches are constantly being developed and we are unable to test all possible methods. Instead, we selected methods that are shown to perform well in the absence of individual level validation data due to the lack of datasets that match the diversity of our target sample and are not included in the GWAS from which the summary statistics are derived. Finally, while Chinese participants were included in our study, the sample size was too small to detect differences in performance across PRS models and none of the models resulted in PRS associated with dementia in this subgroup.

While our findings show that current PRS models have modest predictive value, future Alzheimer's disease GWAS in diverse populations have the potential to enhance their predictive power. Furthermore, the utility of PRS may extend beyond estimating the overall likelihood of disease development. Future research should focus on how PRS can be used to identify differences in disease pathogenesis and clinical trajectories. Incorporation of rare variants and leveraging functional annotation to develop pathway specific scores will further enhance the precision and translational value of PRS.

## References

- [1] Nichols E, Steinmetz JD, Vollset SE, Fukutaki K, Chalek J, Abd-Allah F, et al. Estimation of the global prevalence of dementia in 2019 and forecasted prevalence in 2050: an analysis for the Global Burden of Disease Study 2019. *Lancet Public Health* 2022;7:e105–25. [https://doi.org/10.1016/S2468-2667\(21\)00249-8](https://doi.org/10.1016/S2468-2667(21)00249-8).
- [2] 2024 Alzheimer’s disease facts and figures. *Alzheimer’s & Dementia* 2024;20:3708–821. <https://doi.org/https://doi.org/10.1002/alz.13809>.
- [3] Champion D, Dumanchin C, Hannequin D, Dubois B, Belliard S, Puel M, et al. Early-onset autosomal dominant Alzheimer disease: prevalence, genetic heterogeneity, and mutation spectrum. *The American Journal of Human Genetics* 1999;65:664–70.
- [4] Pericak-Vance MA, Bebout JL, Gaskell PC, Yamaoka LH, Hung W-Y, Alberts MJ, et al. Linkage studies in familial Alzheimer disease: evidence for chromosome 19 linkage. *Am J Hum Genet* 1991;48:1034.
- [5] Bertram L, Lange C, Mullin K, Parkinson M, Hsiao M, Hogan MF, et al. Genome-wide association analysis reveals putative Alzheimer’s disease susceptibility loci in addition to APOE. *The American Journal of Human Genetics* 2008;83:623–32.
- [6] Andrews SJ, Renton AE, Fulton-Howard B, Podlesny-Drabiniok A, Marcora E, Goate AM. The complex genetic architecture of Alzheimer’s disease: novel insights and future directions. *EBioMedicine* 2023;90. <https://doi.org/10.1016/j.ebiom.2023.104511>.
- [7] Leonenko G, Baker E, Stevenson-Hoare J, Sierksma A, Fiers M, Williams J, et al. Identifying individuals with high risk of Alzheimer’s disease using polygenic risk scores. *Nature Communications* 2021 12:1 2021;12:1–10. <https://doi.org/10.1038/s41467-021-24082-z>.
- [8] de Rojas I, Moreno-Grau S, Tesi N, Grenier-Boley B, Andrade V, Jansen IE, et al. Common variants in Alzheimer’s disease and risk stratification by polygenic risk scores. *Nature Communications* 2021;12. <https://doi.org/10.1038/s41467-021-22491-8>.
- [9] Lambert SA, Abraham G, Inouye M. Towards clinical utility of polygenic risk scores. *Hum Mol Genet* 2019;28:R133–42. <https://doi.org/10.1093/hmg/ddz187>.
- [10] Choi SW, Mak TS-H, O’Reilly PF. Tutorial: a guide to performing polygenic risk score analyses. *Nat Protoc* 2020;15:2759–72.
- [11] Escott-Price V, Sims R, Bannister C, Harold D, Vronskaya M, Majounie E, et al. Common polygenic variation enhances risk prediction for Alzheimer’s disease. *Brain* 2015;138:3673–84. <https://doi.org/10.1093/brain/awv268>.

- [12] Leonenko G, Shoai M, Bellou E, Sims R, Williams J, Hardy J, et al. Genetic risk for Alzheimer disease is distinct from genetic risk for amyloid deposition. *Ann Neurol* 2019;86:427–35. <https://doi.org/https://doi.org/10.1002/ana.25530>.
- [13] Zhang Q, Sidorenko J, Couvy-Duchesne B, Marioni RE, Wright MJ, Goate AM, et al. Risk prediction of late-onset Alzheimer’s disease implies an oligogenic architecture. *Nat Commun* 2020;11:4799. <https://doi.org/10.1038/s41467-020-18534-1>.
- [14] Martin AR, Kanai M, Kamatani Y, Okada Y, Neale BM, Daly MJ. Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat Genet* 2019;51:584–91. <https://doi.org/10.1038/s41588-019-0379-x>.
- [15] Privé F, Aschard H, Carmi S, Folkersen L, Hoggart C, O’Reilly PF, et al. Portability of 245 polygenic scores when derived from the UK Biobank and applied to 9 ancestry groups from the same cohort. *The American Journal of Human Genetics* 2022;109:12–23. <https://doi.org/10.1016/j.ajhg.2021.11.008>.
- [16] Ding Y, Hou K, Xu Z, Pimplaskar A, Petter E, Boulier K, et al. Polygenic scoring accuracy varies across the genetic ancestry continuum. *Nature* 2023;618:774–81. <https://doi.org/10.1038/s41586-023-06079-4>.
- [17] Kunkle BW, Grenier-Boley B, Sims R, Bis JC, Damotte V, Naj AC, et al. Genetic meta-analysis of diagnosed Alzheimer’s disease identifies new risk loci and implicates A $\beta$ , tau, immunity and lipid processing. *Nat Genet* 2019;51:414–30. <https://doi.org/10.1038/s41588-019-0358-2>.
- [18] Bellenguez C, Küçükali F, Jansen IE, Kleindam L, Moreno-Grau S, Amin N, et al. New insights into the genetic etiology of Alzheimer’s disease and related dementias. *Nat Genet* 2022;54:412–36. <https://doi.org/10.1038/s41588-022-01024-z>.
- [19] Reitz C, Pericak-Vance MA, Foroud T, Mayeux R. A global view of the genetic basis of Alzheimer disease. *Nat Rev Neurol* 2023;19:261–77.
- [20] Xue D, Blue EE, Conomos MP, Fohner AE. The power of representation: Statistical analysis of diversity in US Alzheimer’s disease genetics data. *Alzheimer’s and Dementia: Translational Research and Clinical Interventions* 2024;10. <https://doi.org/10.1002/trc2.12462>.
- [21] Bild DE, Bluemke DA, Burke GL, Detrano R, Diez Roux A V, Folsom AR, et al. Multi-Ethnic Study of Atherosclerosis: Objectives and Design. *Am J Epidemiol* 2002;156:871–81. <https://doi.org/10.1093/aje/kwf113>.
- [22] Maples BK, Gravel S, Kenny EE, Bustamante CD. RFMix: A discriminative modeling approach for rapid and robust local-ancestry inference. *Am J Hum Genet* 2013;93:278–88. <https://doi.org/10.1016/j.ajhg.2013.06.020>.

- [23] Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* 2020;581:434–43. <https://doi.org/10.1038/s41586-020-2308-7>.
- [24] Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti A, et al. A global reference for human genetic variation. *Nature* 2015;526:68–74. <https://doi.org/10.1038/nature15393>.
- [25] Fujiyoshi A, Jacobs Jr DR, Alonso A, Luchsinger JA, Rapp SR, Duprez DA. Validity of Death Certificate and Hospital Discharge ICD Codes for Dementia Diagnosis: The Multi-Ethnic Study of Atherosclerosis. *Alzheimer Dis Assoc Disord* 2017;31:168–72. <https://doi.org/10.1097/WAD.000000000000164>.
- [26] Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen W-M. Robust relationship inference in genome-wide association studies. *Bioinformatics* 2010;26:2867–73.
- [27] Chang CC, Chow CC, Tellier LCAM, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: Rising to the challenge of larger and richer datasets. *Gigascience* 2015. <https://doi.org/10.1186/s13742-015-0047-8>.
- [28] Ge T, Chen C-Y, Ni Y, Feng Y-CA, Smoller JW. Polygenic prediction via Bayesian regression and continuous shrinkage priors. *Nat Commun* 2019;10:1776. <https://doi.org/10.1038/s41467-019-09718-5>.
- [29] Ruan Y, Lin Y-F, Feng Y-CA, Chen C-Y, Lam M, Guo Z, et al. Improving polygenic prediction in ancestrally diverse populations. *Nat Genet* 2022;54:573–80. <https://doi.org/10.1038/s41588-022-01054-7>.
- [30] Jun GR, Chung J, Logue MW, Sherva R, Farrer LA, Mez J, et al. Transethnic genome-wide scan identifies novel Alzheimer’s disease loci. *Alzheimer’s and Dementia* 2017;13:727–38. <https://doi.org/10.1016/j.jalz.2016.12.012>.
- [31] Kunkle BW, Schmidt M, Klein HU, Naj AC, Hamilton-Nelson KL, Larson EB, et al. Novel Alzheimer Disease risk loci and pathways in african American individuals using the african genome resources panel a meta-analysis. *JAMA Neurol* 2020;Published. <https://doi.org/10.1001/jamaneurol.2020.3536>.
- [32] Hao L, Kraft P, Berriz GF, Hynes ED, Koch C, Korategere V Kumar P, et al. Development of a clinical polygenic risk score assay and reporting workflow. *Nat Med* 2022;28:1006–13.
- [33] Correction to the 2024 Annual Meeting of the International Genetic Epidemiology Society. *Genet Epidemiol* 2024;n/a. <https://doi.org/10.1002/gepi.22599>.

- [34] Gunn S, Wang X, Posner DC, Cho K, Huffman JE, Gaziano M, et al. Comparison of methods for building polygenic scores for diverse populations. *Human Genetics and Genomics Advances* 2025;6. <https://doi.org/10.1016/j.xhgg.2024.100355>.
- [35] Bitarello BD, Mathieson I. Polygenic Scores for Height in Admixed Populations. *G3 Genes|Genomes|Genetics* 2020;10:4027–36. <https://doi.org/10.1534/g3.120.401658>.

## Tables

**Table 1. Demographic Characteristics of MESA participants at Baseline**

	<b>Censored (N=5778)</b>	<b>Affected (N=560)</b>	<b>Overall (N=6338)</b>
<b>Age at baseline (years)</b>			
Mean (SD)	61.3 (9.95)	72.2 (7.71)	62.2 (10.2)
Median [Min, Max]	61.0 [44.0, 84.0]	74.0 [45.0, 84.0]	62.0 [44.0, 84.0]
<b>Sex</b>			
Female	3035 (52.5%)	283 (50.5%)	3318 (52.4%)
Male	2743 (47.5%)	277 (49.5%)	3020 (47.6%)
<b>Race/ethnicity</b>			
Black	1457 (25.2%)	146 (26.1%)	1603 (25.3%)
Chinese	734 (12.7%)	40 (7.1%)	774 (12.2%)
Hispanic/Latino	1327 (23.0%)	115 (20.5%)	1442 (22.8%)
White	2260 (39.1%)	259 (46.3%)	2519 (39.7%)
<b>Level of education</b>			
< High school degree	1030 (17.8%)	126 (22.5%)	1156 (18.2%)
High school degree	1021 (17.7%)	137 (24.5%)	1158 (18.3%)
Some college, no bachelor's degree	1639 (28.4%)	141 (25.2%)	1780 (28.1%)
Bachelor's degree or higher	2069 (35.8%)	155 (27.7%)	2224 (35.1%)
Missing	19 (0.3%)	1 (0.2%)	20 (0.3%)
<b>APOE genotype</b>			
<i>e2/e2</i>	45 (0.8%)	2 (0.4%)	47 (0.7%)
<i>e2/e3</i>	694 (12.0%)	58 (10.4%)	752 (11.9%)
<i>e2/e4</i>	147 (2.5%)	14 (2.5%)	161 (2.5%)
<i>e3/e3</i>	3471 (60.1%)	296 (52.9%)	3767 (59.4%)
<i>e3/e4</i>	1196 (20.7%)	160 (28.6%)	1356 (21.4%)
<i>e4/e4</i>	128 (2.2%)	20 (3.6%)	148 (2.3%)
Missing	97 (1.7%)	10 (1.8%)	107 (1.7%)

**Table 2. Description of PRS models**

Table 2 Legend. This table shows the PRS models being compared. The models differ in computational method, *p*-value threshold and corresponding and the number of SNPs included, and/or GWAS training data. The NIAGADS study number is provided for the GWAS summary statistics used.

Model	Method	P-value cutoff	# SNPs	GWAS ancestry	GWAS Cases/Controls	GWAS study name [NIAGADS study ID]
1	C+T	5e-08	15	NFE-like	21,982/41,944	International Genomics of Alzheimer's Project [NG00075, 2019]
2		1e-05	53			
3		0.1	9,023			
4	PRS-CS	NA	862,647	Multi-ancestry	15,579/17,690	ADGC multi-ancestry [NG00056, 2017]
5			851,128			
6	PRS-CSx		968,595	NFE-like	21,982/41,944	[NG00075, 2019]
				AFR-like	2,748/5,222	African Genome Resources Panel [NG00100, 2021]

### Table 3. Value Added of PRS

Table 3 Legend. We computed the added predictive value of each PRS method compared to baseline models that included sex, age, and *APOEε4* carrier status and used likelihood ratio tests to test the significance of improved model fit when including the PRS. This table shows the change in concordance from adding PRS to the model and the p-value corresponding to the likelihood ratio test. The modest changes in C-index across all models likely reflect that *APOE*, sex, and age already account for a substantial portion of the variance ( $C = 0.84$ ).

PRS Model	$\Delta$ Concordance	p-value <sub>LRT</sub>
C+T, p <5e-08	0.0015	<b>5e-05</b>
C+T, p <1e-05	0.0007	<b>0.008</b>
C+T, p <0.01	0.0001	0.903
PRS-CS, IGAP	0.0001	0.220
PRS-CS, Cross-Pop	0.0002	0.128
PRS-CSx	0.0014	<b>0.001</b>

**Table 4. Spearman Rank Correlation coefficients across different PRS methods.**

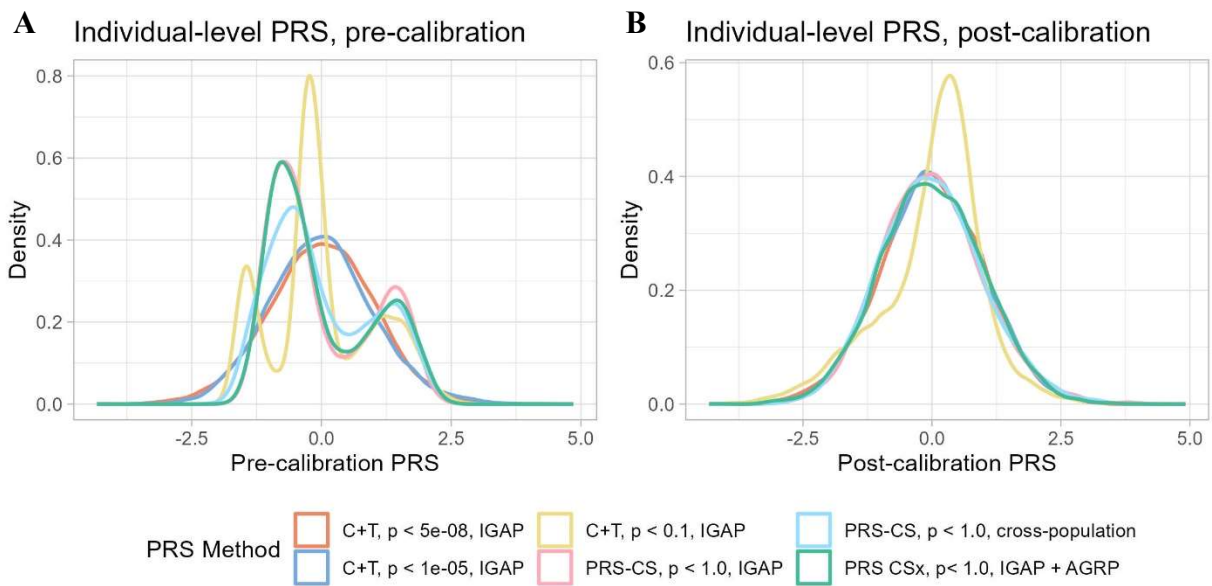
Table 4 Legend. We computed the correlation across scores using Spearman rank correlation coefficients. Correlation values range from -1 to 1, with higher values indicating stronger concordance between the different methods.

	<b>C+T, p &lt;5e-08</b>	<b>C+T, p &lt;1e-05</b>	<b>C+T, p &lt;0.01</b>	<b>PRS-CS, IGAP</b>	<b>PRS-CS, Cross-Pop</b>	<b>PRS-CSx</b>
<b>C+T, p &lt;5e-08</b>	1.00					
<b>C+T, p &lt;1e-05</b>	0.64	1.00				
<b>C+T, p &lt;0.01</b>	0.08	0.12	1.00			
<b>PRS-CS, IGAP</b>	0.18	0.18	0.16	1.00		
<b>PRS-CS, Cross-Pop</b>	0.11	0.11	0.07	0.70	1.00	
<b>PRS-CSx</b>	0.25	0.26	0.24	0.58	0.45	1.00

## Figures

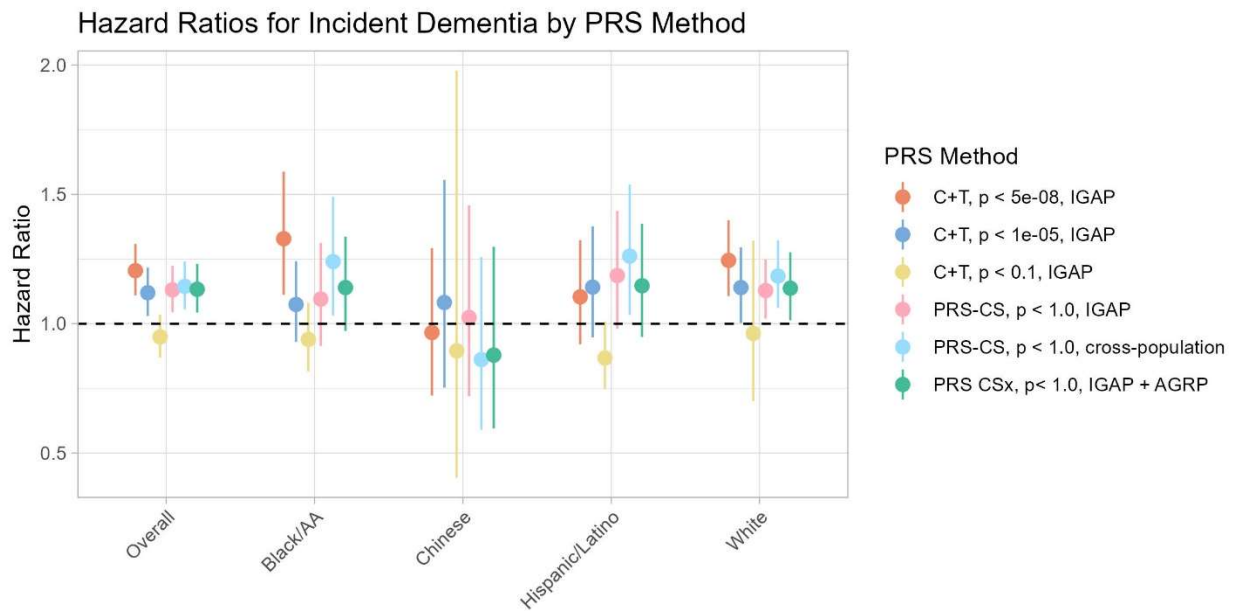
### Figure 1. PRS distributions before and after calibration by principal components.

Figure 1 Legend. Density plots show the distribution of risk scores for each PRS method across all participants. A) Distributions of scores that have been mean-standardized but not adjusted for principal components. B) Distributions of scores after standardization and principal component calibration.



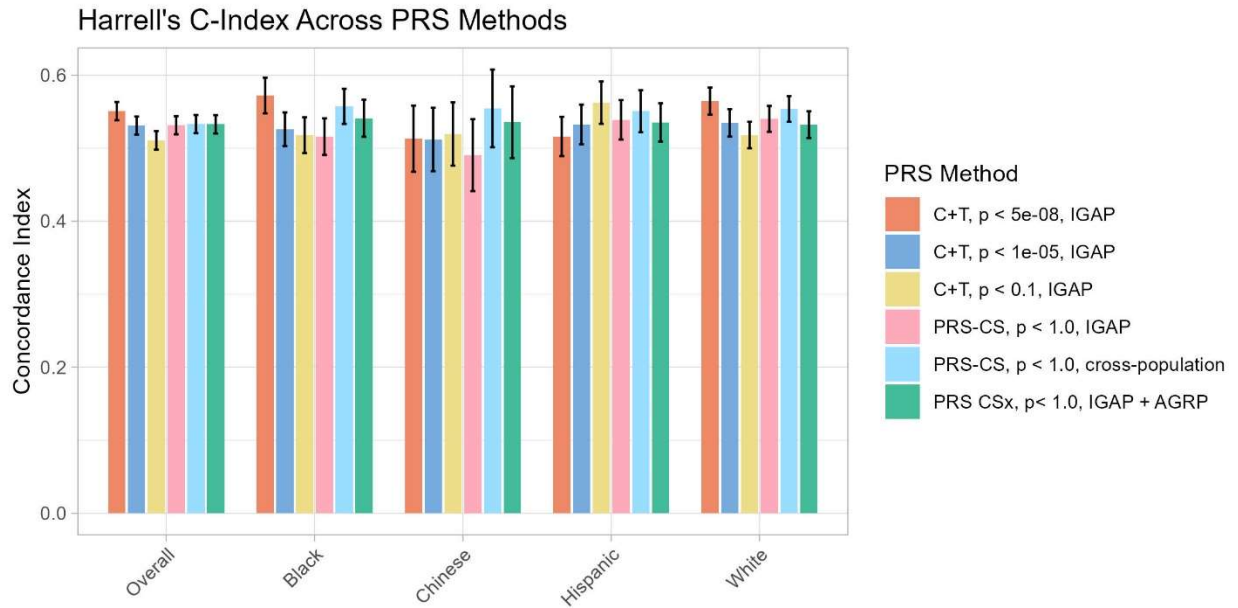
**Figure 2. Association between PRS and incident dementia.**

Figure 2 Legend. The forest plots display the hazard ratio and confidence intervals for univariate models with the PRS as exposure and incident dementia outcome. Results are presented for all participants (Overall) and for groups stratified by self-reported race/ethnicity. Overall, among all participants, only the C+T methods with restrictive P-value cutoffs are significantly associated with dementia.



**Figure 3. PRS predictive performance measured by Harrell's Concordance Index.**

Figure 3 Legend. The bar plots show a comparison of prediction accuracy as measured by the concordance index or Harrell's C. Prediction accuracy is best overall using the C+T genome-wide significant cutoff.



**Figure 4. Inter-model reliability.**

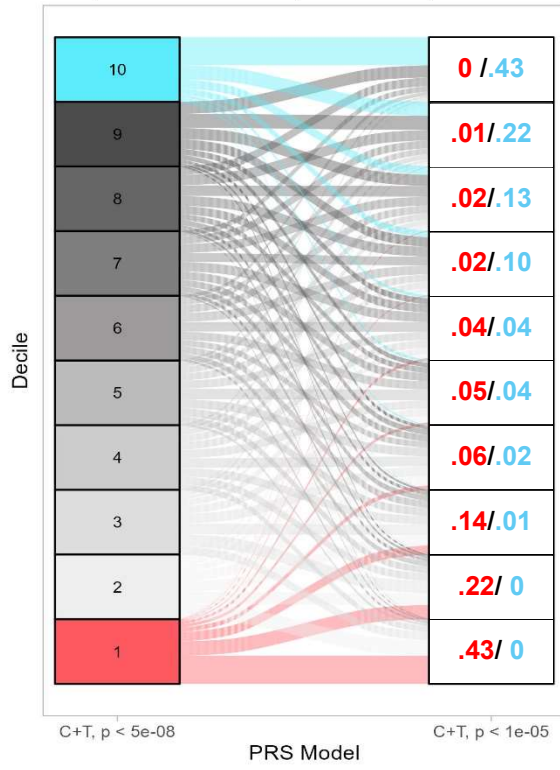
Figure 4 Legend. The Sankey plots display pairwise comparisons for the models constructed using the same GWAS summary statistics but vary in  $p$ -value threshold and method.

**Red** indicates being in the lowest risk decile in the lefthand model.

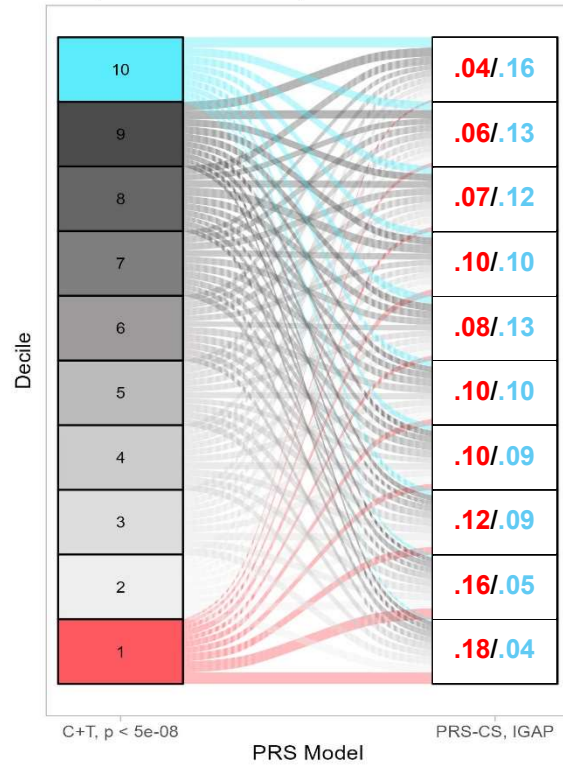
**Blue** indicates being in the highest risk decile in the lefthand model.

The **red** and **blue** proportions reported for each decile indicate the proportion of those in the lowest and highest risk deciles, respectively, that are in each risk decile according to the righthand model.

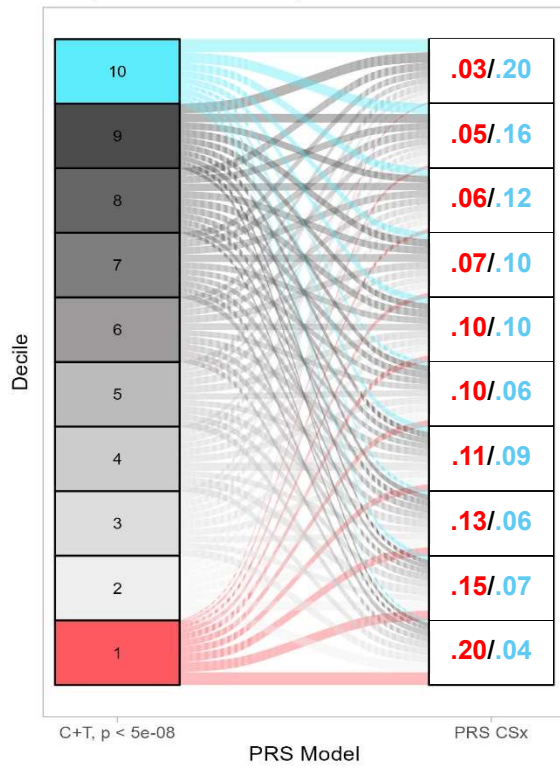
C+T p<5e-08 deciles compared to C+T p<1e-05



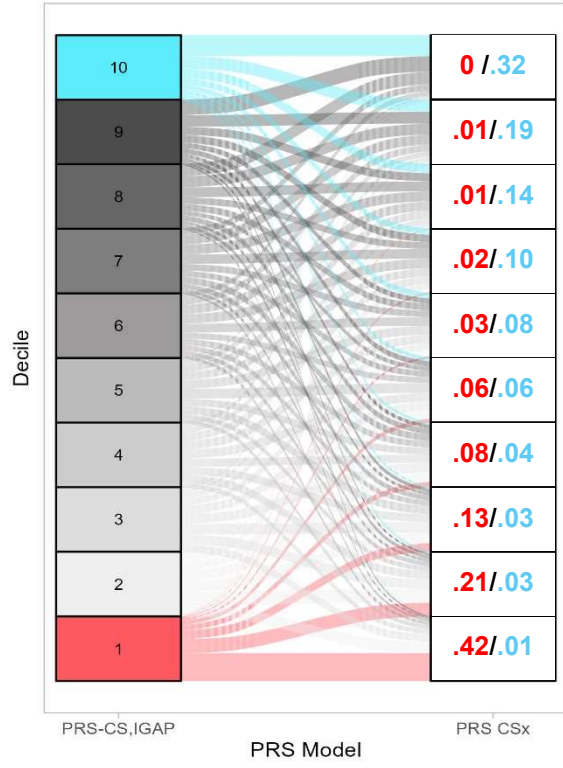
C+T p<5e-08 deciles compared to PRS-CS



C+T p<5e-08 deciles compared to PRS-CSx



PRS-CS deciles compared to PRS-CSx



## Paper 3: Integrating Contextual Determinants and Polygenic Risk to Examine Dementia and Cognition in the Multi-Ethnic Study of Atherosclerosis

### Abstract

**Background:** Alzheimer's disease, the primary cause of dementia, has a complex etiology involving genetic and environmental risk factors. Few studies have examined the joint effects of neighborhood context and genome-wide risk. This study investigates the effect of contextual exposures on dementia and late-life cognition, both independently and in interaction with polygenic risk.

**Methods:** We analyzed data from the Multi-Ethnic Study of Atherosclerosis, assessing the effects of seven contextual determinants: neighborhood socioeconomic status (SES), nitrogen dioxide, particulate matter  $< 2.5\mu\text{g}$ , and distances to the nearest favorable food stores, fast food chains, alcohol establishments, and physical activity facilities. Dementia incidence was tracked from baseline (2000-2002) through 2018 using ICD codes obtained from hospital records and death certificates. Proportional hazards regression and generalized estimating equation models were used to examine the relationships between contextual determinants and incident dementia and late-life cognition. We also tested for interactions between each environmental feature and genetic risk. Finally, we implemented Bayesian Kernel Machine Regression to investigate the joint effects of all seven contextual exposures on late-life cognition.

**Results:** While specific pollutants and built environment characteristics are not associated with dementia after controlling for genetic risk, neighborhood SES is associated with incident dementia ( $\text{HR}_{\text{SES}} = 1.14$ ; 95% CI: 1.05-1.25) and late-life cognition ( $\beta_{\text{SES}} = -0.30$ , 95% CI: -0.49 – -0.11) after controlling for genetic risk factors and other covariates. We also find that neighborhood SES is most strongly associated with hazard of dementia in the highest genetic risk

group. Overall, the effects of other contextual determinants did not differ significantly across groups stratified by genetic risk for Alzheimer's disease.

**Conclusion:** Neighborhood socioeconomic status was significantly associated with incident dementia and late-life cognition, independent of *APOE* and other polygenic risk factors. Our findings suggest that not only is SES associated with dementia, but there is also evidence for gene-environment interaction, with neighborhood SES having stronger effects on dementia risk in groups with high genetic risk. Further research is needed to understand the causal mechanisms linking SES and dementia outcomes.

## **Abbreviations**

Bayesian kernel machine regression (BKMR),

generalized estimating equation (GEE),

Hazard Ratio (HR),

International Classification of Disease (ICD),

interquartile range (IQR),

Multi-Ethnic Study of Atherosclerosis (MESA),

nitrogen dioxide (NO<sub>2</sub>),

odds ratio (OR),

particulate matter <2.5μg (PM<sub>2.5</sub>),

polygenic risk score (PRS),

socioeconomic status (SES)

## Introduction

Dementia is a growing public health burden. In the United States, dementia is the sixth leading cause of death, and estimated health care costs associated with dementia are projected to surpass \$1 trillion in the next 25 years[1,2]. Delaying the progression of dementia or cognitive decline is crucial to alleviating the burden on patients, caregivers, and health care systems. A simulation study found that a five-year delay in age-at-onset of dementia would lower prevalence of dementia in 2050 by 41%[3]. While recent advances in drug development hold promise, population-level prevention strategies that address modifiable risk factors and delay symptom onset are also crucial for impacting the burden of dementia.

Modifiable risk factors vary in definition, encompassing a wide range of interrelated factors that differ in their scale of actionability. The 2024 Lancet Commission on dementia describes 14 modifiable risk factors for dementia prevention that span the life-course, ranging from exogenous factors such as air pollution and education to individual characteristics including hearing loss, physical inactivity, and hypertension[4]. The commission suggests that nearly half of dementia cases could be prevented by targeting these factors, however, the modifiability of these risk factors varies drastically. While some factors, like obtaining hearing aids, are more directly actionable at the individual level, modifying environmental exposures like air pollution require broader systemic interventions.

Given the high heritability of Alzheimer's disease[5], the role of contextual determinants must also be investigated alongside genetic risk for a more complete understanding of disease etiology. Most studies of modifiable risk factors that do include genetic information have only focused on the *APOE* genotype [6,7]. The *APOE**e4* allele is the strongest genetic risk factor for Alzheimer's disease, but less than half of all patients carry a copy of the *e4* allele[8], and the effect of *e4* may differ across populations due to local ancestry effects[9,10]. Beyond *APOE*,

dozens of disease-associated loci have been identified from genome-wide association studies (GWAS)[11].

Polygenic risk scores (PRS) offer a valuable tool for summarizing the genetic risk captured by GWAS. By aggregating the effects of single nucleotide polymorphisms (SNPs), PRS provide a weighted genetic risk score for each individual[12]. The integration of genetic data with contextual factors remains underexplored, but the use of PRS enables studies of gene-environment interactions that are more feasible than investigating individual SNP-environment interaction effects. PRS also allow us to consider genetic risk even among individuals who do not carry the *APOEε4* allele. Genetic predisposition may modify the strength or direction of environmental effects. Integrating PRS and environmental data can reveal environmental factors that have more substantial effects within certain genetic subgroups— effects that may be masked or attenuated when analyzing the sample as whole.

We use data from the Multi-Ethnic Study of Atherosclerosis (MESA), which along with its dozens of ancillary studies, provides a rich resource for investigating the interplay between genetic and contextual factors. In this study we focus on the role of meso- or neighborhood-level modifiable risk factors on dementia incidence and cognition. Rather than focusing on individual activity levels or diet, for example, we investigate the role of proximity to recreational space and favorable food markets. In addition to adjusting for *APOE* genotype we calculated PRS that summarize the genetic effects of common disease-associated variants. We hypothesize that contextual determinants may have differing effects on dementia across groups differing in underlying genetic risk. This study is one of the first to examine neighborhood-level modifiable risk factors in the context of genome-wide genetic risk for Alzheimer's disease beyond *APOE*.

By integrating these factors, we seek to identify how neighborhood-level modifiable risk factors may affect people differently based on genetic predisposition.

## **Methods**

### *Study Population*

Our sample is derived from MESA, a longitudinal study that enrolled 6,814 adults aged 45-84 years at baseline (2000-2002)[13]. Participants were recruited at one of six sites across the United States: Baltimore, Maryland; Chicago, Illinois; Forsyth County, North Carolina; Los Angeles County, California; Northern Manhattan and the Bronx, New York; and Saint Paul, Minnesota. Recruitment was designed to target four racial ethnic groups at approximately equal gender ratios: Black/African American, Chinese, Hispanic/Latino, and White. Individuals with missing values for dementia outcomes were excluded. Informed consent was obtained from each participant at baseline and at each subsequent exam. Approval was received at each site from the local institutional review board for each examination.

### *Outcome Assessment*

#### *A. Dementia*

Participants were followed up by phone call every 9 to 12 months to update hospital admissions and deaths. Incident dementia was ascertained based on date of hospitalization and/or death certificate International Classification of Diseases (ICD) codes. A physician blinded to the ICD codes validated diagnoses by analyzing text from electronic health records for language that pointed to significant cognitive decline and dementia[14].

#### *B. Global Cognition*

The cognitive abilities screening instrument (CASI) score was measured at MESA Exam 5. The CASI has a score range of 0 to 100 and provides a global cognitive assessment covering the following domains: attention, concentration, orientation, short-term memory, long-term

memory, language abilities, visual construction, list-generating fluence, abstraction, and judgement[15].

### *Exposure Assessment*

We selected seven contextual determinants to capture the social environment, chemical environment (traffic- and fossil-fuel related pollution), and built environment (including proximity to food and alcohol businesses and recreational facilities).

#### *A. Neighborhood Socioeconomic Status (SES)*

The social environment is captured by the MESA area-level SES index, what we will refer to as the neighborhood SES. Neighborhood SES was available at the US Census Bureau tract level and was based on 2010–2011 American Community Survey data. The SES index used here was previously derived in a principal components analysis based on seven neighborhood characteristics (created from MESA Neighborhood Study[16]: 1) percent of neighborhood residents with a managerial occupation; 2) percent of residents with a bachelor’s degree; 3) percent of residents with a high school degree; 4) percent of residents with an annual household income >\$50,000; 5) median household income; 6) percent rental income; 7) and median home value. Higher values of the index indicate worse SES.

#### *B. Ambient Air Pollutants*

We assessed the role of long-term exposure levels to nitrogen dioxide (NO<sub>2</sub>) and particulate matter <2.5µm in aerodynamic diameter (PM<sub>2.5</sub>). Measurement of long-term individual-level exposure to ambient air pollutants through the ancillary MESA Air study has been previously described[17]. Briefly, MESA Air used national universal kriging models to estimate ambient concentrations of PM<sub>2.5</sub> and NO<sub>2</sub> localized to each participant’s residential address. The concentrations used here are three-year averages of two-week measurements

spanning from 1999-2002. All pollutants measures were converted to an interquartile range (IQR) scale.

### *C. Built Environment Features*

The built environment measures were chosen to represent the food and physical activity environment for each participant. For each of the following categories, we used the Euclidean distance to the nearest establishment: favorable food stores, fast food chains, alcohol outlets (liquor stores and on-site drinking places), and physical activity facilities (including indoor and outdoor recreational facilities). The distances were  $\log_2$  transformed due to the skewed distributions of the raw values (**Supplementary Figures 1-2**). A one-unit increase in a  $\log_2$ -transformed variables represents a two-fold increase, allowing for greater spread in lower values. Classification of business data for food stores and recreational facilities are derived from the National Establishment Time Series (NETS) data for all zip codes within a 5-mile buffer of MESA address of record. The data are linked to the MESA address provided by participants at the year of the exam.

### *D. Genetics*

We modeled *APOEε4* carrier status separately from the polygenic risk score, which excludes the *APOE* region. The PRS is derived from summary statistics from the largest genome-wide association study of clinically adjudicated AD, the International Genomics of Alzheimer's Project Stage I results[18]. The PRS has been adjusted for population structure, captured by the first three principal components. PLINK 1.9 was used to calculate the PRS using clumping and thresholding with a genome-wide significant  $p$ -value threshold ( $5e-08$ ) and  $r$ -squared  $< 0.01$ .

### *Covariates*

In addition to the primary contextual exposure of interest, all models adjusted for age, sex, site, income, education, smoking, marital status, and race/ethnicity. All models for ambient air pollutants and built environment features were also adjusted for neighborhood SES, as SES can affect variations in these factors and have direct associations with dementia and cognition, independent of the pathways involving pollution or built environment characteristics.

### *Statistical Analysis*

#### *A. Generalized Linear Models*

To examine the effects of contextual determinants on both incident dementia and cognition, adjusted for genetic risk, we employed two types of regression models: Cox proportional hazards models for incident dementia and linear regression models with generalized estimating equations (GEE) for cognition. The Cox proportional hazards models estimate the hazard ratio for exposures at baseline on time-to-hospitalization or death due to dementia. The GEE models allow for adjustments for neighborhood-level clustering and repeated measures for MESA participants across five time points. For both incident dementia and late-life cognition, we modeled each contextual exposure individually using two models: **Model 1:** Without genetic risk factors; **Model 2:** Adjusted for *APOEε4* carrier status and PRS.

We also tested for linear and non-linear interactions. For each of the contextual determinants, we fit models stratified by genetic risk group using the following genetic risk categories: **Category 1:** Carriers of the *APOEε4* allele; **Category 2:** Non-carriers of the *APOEε4* allele with PRS in the top quartile; **Category 3:** Non-carriers of the *APOEε4* allele with PRS in the interquartile range (25th to 75th percentile); **Category 4:** Non-carriers of the *APOEε4* allele with PRS in the lowest quartile. We similarly tested for nonlinear interactions by fitting each contextual exposure as a cubic spline with two degrees of freedom and tested interaction with

genetic risk modeled by genetic risk category. ANOVA tests were used to assess the statistical significance of interactions between non-linear exposures and genetic risk category.

As a sensitivity analysis, we used GEE to assess the association between contextual determinants and binary dementia status. This approach avoids relying on time-to-event data, which may inaccurately estimate disease onset, and instead focuses on binary case-control status. It also accounts for neighborhood-level clustering and incorporates repeated measures across five MESA exams for the exposures. We used the GEE models to assess the association of each exposure with dementia status in all MESA participants as well as in the four genetically stratified risk groups.

### *B. Bayesian Kernel Machine Regression*

Bayesian Kernel Machine Regression (BKMR) is an approach for modeling the effects of a mixture of exposures, allowing for nonlinear effects and interactions among mixture components[19]. BKMR also estimates the effects of each of the individual exposures and determines which of the exposures contributes to the mixture's overall effect on the outcome. We used BKMR to jointly model the effects of the seven contextual determinants on cognition. We fit separate BKMR models for the four different levels of genetic risk for AD to determine if the mixture of contextual determinants has different effects depending on genetic risk. All seven contextual determinants were standardized. Standardized distances to alcohol establishments and fast food chains were flipped based on our hypothesis that proximity to these establishments have opposing effects to the other features in the model. Only complete cases were retained for the BKMR analysis.

## **Results**

### *Study Population*

We excluded individuals missing dementia follow-up or genotyping data. The final analytical sample consisted of 5,876 individuals. Participants had a mean age of 62 years at baseline (standard deviation [SD] = 10.1), and 52% were women. 23.8% were *APOEε4* carriers, with 2.4% homozygous for the *e4* allele. 40% of the population was white, 25.8% Black, 22.1% Hispanic and/or Latino, and 12.1% Chinese. After a median follow-up time of 16.8 years, 569 MESA participants developed dementia according to hospital and death certificate ICD codes. Complete baseline demographics and case rates per genetic risk category are presented in **Table 1** and **Supplementary Table 1**, respectively.

### *Linear Models*

#### *A. Cox proportional hazards regression to model incident dementia*

When comparing participants with a one-unit difference in neighborhood SES, adjusting for *APOE* genotype, PRS, age, sex, site, income, education, smoking, marital status, and self-reported race/ethnicity, worse neighborhood SES was associated with a 14% greater hazard of dementia ( $HR_{SES} = 1.14$ ; 95% CI: 1.05-1.25; **Figure 1A**). There was no evidence of association with dementia incidence for ambient pollutants or built environment features (**Table 2**).

Cox proportional hazards models stratified by genetic risk category revealed that neighborhood SES was only associated with increased hazard of dementia among *APOEε4* carriers, where a one-unit difference in SES index was associated with associated with a 25% greater hazard of dementia ( $HR_{SES\_CI} = 1.25$ ; 95% CI: 1.06 – 1.47). While there was some evidence of a dose-response relationship between genetic risk category and SES, as the average hazard ratio of SES increased alongside increasing genetic risk (**Figure 2**), there was no evidence of significant effect modification (p-value > 0.5). The remaining contextual determinants were not associated with incident dementia in the stratified groups (**Table 3**). Among the genetic risk category 3 participants, we observed a high point estimate per IQR difference in NO<sub>2</sub>, but a

substantial degree of uncertainty ( $HR_{NO2\_C3} = 2.08$ ; 95% CI: 0.85 – 5.07). Sensitivity analysis modeling dementia status as a binary trait revealed similar results (**Table 2, Supplementary Table 2**). Across all MESA participants, neighborhood SES was associated with odds of dementia ( $OR_{SES} = 1.1$ ; 95%CI: 1.04 – 1.25). The other contextual determinants were not associated with odds of dementia. When stratified by genetic risk, worse neighborhood SES was only associated with increased odds of dementia among *APOEε4* carriers ( $OR_{SES\_C1} = 1.38$ ; 95% CI: 1.15 – 1.65).

Additionally, all contextual determinants were modeled with a two degree of freedom cubic spline to test for non-linear interactions. The model included interaction terms between the splines and four levels of genetic risk. ANOVA tests of cubic spline ( $df = 2$ ) effect modification models showed that no individual contextual determinants had significantly different effects across genetic risk groups ( $p\text{-value} > 0.5$ ; **Supplementary Figures 3, 5-10**)

#### *B. Generalized estimating equations to model cognition*

After controlling for *APOE* genotype, PRS, age, sex, site, income, education, smoking, marital status, and self-reported race/ethnicity, a one-unit increase in neighborhood SES index, indicating worse SES, was associated with lower CASI score, ( $\beta_{SES} = -0.30$ , 95% CI: -0.49 – -0.11; **Figure 1B**). After adjusting for SES in addition to the covariates, a two-fold greater distance from alcohol establishment was associated with a 0.17 difference in CASI score ( $\beta_{ALC} = 0.17$ , 95% CI: 0.02 – 0.33). There was no evidence of a difference in cognition for other contextual determinants (**Table 2**).

Models stratified by genetic risk category reveal that neighborhood SES was associated with cognition only among those with moderate genetic risk ( $\beta_{SES\_C3} = -0.35$ , 95% CI: -0.69 – -0.01). In addition, distance from alcohol establishments was only associated with cognition among *APOEε4* carriers ( $\beta_{ALC\_C1} = 0.34$ , 95% CI: 0.02 – 0.33). There was no evidence of linear

or nonlinear effect modification between the contextual determinants and genetic risk group on cognition (ANOVA p-value > 0.5).

### C. Bayesian Kernel Machine Regression

We used BKMR to model the joint effects of all seven contextual determinants on cognition. Separate models were fit in four genetic risk categories. 960 individuals were in the *e4* group, 676 individuals were in the non-*e4*, High PRS group, 1368 individuals in the non-*e4*, Intermediate PRS group, and 683 individuals in the non-*e4*, Low PRS group. Across all groups, increasing quantiles of the mixture was associated with worse cognition, although the effects differ in strength depending on the genetic risk group (**Figure 3**). The overall effect of the mixture on cognition is most pronounced among the Category 3 participants with intermediate genetic risk. Posterior inclusion probabilities and visual inspection of the effects of single-exposures on the mixture indicate that these differences are driven by the differing effects of NO<sub>2</sub> across groups (**Figure 4, Supplementary Table 3, Supplementary Figure 4**).

## Discussion

This is the first study in MESA to integrate polygenic risk scores and contextual determinants to evaluate the association between social, chemical, and built environmental context and late-life dementia and cognition. We found that lower neighborhood SES is associated with both higher incident dementia and cognition, even after accounting for individual-level genetic risk and demographics including income, race, age, and sex.

We also observed evidence of gene-environment interaction between neighborhood SES and genetic risk for Alzheimer's disease. When stratified by genetic risk, neighborhood SES was only significantly associated with incident dementia among *APOEε4* carriers. In this group with elevated genetic risk, living in more socioeconomically advantaged areas was associated with lower dementia risk. Associations between lower SES and greater dementia risk are well known

[20,21]. SES affects access to education and other resources that are known to lower dementia risk. Our study is consistent with these results and suggests the strength of the association may vary based on underlying genetics.

Living in socioeconomically disadvantaged neighborhoods has been associated with lower hippocampal volume in addition to dementia incidence [22,23]. Possible explanations include the role of inflammation and oxidative stress, where neighborhood SES contributes to greater risk for diseases that dysregulate inflammatory responses such as cardiovascular disease and hypertension [24–26]. The inflammatory pathway also aligns with findings from genetic studies, as many implicated genes for Alzheimer’s disease are primarily expressed in immune cells [27–29]. Future studies must continue to leverage both genetic and contextual approaches to identify specific mechanisms that explain why contextual determinants are consistently linked to brain aging and dementia outcomes.

In addition to the tests of individual contextual determinants, the BKMR model also suggests that NO<sub>2</sub> may be an important feature with late-life cognition when considering all contextual determinants jointly. Greater exposure to NO<sub>2</sub> was associated with lower values of cognition when genetic risk is moderate and other contextual determinants are at low quantiles. While NO<sub>2</sub> did not have a significant linear association with dementia, the high amount of uncertainty suggest that this study may be underpowered to detect an association, with our model showing that a per-IQR increase in NO<sub>2</sub> levels is associated with a hazard ratio of 1.16 (95% CI: 0.70, 1.94). Larger sample sets may be able to capture an increased hazard dementia associated with more precision. Furthermore, the increased effect of NO<sub>2</sub> on cognition in groups with moderate genetic risk seen in our BKMR results is echoed by our stratified Cox proportional

hazards models, where NO<sub>2</sub> also has the highest point estimate on incident dementia in the group with moderate genetic risk (HR<sub>NO<sub>2</sub>\_C3</sub>: 2.08; 95% CI: 0.85 – 5.07).

There are several limitations to this study. First, as this is an observational study, the results should not be interpreted as establishing causal relationships. Additionally, the use of ICD codes to identify dementia diagnoses, are known to be specific but not sensitive, thus we likely missed some participants with dementia. Additionally, we were unable to differentiate between dementia subtypes. There may be varying effects between contextual determinants on Alzheimer's disease vs. vascular dementia. Our PRS is based on genome-wide studies of Alzheimer's disease. While there may be some overlap in genetic risk factors across dementia subtypes, we are likely missing genetic risk factors for non-Alzheimer's disease dementia. Furthermore, our measures of the built environment and ambient pollutants are based on residential addresses provided by MESA participants, but people often spend significant time in neighborhoods outside of their place of residence, which could affect the accuracy of environmental exposure assessments.

Despite these limitations, our study provides important insights into the complex interactions between genetic risk, neighborhood context, and dementia. The modifiable factors we investigated here related to the built, social, and chemical environment are actionable through policy changes at various levels of governance, offering potential avenues for public health interventions aimed at reducing the burden of dementia. Additionally, by incorporating genetic information, we identified populations where neighborhood context may have stronger influence on dementia risk and late-life cognition, providing support for future studies that aim toward more precise interventions.

## References

- [1] Skaria AP. The economic and societal burden of Alzheimer disease: managed care considerations. *Am J Manag Care* 2022;28:S188–96.  
<https://doi.org/10.37765/ajmc.2022.89236>.
- [2] 2024 Alzheimer’s disease facts and figures. *Alzheimer’s & Dementia* 2024;20:3708–821.  
<https://doi.org/10.1002/alz.13809>.
- [3] Zissimopoulos J, Crimmins E, St. Clair P. The value of delaying Alzheimer’s disease onset. *Forum Health Econ Policy*, vol. 18, De Gruyter; 2015, p. 25–39.  
<https://doi.org/10.1515/fhep-2014-0013>.
- [4] Livingston G, Huntley J, Liu KY, Costafreda SG, Selbæk G, Alladi S, et al. Dementia prevention, intervention, and care: 2024 report of the Lancet standing Commission. *The Lancet* 2024;404:572–628. [https://doi.org/10.1016/S0140-6736\(24\)01296-0](https://doi.org/10.1016/S0140-6736(24)01296-0).
- [5] Gatz M, Reynolds CA, Fratiglioni L, Johansson B, Mortimer JA, Berg S, et al. Role of genes and environments for explaining Alzheimer disease. *Arch Gen Psychiatry* 2006;63:168–74. <https://doi.org/10.1001/archpsyc.63.2.168>.
- [6] Besser L, Galvin JE, Rodriguez D, Seeman T, Kukull W, Rapp SR, et al. Associations between neighborhood built environment and cognition vary by apolipoprotein E genotype: Multi-Ethnic Study of Atherosclerosis. *Health Place* 2019;60:102188.  
<https://doi.org/10.1016/j.healthplace.2019.102188>.
- [7] M. SR, N. BM, Ge L, D. AS, Marco C, A. SA, et al. Fine Particulate Matter and Dementia Incidence in the Adult Changes in Thought Study. *Environ Health Perspect* 2021;129:87001. <https://doi.org/10.1289/EHP9018>.
- [8] Emrani S, Arain HA, DeMarshall C, Nuriel T. APOE4 is associated with cognitive and pathological heterogeneity in patients with Alzheimer’s disease: a systematic review. *Alzheimers Res Ther* 2020;12:141. <https://doi.org/10.1186/s13195-020-00712-4>.
- [9] Blue EE, Horimoto ARVR, Mukherjee S, Wijsman EM, Thornton TA. Local ancestry at APOE modifies Alzheimer’s disease risk in Caribbean Hispanics. *Alzheimers Dement* 2019;15:1524–32. <https://doi.org/10.1016/j.jalz.2019.07.016>.

- [10] Hohman TJ, Cooke-Bailey JN, Reitz C, Jun G, Naj A, Beecham GW, et al. Global and local ancestry in African-Americans: Implications for Alzheimer's disease risk. *Alzheimer's and Dementia* 2016;12:233–43. <https://doi.org/10.1016/j.jalz.2015.02.012>.
- [11] Andrews SJ, Renton AE, Fulton-Howard B, Podlesny-Drabiniok A, Marcora E, Goate AM. The complex genetic architecture of Alzheimer's disease: novel insights and future directions. *EBioMedicine* 2023;90. <https://doi.org/10.1016/j.ebiom.2023.104511>.
- [12] Choi SW, Mak TS-H, O'Reilly PF. Tutorial: a guide to performing polygenic risk score analyses. *Nat Protoc* 2020;15:2759–72. <https://doi.org/10.1038/s41596-020-0353-1>.
- [13] Bild DE, Bluemke DA, Burke GL, Detrano R, Diez Roux A V, Folsom AR, et al. Multi-Ethnic Study of Atherosclerosis: Objectives and Design. *Am J Epidemiol* 2002;156:871–81. <https://doi.org/10.1093/aje/kwf113>.
- [14] Fujiyoshi A, Jacobs Jr DR, Alonso A, Luchsinger JA, Rapp SR, Duprez DA. Validity of Death Certificate and Hospital Discharge ICD Codes for Dementia Diagnosis: The Multi-Ethnic Study of Atherosclerosis. *Alzheimer Dis Assoc Disord* 2017;31:168–72. <https://doi.org/10.1097/WAD.0000000000000164>.
- [15] Teng EL, Hasegawa K, Homma A, Imai Y, Larson E, Graves A, et al. The Cognitive Abilities Screening Instrument (CASI): A Practical Test for Cross-Cultural Epidemiological Studies of Dementia. *Int Psychogeriatr* 1994;6:45–58. <https://doi.org/10.1017/S1041610294001602>.
- [16] Diez Roux A V, Mujahid MS, Hirsch JA, Moore K, Moore L V. The Impact of Neighborhoods on CV Risk. *Glob Heart* 2016;11:353–63. <https://doi.org/10.1016/j.gheart.2016.08.002>.
- [17] Kaufman JD, Adar SD, Allen RW, Barr RG, Budoff MJ, Burke GL, et al. Prospective Study of Particulate Air Pollution Exposures, Subclinical Atherosclerosis, and Clinical Cardiovascular Disease: The Multi-Ethnic Study of Atherosclerosis and Air Pollution (MESA Air). *Am J Epidemiol* 2012;176:825–37. <https://doi.org/10.1093/aje/kws169>.
- [18] Kunkle BW, Grenier-Boley B, Sims R, Bis JC, Damotte V, Naj AC, et al. Genetic meta-analysis of diagnosed Alzheimer's disease identifies new risk loci and implicates A $\beta$ , tau,

- immunity and lipid processing. *Nat Genet* 2019;51:414–30.  
<https://doi.org/10.1038/s41588-019-0358-2>.
- [19] Bobb JF, Claus Henn B, Valeri L, Coull BA. Statistical software for analyzing the health effects of multiple concurrent exposures via Bayesian kernel machine regression. *Environmental Health* 2018;17:1–10. <https://doi.org/10.1515/fhep-2014-0013>.
- [20] Mortimer JA, Graves AB. Education and other socioeconomic determinants of dementia and Alzheimer’s disease. *NEUROLOGY-MINNEAPOLIS-* 1993;43:39.
- [21] Li R, Li R, Xie J, Chen J, Liu S, Pan A, et al. Associations of socioeconomic status and healthy lifestyle with incident early-onset and late-onset dementia: a prospective cohort study. *Lancet Healthy Longev* 2023;4:e693–702. [https://doi.org/10.1016/S2666-7568\(23\)00211-8](https://doi.org/10.1016/S2666-7568(23)00211-8).
- [22] Hunt JF V, Buckingham W, Kim AJ, Oh J, Vogt NM, Jonaitis EM, et al. Association of Neighborhood-Level Disadvantage With Cerebral and Hippocampal Volume. *JAMA Neurol* 2020;77:451–60. <https://doi.org/10.1001/jamaneurol.2019.4501>.
- [23] Dintica CS, Bahorik A, Xia F, Kind A, Yaffe K. Dementia Risk and Disadvantaged Neighborhoods. *JAMA Neurol* 2023;80:903–9. <https://doi.org/10.1001/jamaneurol.2023.2120>.
- [24] Yaffe K, Falvey C, Harris TB, Newman A, Satterfield S, Koster A, et al. Effect of socioeconomic disparities on incidence of dementia among biracial older adults: prospective study. *Bmj* 2013;347.
- [25] Faraco G, Iadecola C. Hypertension: a harbinger of stroke and dementia. *Hypertension* 2013;62:810–7. <https://doi.org/10.1161/HYPERTENSIONAHA.113.01063>.
- [26] Luchsinger JA, Mayeux R. Cardiovascular risk factors and Alzheimer’s disease. *Curr Atheroscler Rep* 2004;6:261–6. <https://doi.org/10.1007/s11883-004-0056-z>.
- [27] Bertram L, Lange C, Mullin K, Parkinson M, Hsiao M, Hogan MF, et al. Genome-wide association analysis reveals putative Alzheimer’s disease susceptibility loci in addition to APOE. *The American Journal of Human Genetics* 2008;83:623–32. <https://doi.org/10.1016/j.ajhg.2008.10.008>.

- [28] Jonsson T, Stefansson H, Steinberg S, Jonsdottir I, Jonsson P V, Snaedal J, et al. Variant of TREM2 associated with the risk of Alzheimer's disease. *New England Journal of Medicine* 2013;368:107–16. <https://doi.org/10.1056/NEJMoa1211103>.
- [29] Romero-Molina C, Garretti F, Andrews SJ, Marcora E, Goate AM. Microglial efferocytosis: diving into the Alzheimer's disease gene pool. *Neuron* 2022;110:3513–33. <https://doi.org/10.1016/j.neuron.2022.10.015>.

## Tables

**Table 1. Baseline Characteristics of MESA Participants.**

	<b>Censored</b> (N=5352)	<b>Affected</b> (N=518)	<b>Overall</b> (N=5870)
<b>Age at baseline (years)</b>			
Mean (SD)	61.0 (9.82)	71.9 (7.75)	62.0 (10.1)
Median [Min, Max]	61.0 [44.0, 84.0]	74.0 [45.0, 84.0]	62.0 [44.0, 84.0]
<b>Site of clinic visit</b>			
WFU	829 (15.5%)	104 (20.1%)	933 (15.9%)
COL	801 (15.0%)	88 (17.0%)	889 (15.1%)
JHU	854 (16.0%)	72 (13.9%)	926 (15.8%)
UMN	861 (16.1%)	97 (18.7%)	958 (16.3%)
NWU	909 (17.0%)	80 (15.4%)	989 (16.8%)
UCLA	1098 (20.5%)	77 (14.9%)	1175 (20.0%)
<b>Sex</b>			
Female	2794 (52.2%)	256 (49.4%)	3050 (52.0%)
Male	2558 (47.8%)	262 (50.6%)	2820 (48.0%)
<b>Level of education</b>			
< High school degree	873 (16.3%)	114 (22.0%)	987 (16.8%)
High school degree	946 (17.7%)	125 (24.1%)	1071 (18.2%)
Some college	1526 (28.5%)	131 (25.3%)	1657 (28.2%)
Bachelor's degree or higher	1991 (37.2%)	146 (28.2%)	2137 (36.4%)
Missing	16 (0.3%)	2 (0.4%)	18 (0.3%)
<b>Marital status</b>			
Divorced/ separated/ never married	1981 (37.0%)	255 (49.2%)	2236 (38.1%)
Married/Living as Married	3356 (62.7%)	261 (50.4%)	3617 (61.6%)
Missing	15 (0.3%)	2 (0.4%)	17 (0.3%)
<b>Household income</b>			
< \$30,000	1776 (33.2%)	252 (48.6%)	2028 (34.5%)
> \$30,000	3398 (63.5%)	234 (45.2%)	3632 (61.9%)
Missing	178 (3.3%)	32 (6.2%)	210 (3.6%)
<b>Race/ethnicity</b>			
Black	1380 (25.8%)	133 (25.7%)	1513 (25.8%)
Chinese	673 (12.6%)	36 (6.9%)	709 (12.1%)
Hispanic/Latino	1191 (22.3%)	107 (20.7%)	1298 (22.1%)
White	2108 (39.4%)	242 (46.7%)	2350 (40.0%)
<b>Pack-years of cigarette smoking</b>			
Mean (SD)	11.0 (20.7)	12.7 (22.3)	11.1 (20.8)
Median [Min, Max]	0 [0, 265]	0 [0, 140]	0 [0, 265]
Missing	73 (1.4%)	5 (1.0%)	78 (1.3%)
<b>APOE genotype</b>			
<i>e2/e2</i>	42 (0.8%)	2 (0.4%)	44 (0.7%)
<i>e2/e3</i>	625 (11.7%)	54 (10.4%)	679 (11.6%)
<i>e2/e4</i>	133 (2.5%)	14 (2.7%)	147 (2.5%)
<i>e3/e3</i>	3214 (60.1%)	273 (52.7%)	3487 (59.4%)
<i>e3/e4</i>	1116 (20.9%)	142 (27.4%)	1258 (21.4%)
<i>e4/e4</i>	119 (2.2%)	20 (3.9%)	139 (2.4%)
Missing	103 (1.9%)	13 (2.5%)	116 (2.0%)
<b>Cognitive Assessment Score at Exam 5</b>			
Mean (SD)	88.0 (8.61)	80.6 (12.1)	87.5 (9.03)
Median [Min, Max]	89.5 [23.2, 100]	83.0 [28.1, 98.4]	89.0 [23.2, 100]
Missing	1345 (25.1%)	259 (50.0%)	1604 (27.3%)

**Table 2. Effects of contextual determinants estimated by linear models.**

Table 2 Legend. The effect size of association of ambient pollutants and built environment features with dementia and cognitive decline after controlling for neighborhood SES, *APOEε4*, PRS, carrier status, site, sex, race/ethnicity, education, household income, smoking, and age.

*\*Estimates represent per IQR difference.*

*\*\*Estimates represent per two-fold difference in Euclidean distance to the nearest.*

*Significant associations are in boldface.*

	<b>Dementia Hazard Ratio (95% CI)</b>	<b>Dementia Marginal Odds Ratio (95% CI)</b>	<b>Cognition Estimate (95% CI)</b>
<b>SES</b>	<b>1.14 (1.05, 1.25)</b>	<b>1.14 (1.04, 1.25)</b>	<b>-0.30 (-0.49, -0.11)</b>
<b>PM2.5*</b>	1.04 (0.86, 1.27)	1.00 (0.98, 1.01)	0.16 (-0.30, 0.61)
<b>NO2*</b>	1.16 (0.70, 1.94)	1.00 (0.96, 1.03)	0.17 (-0.98, 1.32)
<b>Favorable Food Store**</b>	0.97 (0.88, 1.06)	1.00 (0.99, 1.00)	0.15 (-0.03, 0.33)
<b>Fast Food Chain**</b>	1.00 (0.91,1.09)	1.00 (0.99, 1.04)	0.09 (-0.07, 0.25)
<b>Alcohol Establishment**</b>	0.98 (0.89,1.05)	1.00 (0.99, 1.01)	<b>0.17 (0.02, 0.33)</b>
<b>Physical Activity Facility**</b>	0.98 (0.90,1.07)	1.00 (0.99, 1.00)	-0.02 (-0.17, 0.14)

**Table 3. Effect sizes and confidence intervals of contextual determinants on incident dementia in genetically stratified groups.**

Table 3 Legend. This table shows the hazard ratio of contextual determinants on incident dementia in genetically stratified groups. Bolded cells indicate statistically significant associations with dementia.

*\* models not adjusted for SES*

	<b>APOEε4 carriers</b>	<b>Non-ε4 carriers (High PRS)</b>	<b>Non-ε4 carriers (Mid. PRS)</b>	<b>Non-ε4 carriers (Low PRS)</b>
<b>SES*</b>	<b>1.25 (1.06,1.47)</b>	1.19 (0.96, 1.47)	1.10 (0.94, 1.28)	1.05 (0.85, 1.32)
<b>PM2.5</b>	1.35 (0.93, 1.95)	0.65 (0.43, 0.98)	1.37 (0.94, 1.99)	1.06 (0.61, 1.84)
<b>NO2</b>	1.70 (0.66, 4.37)	0.68 (0.22, 2.14)	2.08 (0.85, 5.07)	0.64 (0.16, 2.48)
<b>Favorable Food Store</b>	0.98 (0.83, 1.16)	0.93 (0.75, 1.16)	0.99 (0.84, 1.16)	1.17 (0.88, 1.56)
<b>Fast Food Chain</b>	1.01 (0.87, 1.17)	0.98 (0.78, 1.23)	1.01 (0.86, 1.18)	0.99 (0.79, 1.26)
<b>Alcohol Establishment</b>	0.93 (0.81, 1.07)	1.03 (0.84, 1.27)	0.96 (0.83, 1.11)	0.98 (0.79, 1.23)
<b>Physical Activity Facility</b>	1.07 (0.91, 1.26)	0.94 (0.79, 1.14)	0.96 (0.84, 1.10)	1.01 (0.80, 1.29)

**Table 4. Effect sizes and confidence intervals of contextual determinants on late-life cognition in genetically stratified groups.**

Table 4 Legend. This table shows the estimate of contextual determinants on CASI score in genetically stratified groups. Bolded cells indicate statistically significant associations with dementia.

*\* models not adjusted for SES*

	<b>APOEε4 carriers</b>	<b>Non-ε4 carriers (High PRS)</b>	<b>Non-ε4 carriers (Mid. PRS)</b>	<b>Non-ε4 carriers (Low PRS)</b>
<b>SES*</b>	-0.23 (-0.57, 0.11)	-0.29 (-0.72, 0.13)	<b>-0.35 (-0.69, -0.01)</b>	-0.17 (-0.58, 0.23)
<b>PM2.5</b>	-0.09 (-1.01, 0.81)	-0.26 (-1.31, 0.78)	0.30 (-0.50, 1.11)	0.29 (-0.53, 1.10)
<b>NO2</b>	1.65 (-0.53, 3.83)	-1.32 (-4.27, 1.62)	-0.29 (-2.16, 1.57)	-0.04 (-2.40, 2.32)
<b>Favorable Food Store</b>	0.25 (-0.09, 0.61)	0.04 (-0.40, 0.48)	0.12 (-0.18, 0.42)	0.15 (-0.20, 0.51)
<b>Fast Food Chain</b>	0.04 (-0.25, 0.32)	0.07 (-0.33, 0.47)	0.19 (-0.07, 0.47)	-0.05 (-0.40, 0.29)
<b>Alcohol Establishment</b>	<b>0.34 (0.05, 0.62)</b>	-0.01 (-0.40, 0.38)	0.17 (-0.08, 0.43)	0.22 (-0.14, 0.58)
<b>Physical Activity Facility</b>	0.20 (-0.11, 0.51)	-0.28 (-0.63, 0.07)	-0.08 (-0.31, 0.15)	0.11 (-0.25, 0.47)

## Figures

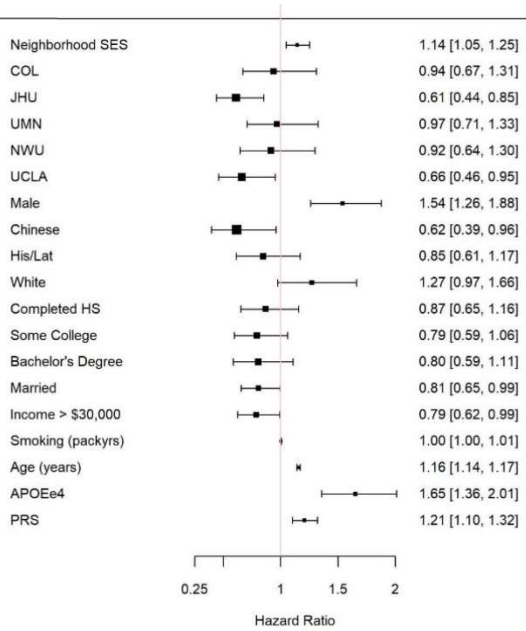
### Figure 1. Association of Neighborhood with incident dementia and CASI

Figure 1 Legend.

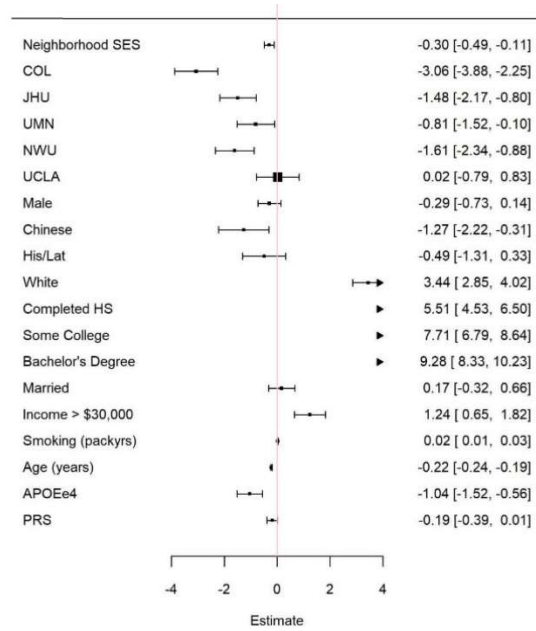
**A)** The Cox proportional hazards regression model estimated hazard ratios (HR) for the risk of developing dementia. After controlling for genetic risk and other covariates, we find that neighborhood SES is associated with dementia risk (HR = 1.14, 95% CI: 1.05, 1.25). The plot depicts hazard ratios with 95% confidence intervals, with values above 1 indicating increased risk and values below 1 indicating decreased risk. Higher values of neighborhood indicate worse SES.

**B)** The generalized estimating equation linear model estimated the effect of SES on late-life global cognition, as measured by the CASI at MESA Exam 5. After controlling for genetic risk and other covariates, neighborhood SES is associated with cognition score, where on average, a one unit increase in SES index (indicating worse SES) is associated with lower performance on the CASI. ( $\beta = -0.30$ , 95% CI: -0.49, -0.11).

**A**

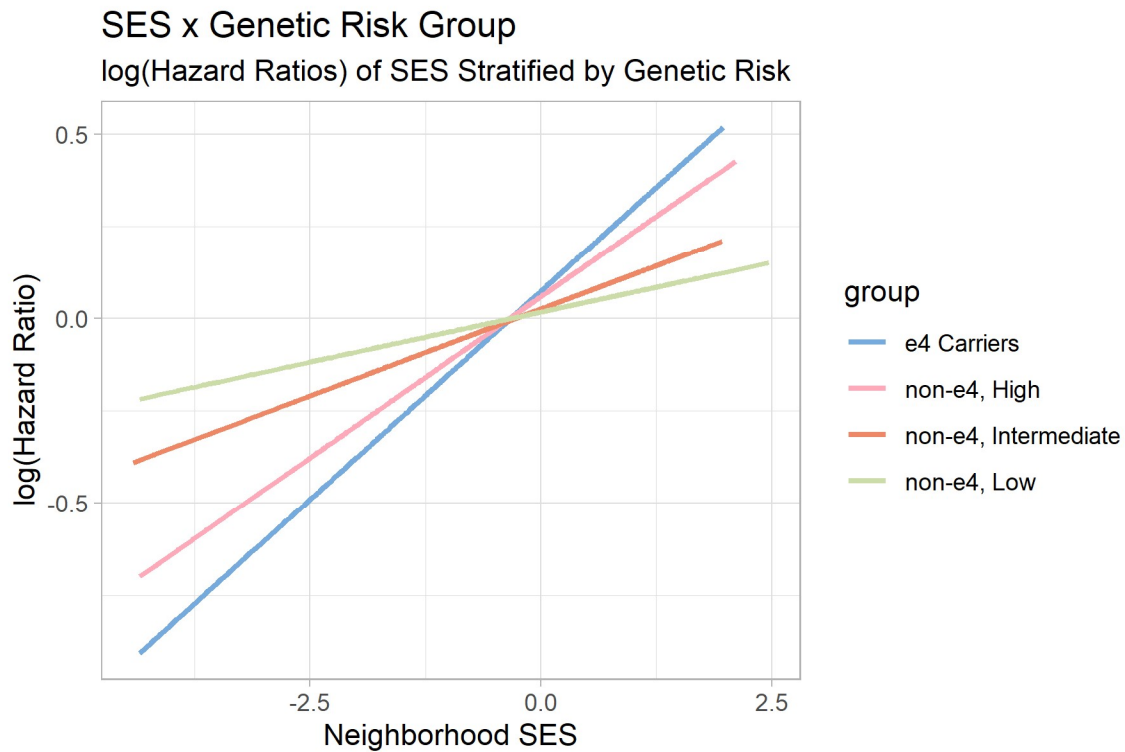


**B**



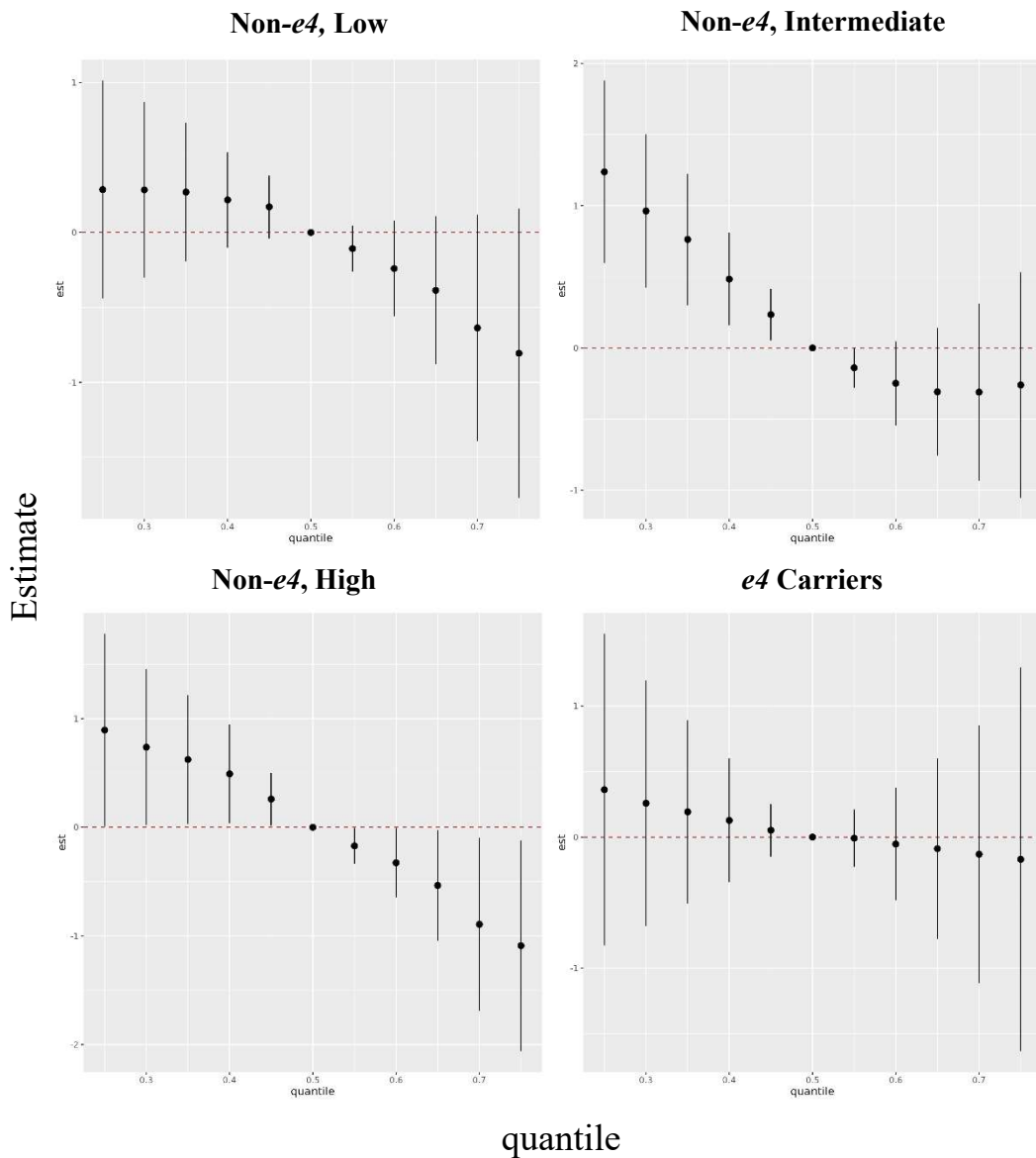
## Figure 2. Linear Interaction between Neighborhood SES and Genetic Risk on Incident Dementia

Figure 2 Legend. The MESA participants were stratified by genetic risk into four groups. The *e4* carriers are considered the highest genetic risk group. The non-*e4* carriers are then split into three groups, with the high and low groups being the top and bottom quartile of polygenic risk scores among the non-*e4* carriers. Cox proportional hazards models were fit in each genetic risk group. We observe a dose-response relationship here where SES has the highest log-HR in the highest genetic risk group.



**Figure 3. Estimated joint effects of contextual determinants on cognition by BKMR.**

Figure 3 Legend. Instead of analyzing the effects of each contextual exposure individually, Bayesian Kernel Machine Regression allows for the assessment of the effect of a mixture of exposures on late-life cognition (CASI score). We used BKMR to identify the joint effects of all seven contextual determinants among groups stratified by genetic risk. For each quantile, the estimate is showing the predictor-response function or the effect of the total mixture when compared to all exposures being set at their median (50% percentile).



## Supplementary Materials

Paper 1 Supplement: The power of representation: Statistical analysis of diversity in US Alzheimer's disease genetics data

**Supplementary Table 1. Study-specific sample sizes for array datasets**

N: Total sample size, pCases: Proportion of samples that are cases affected by AD as defined by the study.

<b>Dataset</b>	<b>Abbreviation</b>	<b>N</b>	<b>pCases</b>
<b>African American Alzheimer's Disease Genetics</b>	AAG	785	0.79
<b>Adult Changes in Thought</b>	ACT	2462	0.21
<b>Alzheimer Disease Center Dataset 1</b>	ADC1	2702	0.73
<b>Alzheimer Disease Center Dataset 2</b>	ADC2	926	0.74
<b>Alzheimer Disease Center Dataset 3</b>	ADC3	1510	0.55
<b>Alzheimer Disease Center Dataset 4</b>	ADC4	1054	0.41
<b>Alzheimer Disease Center Dataset 5</b>	ADC5	1223	0.33
<b>Alzheimer Disease Center Dataset 6</b>	ADC6	1333	0.42
<b>Alzheimer Disease Center Dataset 7</b>	ADC7	1462	0.36
<b>Alzheimer's Disease Neuroimaging Initiative</b>	ADNI	441	0.61
<b>Atherosclerosis Risk in Communities</b>	ARIC	11561	0.19
<b>Biomarkers for Older Controls at Risk for Dementia</b>	BIOCARD	118	0.05
<b>Chicago Health and Aging Project</b>	CHAP	171	0.16
<b>Cardiovascular Health Study</b>	CHS	3359	0.21
<b>Einstein Aging Study</b>	EAS	150	0.06
<b>Estudio Familiar de Influencia Genetica en Alzheimer</b>	EFIGA	683	0.65
<b>Framingham Heart Study</b>	FHS	4240	0.08
<b>Genetic Alzheimer's Disease Associations</b>	GenADA	1378	0.48
<b>Genetic and Environmental Risk Factors for Alzheimer's Disease Among African Americans</b>	GenerAAtions	446	0.54
<b>Indianapolis-Ibadan African Americans</b>	INDY AA	1175	0.15
<b>Knight Alzheimer's Disease Research Center</b>	KnightADRC	4445	0.44
<b>Harmonized Diagnostic Assessment of Dementia for the Longitudinal Aging Study of India</b>	LASI-DAD	960	NA
<b>Mayo Clinic</b>	MAYO	2209	0.57
<b>Multi Institutional Research of Alzheimer Genetic Epidemiology</b>	MIRAGE	1769	0.59
<b>National Centralized Repository for Alzheimer's Disease and Related Dementias Family</b>	NCRAD Families GWAS	688	0.37
<b>National Institute on Aging Late Onset of Alzheimer's Disease Family</b>	NIALOAD	4419	0.48
<b>Oregon Health and Science University</b>	OHSU	285	0.46
<b>Puerto Rican Alzheimer's Disease Initiative</b>	PRADI	674	0.42
<b>Religious Orders Study/Memory and Aging Project</b>	ROSMAP	1202	0.27
<b>Religious Orders Study/Memory and Aging Project 2</b>	ROSMAP2	276	0.21

<b>Texas Alzheimer's Research and Care Consortium</b>	TARCC	2163	0.61
<b>Translational Genomics Research Institute</b>	TGENII	1033	0.63
<b>Miami, Vanderbilt, and Medical School of Mount Sinai</b>	UMVUMSSM	2645	0.48
<b>University of Pittsburgh</b>	UPITT	2277	0.63
<b>Washington University</b>	WASHU	775	0.36
<b>Washington Heights and Inwood Community Aging Project</b>	WHICAP	2801	0.27

**Supplementary Table 2. Study-specific sample sizes for sequencing datasets**

N: total sample size, pCases: proportion of sample that are “cases” affected by AD as defined by the study, Data Type: whole-exome sequencing (WES) or whole-genome sequencing (WGS)

Dataset	N	pCases	WGS or WES
<b>ADGC-TARCC</b>	1017	0.13	WGS
<b>ADNI-WGS-1</b>	809	0.3	WGS
<b>ADSP Discovery-WGS</b>	578	0.65	WGS
<b>ADSP Extension</b>	3399	0.42	WGS
<b>ADSP-FUS1</b>	8159	0.37	WGS
<b>ADSP-FUS2</b>	12506	0.27	WGS
<b>AMP-AD</b>	1318	0.56	WGS
<b>Cache County</b>	207	0	WGS
<b>EOAD1</b>	3132	0.45	WGS
<b>FASe Families WGS</b>	91	0.81	WGS
<b>Knight ADRC WGS</b>	77	0.7	WGS
<b>LASI-DAD</b>	2686	0.07	WGS
<b>NACC-Genentech</b>	137	0.94	WGS
<b>PSP UCLA</b>	408	0.69	WGS
<b>PSP-CurePSP-Tau</b>	869	1	WGS
<b>PSP-NIH-CurePSP-Tau</b>	612	1	WGS
<b>UPitt-Kambohl</b>	209	1	WGS
<b>ADGC AA</b>	3226	0.38	WES
<b>ADSP Discovery</b>	10304	0.43	WES
<b>Brkanac Families</b>	75	1	WES
<b>CBD</b>	335	1	WES
<b>Columbia WHICAP</b>	3814	0.21	WES
<b>FASe Families</b>	1110	0.65	WES
<b>Knight ADRC</b>	650	0.39	WES
<b>Miami Families</b>	108	0.8	WES
<b>PSP</b>	550	1	WES

### Supplementary Table 3. Case Rates for Power Simulations

Case rates used for power simulations within each population group. Case rates are based on available data. Power simulations were conducted separately for array, whole genome sequencing (WGS) data that is currently available, and WGS data that is projected to be available through the ADSP within the next five years.

	<b>Array</b>	<b>WGS (Current)</b>	<b>WGS (Future)</b>
<b>Asian</b>	0.40	0.07	0.30
<b>Black</b>	0.33	0.33	0.30
<b>Hispanic</b>	0.50	0.30	0.32
<b>White</b>	0.33	0.54	0.35

**Supplementary Table 4. Age-adjusted mortality rates for AD per racial/ethnic group**

<b>Race/ Ethnicity</b>	<b>Crude Mortality Rate per 100,000</b>	<b>Age-adjusted Mortality Rate per 100,000</b>
<b>AI/AN</b>	117.7	151.8
<b>Asian</b>	113.9	125.6
<b>Black</b>	189.2	223.9
<b>Hispanic</b>	191.3	213.7
<b>White</b>	250.8	254.2

*AI/AN: American Indian or Alaskan Native.*

**Supplementary Table 5. Minimum sample size to detect significant loci**

This table shows the minimum sample size needed for 80% power to detect genome-wide significant ( $\alpha = 5e-08$ ) loci given a range of allele frequencies and effect sizes. These sample sizes assume a case rate of 0.4.

		<b>OR</b>				
		<b>1.05</b>	<b>1.1</b>	<b>1.2</b>	<b>1.5</b>	<b>2.0</b>
<b>MAF</b>	<b>0.005</b>	6928996	1804832	488457	97607	33414
	<b>0.006</b>	5779929	1505489	407815	81414	27877
	<b>0.007</b>	4970833	1292378	350020	69852	23923
	<b>0.008</b>	4349335	1131459	306619	61179	20957
	<b>0.009</b>	3864080	1006855	272813	54436	18651
	<b>0.01</b>	3481157	907001	245797	49042	16805
	<b>0.02</b>	1757826	458209	124206	24802	8503
	<b>0.03</b>	1184131	308698	83698	16724	5737
	<b>0.04</b>	897796	234009	63457	12688	4355
	<b>0.05</b>	725755	189201	51325	10269	3527
	<b>0.06</b>	611240	159444	43246	8657	2976
	<b>0.07</b>	529583	138158	37485	7509	2583
	<b>0.08</b>	468495	122230	33179	6649	2289
	<b>0.09</b>	421100	109896	29830	5982	2061
	<b>0.1</b>	383243	100005	27157	5450	1879
	<b>0.15</b>	270692	70687	19218	3868	1339
	<b>0.2</b>	215835	56396	15351	3100	1078
	<b>0.25</b>	184301	48185	13132	2660	928
	<b>0.3</b>	164657	43077	11753	2388	837
	<b>0.35</b>	152097	39815	10876	2216	780
<b>0.4</b>	144277	37789	10334	2112	746	
<b>0.45</b>	139984	36691	10045	2059	729	
<b>0.5</b>	138683	36371	9969	2049	728	

### Supplementary Table 6. AD loci global allele frequencies

This table adapted from Kamboh *et al.* (2021, PMID: PMC913039) shows all genome-wide significant loci associated with AD (or proxy-AD) through 2021. Minor allele frequencies are provided for the European (EUR) reference in gnomAD (v4). Relative frequencies (Rel Freq) in the African American (AFA), Native American (AMR), East Asian (EAS), and South Asian (SAS) populations are provided, defined as relative minor allele frequency compared to EUR MAF.

\*OR reported is primarily from NHW GWAS. Exceptions are noted as follows:

Blue: from GWAS among Japanese population

Green: from cross-population GWAS

Purple: from GWAS among African American population

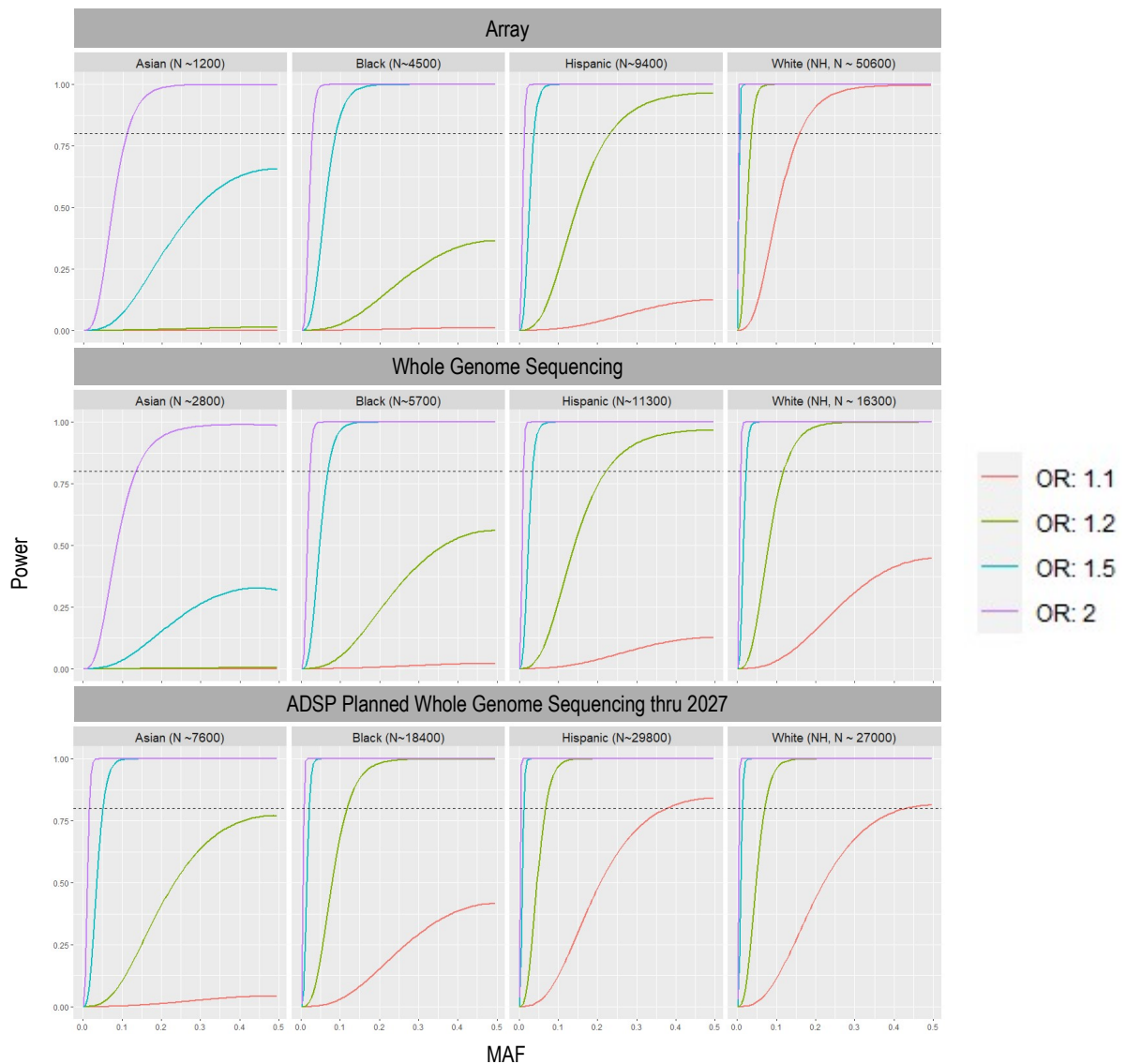
Chr	Position (bp) <sup>a</sup>	Lead SNP	Major/ minor	OR*	EUR MAF	AFA Rel Freq	AMR Rel Freq	EAS Rel Freq	SAS Rel Freq
1	1,049,997	rs113020870	C/T	1.07	0.005	0.17	0.412	0	0.105
1	161,185,602	rs4575098 (3'UTR)	G/A	1.02	0.235	0.252	1.07	1.422	0.508
1	109,345,810	rs141749679 (p.Lys165Glu)	T/C	1.38	0.003	0.135	0.468	0	0.171
1	207,518,704	rs6656401	G/A	1.18	0.188	0.18	0.681	0.184	0.414
2	9,558,882	rs72777026	A/G	1.06	0.143	2.435	1.028	1.102	1.135
2	37,304,796	rs17020490	T/C	1.06	0.146	1.152	2.584	3.787	2.431
2	65,381,229	rs268134	A/G	1.06	0.251	1.057	0.887	0.34	0.645
2	105,749,599	rs143080277	T/C	1.45	0.005	0.177	0.265	0.035	0.115
2	127,135,234	rs6733839	C/T	1.2	0.385	1.029	1.142	1.035	1.036
2	202,878,716	rs139643391 (3'UTR)	TC/T	1.06	0.125	0.262	0.768	0.119	0.246
2	233,117,202	rs10933431	C/G	1.1	0.229	2.766	1.546	1.458	1.64
3	57,192,122	rs184384746	C/T	1.21	0.002	0.117	0.495	0	0.56
3	155,069,722	rs16824536	G/A	1.09	0.053	4.807	0.947	0.372	1.468
	155,084,189	rs61762319 (Met8Val)	A/G	1.16	0.027	0.162	0.42	0.002	0.401
4	993,555	rs3822030	T/G	1.05	0.429	1.79	1.042	1.225	0.749
4	11,025,995	rs4351014	T/C	1.08	0.269	1.388	1.133	0.782	0.815
4	40,197,226	rs2245466 (5'UTR)	C/G	1.05	0.323	0.768	1.071	1.655	0.834
4	76,217,307	rs920608	A/C	1.54	0.009	9.911	2.167	4.822	2.167
5	14,724,304	rs112403360	T/A	1.09	0.078	1.347	0.634	0.387	1.09
5	86,927,378	rs62374257	T/C	1.07	0.226	0.194	1.048	1.205	0.969
5	88,927,603	rs190982	A/G	1.09	0.391	0.26	0.621	0.354	0.954
5	140,335,105	rs2074612	C/T	1.08	0.463	0.561	0.962	1.036	1.177
5	151,052,827	rs871269	C/T	1.02	0.329	1.204	1.551	1.707	1.189

5	157,099,320	rs6891966	G/A	1.02	0.23	0.764	1.111	0.31	0.719
5	180,201,150	rs113706587	G/A	1.09	0.11	0.247	0.857	0.04	1.178
6	27,915,491	rs1497525	C/A	1.34	0.037	5.768	2.215	1.654	1.509
6	32,592,048	rs34855541	A/G	1.11	0.16	0.354	1.389	0.957	0.556
6	41,161,514	rs75932628 (p.Arg47His)	C/T	2.01	0.003	0.146	1.321	0.016	0.898
6	47,520,026	rs10948363	A/G	1.1	0.271	0.67	0.793	0.51	0.758
6	114,291,731	rs785129	C/T	1.04	0.337	0.692	1.238	1.873	1.107
7	7,817,263	rs6943429	C/T	1.05	0.415	1.166	1.226	0.975	1.016
7	8,204,382	rs10952097	C/T	1.07	0.097	5.638	1.275	1.026	0.895
7	12,229,132	rs5011436	A/C	1.02	0.416	1.598	1.307	1.547	1.474
7	28,129,127– 28,129,134	rs1160871	GTCTT/G	1.05	0.219	3.216	1.959	4.43	2.334
7	37,801,932	rs2718058	A/G	1.08	0.365	1.442	0.99	0.541	0.435
7	54,873,635	rs76928645	C/T	1.08	0.107	0.169	0.656	0.004	0.223
7	100,374,211	rs1859788 (p.Gly78Arg)	G/A	1.02	0.316	0.422	1.672	1.883	0.903
7	143,402,040	rs10808026	C/A	1.11	0.203	0.882	0.951	0.735	1.371
7	146,252,937	rs114360492	C/T	1.19	0	106.892	13.199	0.439	0
8	11,844,613	rs1065712 (3'UTR)	G/C	1.09	0.06	0.137	0.339	0.006	0.264
8	27,362,470	rs73223431	C/T	1.1	0.357	0.67	0.782	0.722	1.416
8	27,610,169	rs9331896	T/C	1.14	0.403	1.392	0.898	0.557	0.749
8	144,103,704	rs34173062	G/A	1.13	0.08	0.175	0.608	0.007	0.536
9	104,903,697	rs1800978	C/G	1.06	0.125	0.18	1.448	1.544	1.803
10	11,678,309	rs7920721	A/G	1.08	0.379	0.404	0.993	0.63	0.886
10	60,025,170	rs7068231	G/T	1.05	0.403	0.612	1.276	1.759	1.202
10	59,886,075	rs1171814	G/T	1.05	0.483	0.497	0.994	1.299	1.018
10	80,520,381	rs1878036	T/G	1.07	0.205	0.248	0.605	0.098	0.354
10	96,266,650	rs6584063	A/G	1.12	0.042	1.352	0.695	0.329	0.137
10	122,413,396	rs7908662	A/G	1.04	0.482	0.782	1.155	1.294	1.128
11	47,358,789	rs3740688	T/G	1.09	0.454	0.586	0.988	0.799	0.772
11	60,169,453	rs7933202	A/C	1.12	0.395	0.204	0.708	0.174	1.155
11	86,156,833	rs10792832	G/A	1.15	0.364	0.385	0.998	1.093	1.064
11	121,564,878	rs11218343	T/C	1.25	0.039	2.057	1.764	7.675	1.857
12	113,281,983	rs6489896	T/C	1.08	0.066	2.416	1.297	3.025	1.183
14	52,924,962	rs17125924	A/G	1.12	0.093	0.711	1.094	2.452	1.265
14	92,460,608	rs10498633	G/T	1.1	0.227	0.595	0.852	0.41	0.789
14	105,761,758	rs7157106	G/A	1.05	0.301	2.736	1.915	2.961	1.607
	106,665,591	rs10131280	G/A	1.06	0.131	1.62	1.244	1.33	1.197
15	50,709,337	rs59685680	T/G	1.09	0.204	0.314	0.868	2.259	0.635
15	58,753,575	rs593742	A/G	1.08	0.299	1.003	1.418	2.616	1.71
15	63,277,703	rs117618017 (p.Thr27Ile)	C/T	1.09	0.137	0.158	0.579	0.002	0.183
15	64,131,307	rs3848143	A/G	1.05	0.207	1.896	1.614	0.928	2.749
15	64,433,291	rs74615166	T/C	1.31	0.022	0.18	0.724	0.009	1.133
15	78,936,857	rs12592898	G/A	1.06	0.13	1.785	0.737	0.476	1.144
15	97,449,455	rs570487962	Rare variants	10	0	40.918	6.667	0	0
16	19,796,841	rs7185636	T/C	1.09	0.091	3.046	0.34	0.22	0

16	30,010,081	rs1140239	C/T	1.06	0.397	0.43	0.761	0.926	1.107
16	31,111,250	rs889555	C/T	1.05	0.283	1.455	1.312	0.339	2.329
16	70,660,097	rs4985556 (p.Tyr213Ter)	C/A	1.09	0.118	0.163	0.41	0.411	0.644
16	79,321,960	rs62039712	G/A	1.16	0.117	0.197	0.737	0.008	0.138
16	79,574,511	rs450674	T/C	1.04	0.377	0.674	0.762	0.098	0.755
16	81,908,423	rs72824905 (p.Pro522Arg)	C/G	1.47	0.008	0.142	0.439	0	0.073
16	86,420,604	rs16941239	T/A	1.13	0.025	7.831	2.399	5.724	1.669
16	90,103,687	rs56407236	G/A	1.11	0.069	1.007	1.062	0.369	1.055
17	1,728,047	rs35048651 (5'UTR)	TGAG/T	1.06	0.218	0.572	1.173	0.791	1.755
17	5,233,752	rs7225151	G/A	1.1	0.121	1.913	0.782	0.397	0.334
17	18,156,140	rs2242595	G/A	1.06	0.119	1.23	2.794	4.159	1.299
17	44,364,976	rs708382	T/C	1.02	0.392	0.99	0.966	1.108	0.896
17	49,219,935	rs616338 (p.Ser209Phe)	C/T	1.43	0.011	0.159	0.215	0	0.009
17	58,320,645	rs2526380	G/C	1.03	0.433	1.188	0.749	0.711	0.878
	58,331,728	rs2632516	G/C	1.09	0.444	1.373	0.964	1.142	1.145
17	63,460,787	rs138190086	G/A	1.25	0.017	0.212	0.674	0.995	1.97
18	58,522,227	rs76726049	T/C	1.06	0.014	0.176	0.459	0	0.173
19	1,056,493	rs3752246	C/G	1.15	0.175	0.227	0.624	2.03	4.759
19	1,854,255	rs149080927 (3'UTR)	GC/G	1.05	0.475	1.847	1.319	1.618	1.216
19	3,405,594	rs9749589	T/A	1.32	0.152	1.451	0.985	0.032	1.044
19	44,908,684	rs429358 (p.Cys112Arg)	T/C	3.32	0.151	1.464	0.71	0.65	0.668
19	45,738,583	rs76320948	C/T	1.03	0.036	0.163	0.289	0.011	0.415
19	48,710,247	rs2452170	A/G	1.01	0.499	1.099	1.253	1.995	1.439
19	49,950,060	rs9304690	C/T	1.05	0.243	0.378	0.568	1.119	1.245
19	51,224,706	rs3865444 (5'UTR)	C/A	1.01	0.316	0.275	1.356	0.594	0.496
19	54,313,903	rs1761461	A/C	1.01	0.495	1.49	1.104	1.653	1.236
20	413,334	rs1358782	G/A	1.05	0.235	0.7	0.819	0.336	0.858
20	56,423,488	rs6014724	A/G	1.12	0.089	1.119	1.931	3.979	1.134
20	63,743,088	rs6742	C/T	1.05	0.208	1.442	0.781	0.138	1.345
21	26,784,537	rs2830500	C/A	1.08	0.308	0.309	0.629	0.148	0.779

### Supplementary Figure 1. Simulated power curves for suggestive hits by effect size and allele frequency

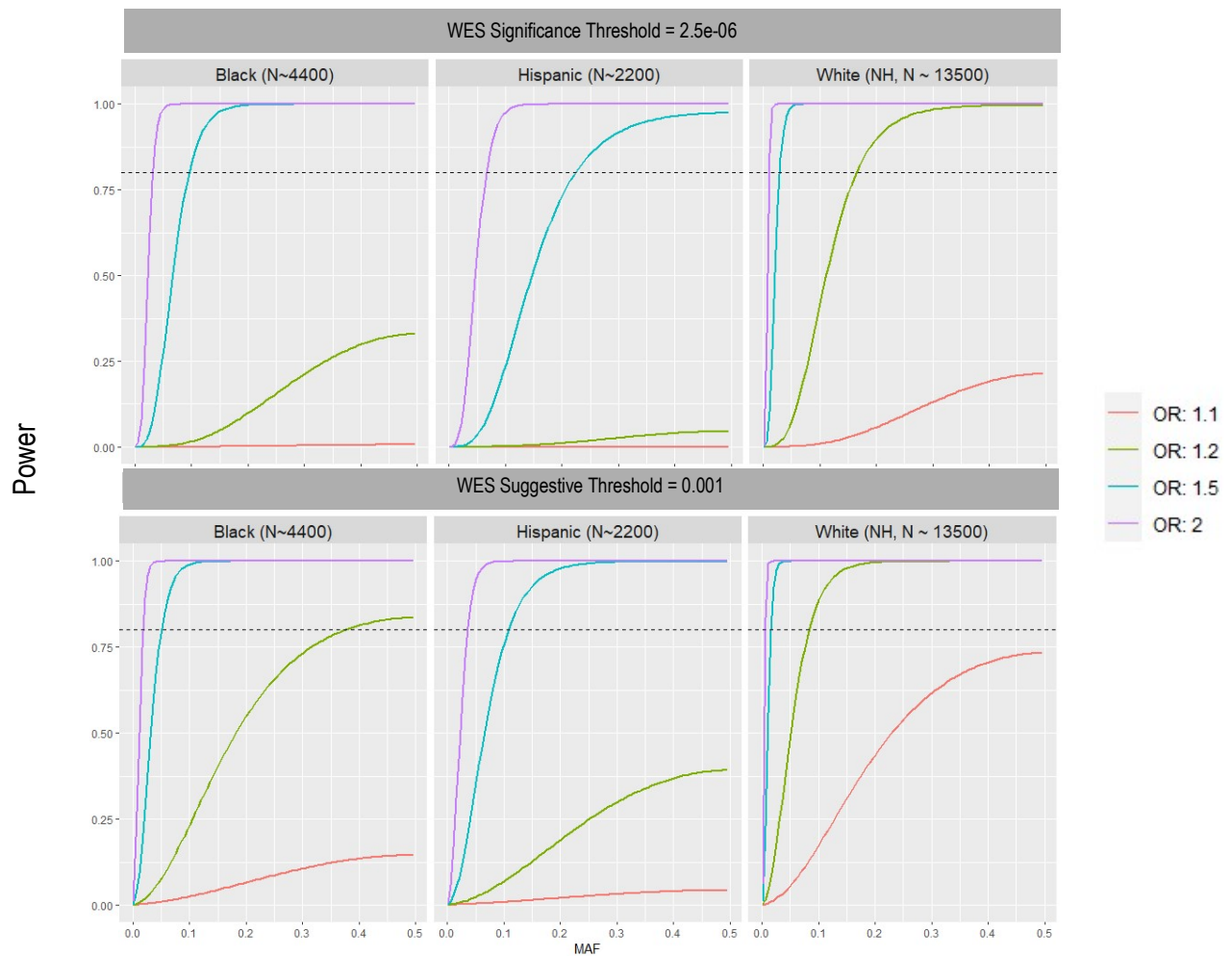
Each panel shows power to detect a suggestively associated SNP ( $\alpha = 1e-05$ ) for a different racial/ethnic group using GWAS array data based on current sample sizes. Power was simulated based on a set of effect sizes (Odds Ratios, OR = 1.1, 1.2, 1.5, 2) and minor allele frequencies (MAF) ranging from 0.001 to 0.5. The dashed line represents power = 0.80. American Indian and Alaskan Native samples are not included because current sample sizes are too small to detect any SNPs regardless of MAF or effect size.



## Supplementary Figure 2. Simulated power curves for whole exome sequencing data.

The top panel displays the power to detect significantly associated SNPs using a whole exome sequencing threshold of  $2.5 \times 10^{-6}$ . The lower panel displays the power to detect SNPs that meet the suggestive threshold of 0.001. Power was simulated based on a set of effect sizes (OR = 1.1, 1.2, 1.5, 2) and MAF ranging from 0.001 to 0.5. The dashed line represents power = 0.80.

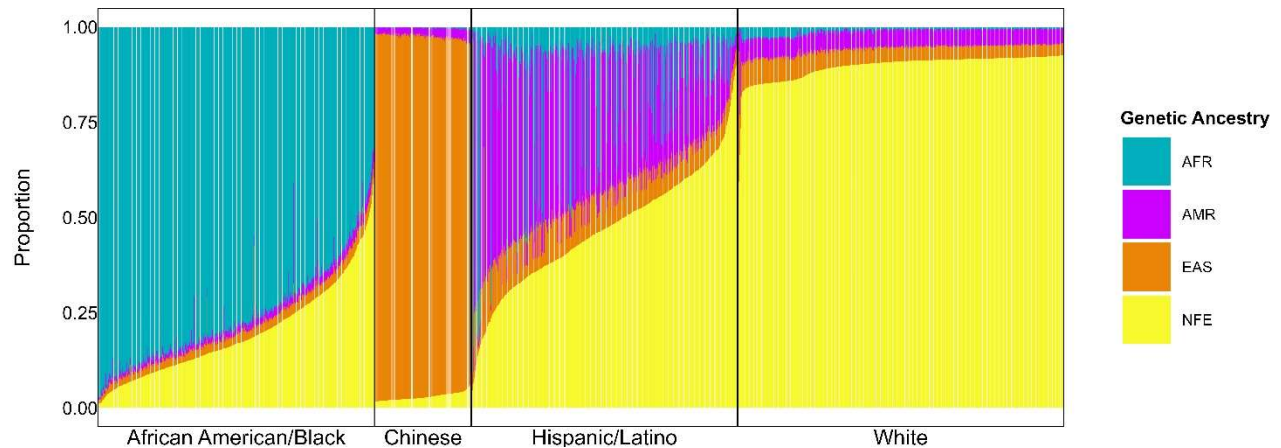
American Indian/Alaska Native and Asian populations are not included because current sample sizes are too small to detect any SNPs regardless of MAF or effect size.



## Paper 2 Supplement: Polygenic risk scores for incident dementia in the Multi-Ethnic Study of Atherosclerosis

### Supplementary Figure 1. Global Ancestry Proportions for MESA participants

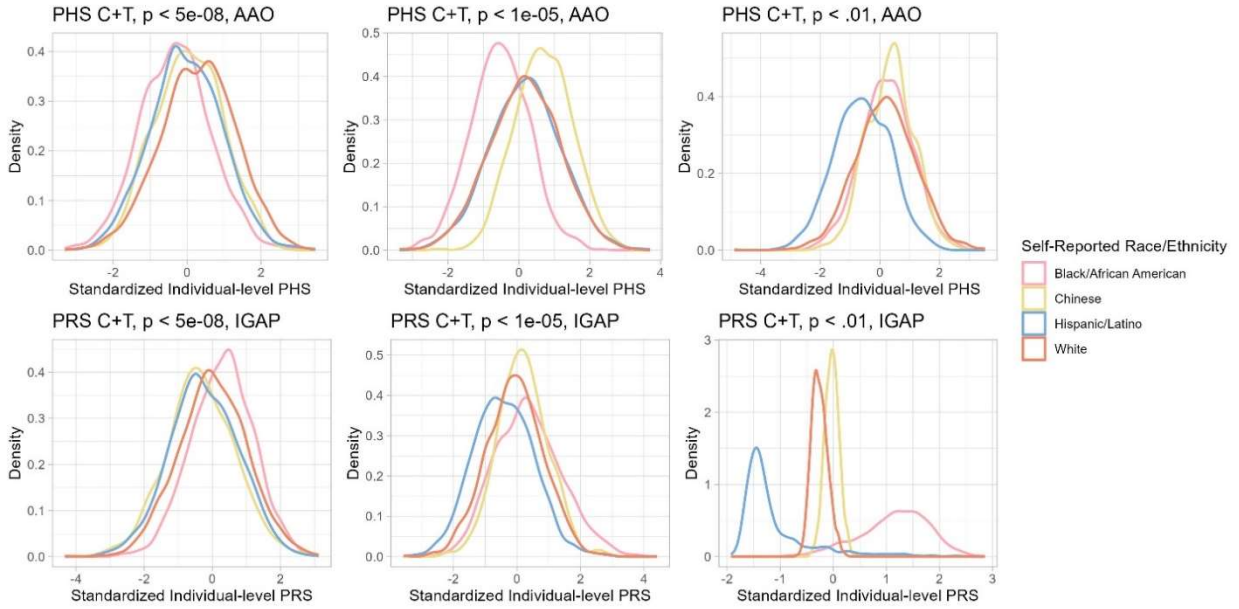
African Americans are an admixed group. This plot represents the proportion of African (AFR), Amerindian (AMR), East Asian (EAS) and Non-Finnish European (NFE) for each participant. Each participant represents one column. Reference samples are from gnomAD v3.1. I estimated the global ancestry proportions to test the hypothesis that the performance of PRS models may differ based on the proportion of European ancestry (or similarity to the GWAS summary statistics).



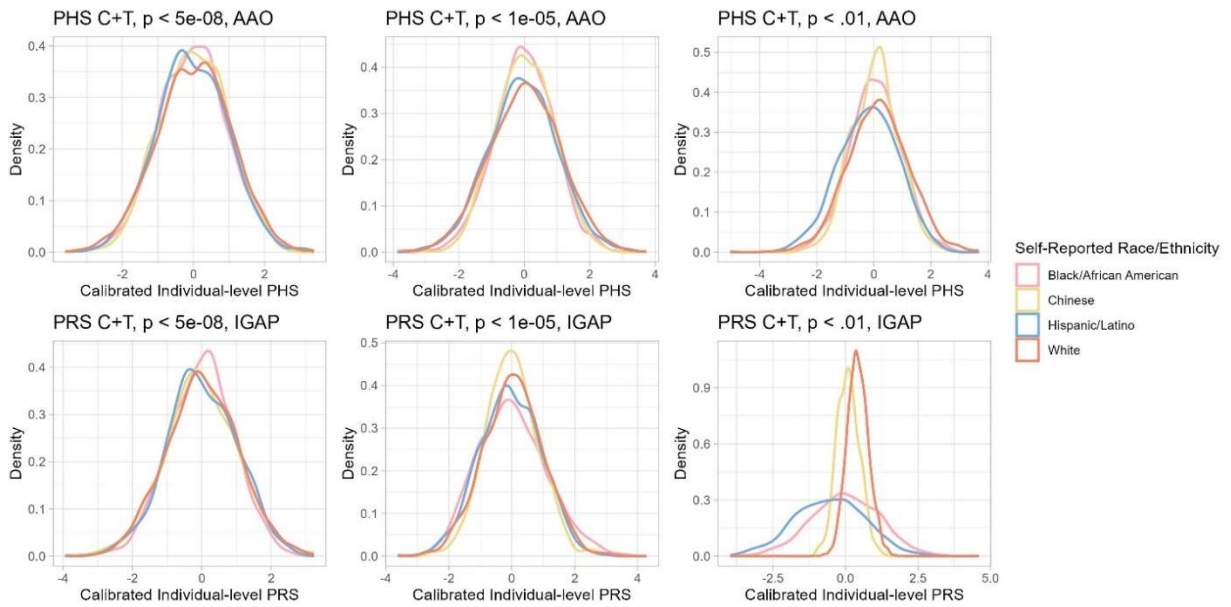
## **Supplementary Figure 2. Clumping and Thresholding Polygenic Hazard Scores and Polygenic Risk Scores by Race**

These distributions are only for the clumping and thresholding scores. The top row of distributions in both Pre-calibrated and Calibrated are calculated from the multi-ethnic age-at-onset GWAS meta-analysis. The bottom row distributions are calculated from the IGAP GWAS conducted in those with European ancestry.

## Pre-Calibrated

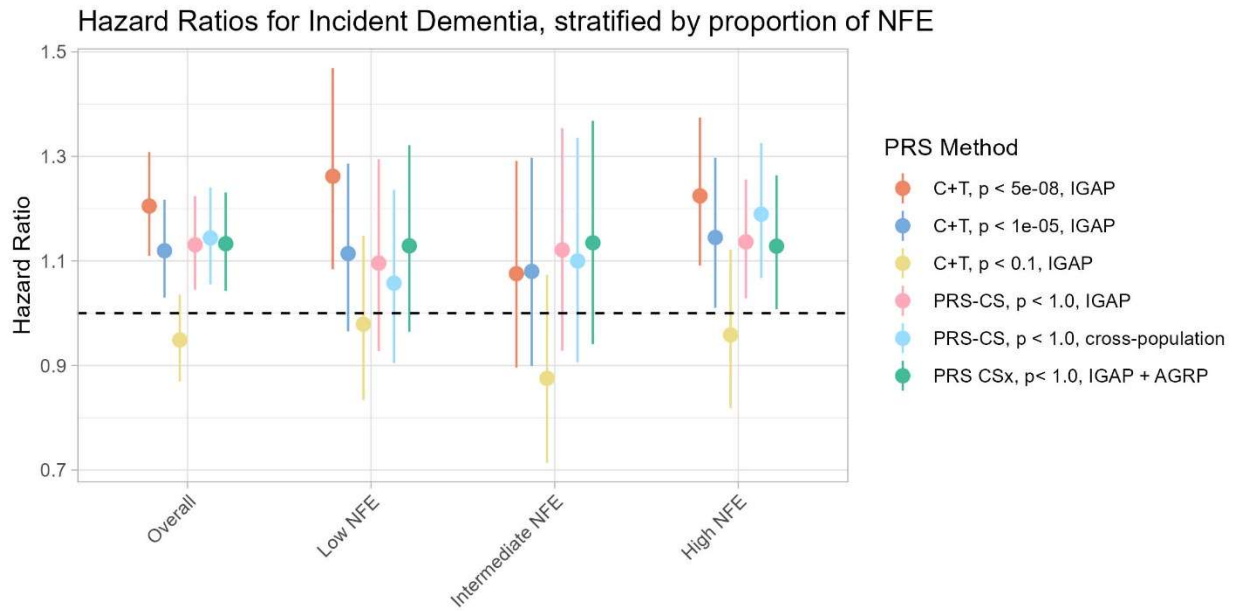


## Calibrated



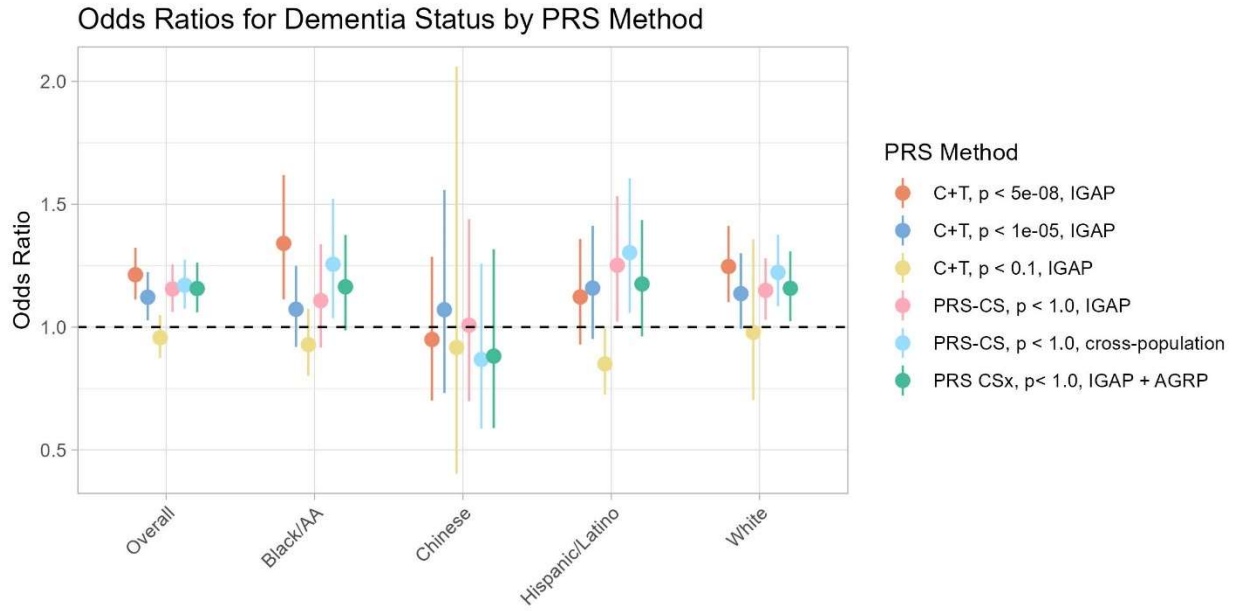
### Supplementary Figure 3. Association between adjusted PRS and incident dementia stratified by proportion of NFE ancestry

MESA participants were stratified based on their proportion of NFE ancestry. Low NFE includes those in the bottom tertile of NFE proportions and High NFE includes those in the top tertile of NFE proportions.



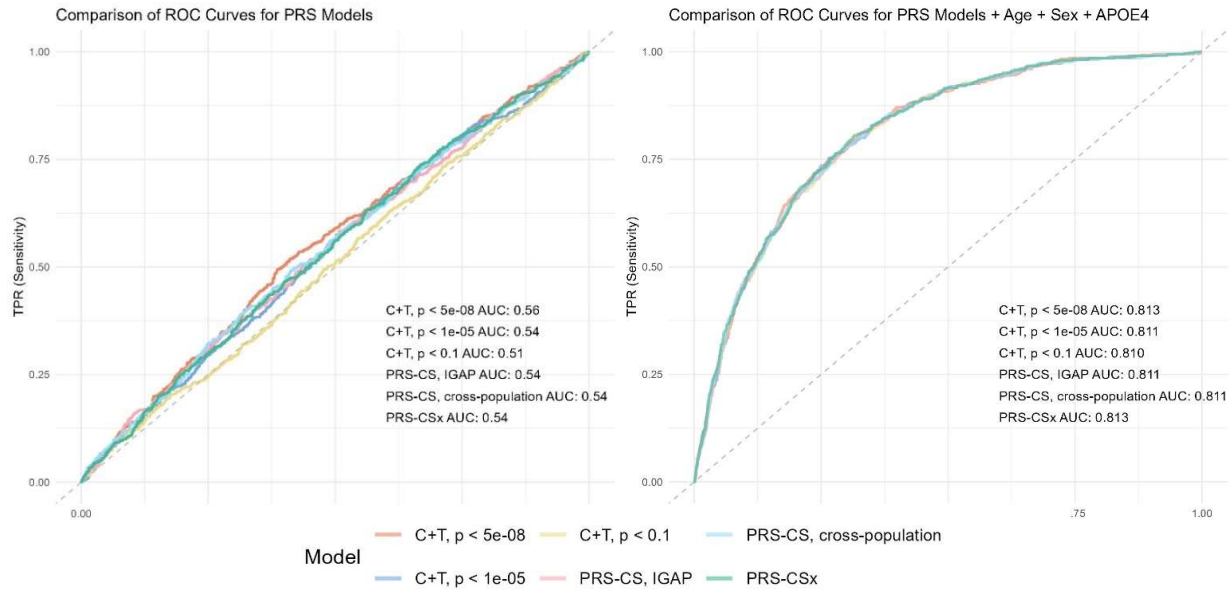
### Supplementary Figure 4. Association between adjusted PRS and dementia case-control status

As a sensitivity analysis, association between PRS and dementia was examined by estimating marginal odds ratios using generalized estimating equations.



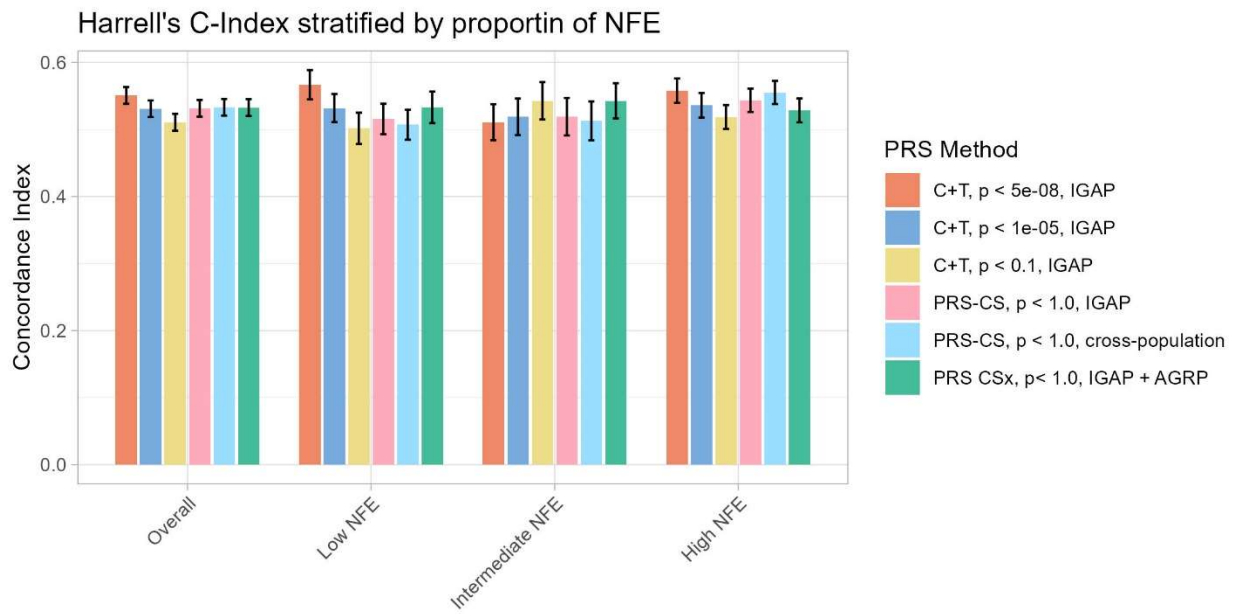
### Supplementary Figure 5. PRS predictive performance measured by AUC

The left-hand panel represents the predictive performance of a univariate model where the PRS alone are used to predict dementia status. The right-hand panel represents the predictive performance of a model that includes the PRS alongside age, sex, and *APOEε4* dosage.



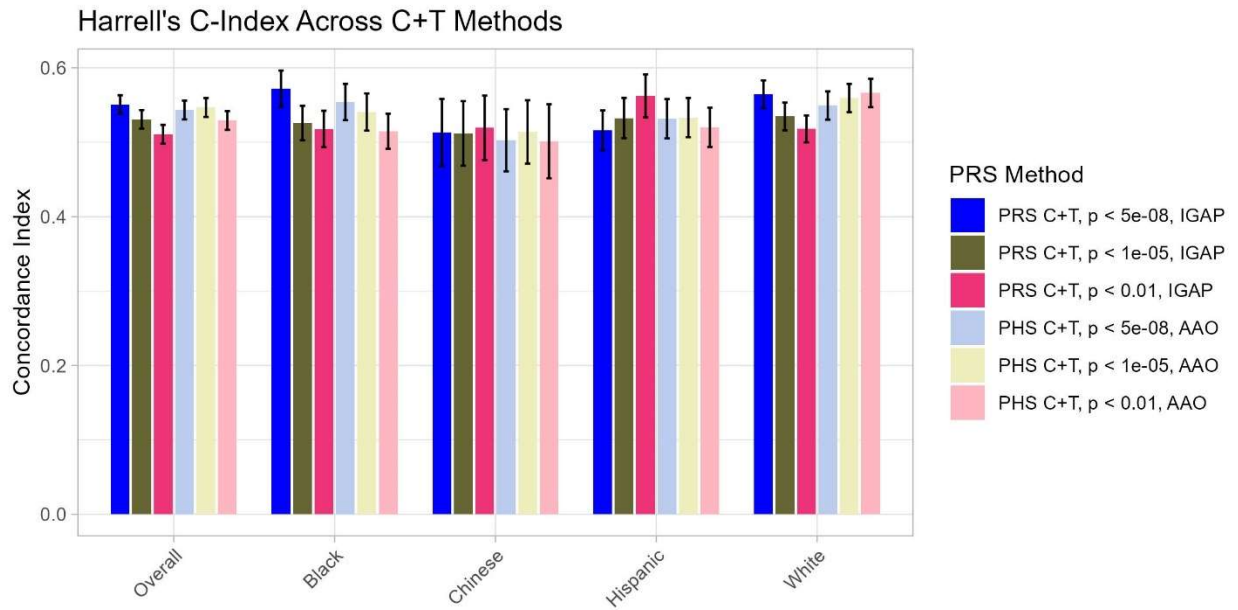
### Supplementary Figure 6. PRS predictive performance comparisons stratified by proportion of European ancestry

We compared PRS predictive performance across tertiles of NFE ancestry because of the large body of work that has shown that PRS performance deteriorates as the genetic distance between training and target populations grow.



### Supplementary Figure 7. Predictive performance of C+T Models, including PHS

These barplots compare the predictive performance of various C+T models that differ in their SNP inclusion threshold and source GWAS. Polygenic risk scores were constructed using the IGAP GWAS. Polygenic hazard scores were constructed from a multi-ethnic GWAS of age-at-onset of AD.



### **Supplementary Figure 8. Inter-model reliability among self-reported white MESA participants**

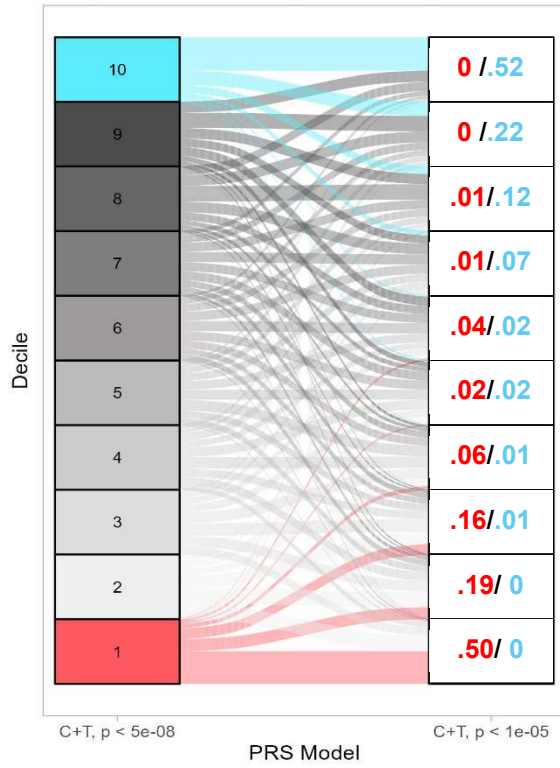
The Sankey plots display pairwise comparisons for the models constructed using the same GWAS summary statistics and have strong Spearman rank correlations. These comparisons are restricted to those who have self-reported non-Hispanic white race/ethnicity. We conducted this additional test to ensure that patterns seen in Figure 4 are not due to inconsistent PRS performance across those with different ancestral backgrounds.

*Red* indicates being in the lowest risk decile in the lefthand model.

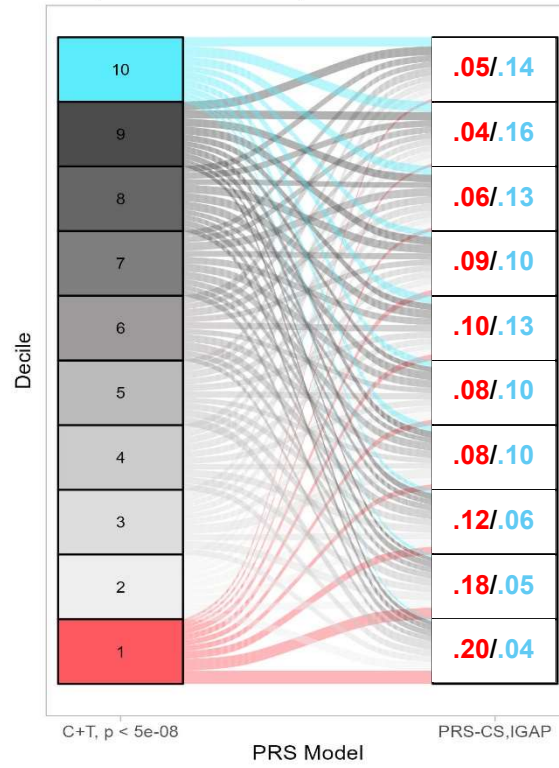
*Blue* indicates being in the highest risk decile in the lefthand model.

The *red* and *blue* proportions indicate the proportion of those in the lowest and highest risk deciles, respectively, that are in each risk decile according to the righthand model.

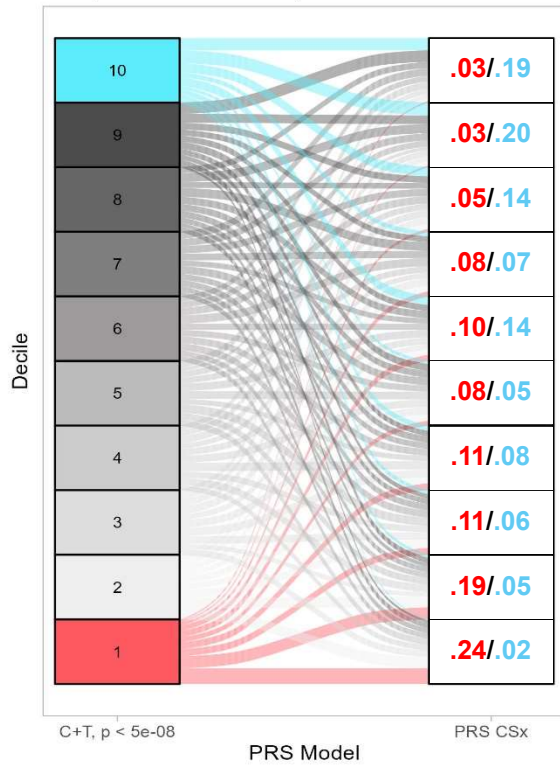
C+T p<5e-08 deciles compared to C+T p<1e-05



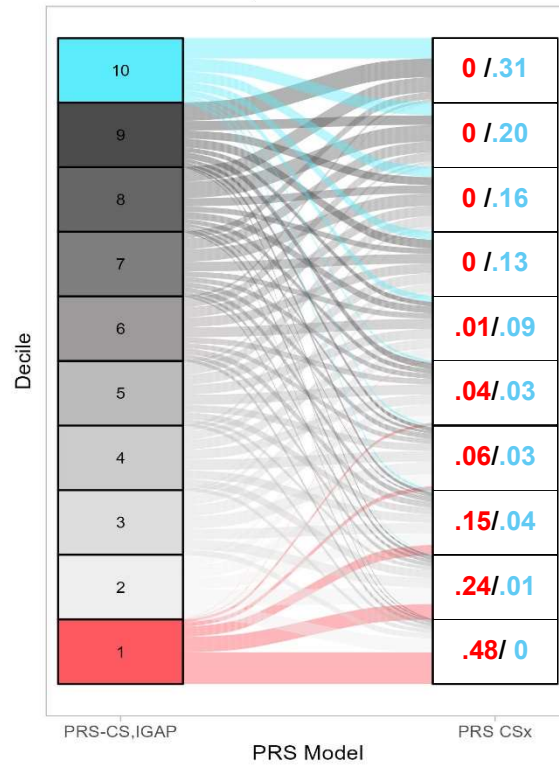
C+T p<5e-08 deciles compared to PRS-CS



C+T p<5e-08 deciles compared to PRS-CSx



PRS-CS deciles compared to PRS-CSx



Paper 3 Supplement: Integrating Contextual Determinants and Polygenic Risk to Examine Dementia and Cognition in the Multi-Ethnic Study of Atherosclerosis

**Supplementary Table 1. Sample size and dementia case rate in groups stratified by genetic risk**

<b>Genetic risk Category</b>	<b>Sample Size</b>	<b>Percentage of dementia cases</b>
Category 1: <i>e4</i> carriers	1368	11.4%
Category 2: non- <i>e4</i> , high	1042	9.2%
Category 3: non- <i>e4</i> , mid	2083	7.8%
Category 4: non- <i>e4</i> , low	1042	6.5%

**Supplementary Table 2. Effect sizes and confidence intervals of contextual determinants on dementia status in genetically stratified groups**

This table shows the odds ratio of contextual determinants on dementia status in genetically stratified groups. Bolded cells indicate statistically significant associations with dementia.

*\* models not adjusted for SES*

	<b>APOEε4 carriers</b>	<b>Non-ε4 carriers (High PRS)</b>	<b>Non-ε4 carriers (Mid. PRS)</b>	<b>Non-ε4 carriers (Low PRS)</b>
<b>SES*</b>	<b>1.38 (1.15, 1.65)</b>	1.04 (0.85, 1.26)	1.08 (0.94, 1.23)	1.08 (0.86, 1.36)
<b>PM2.5</b>	1.18 (0.82, 1.71)	0.58 (0.32, 1.06)	1.29 (0.86, 1.94)	0.94 (0.50, 1.74)
<b>NO2</b>	1.08 (0.40, 2.91)	0.56 (0.14, 2.32)	2.23 (0.87, 5.72)	0.56 (0.13, 2.40)
<b>Favorable Food Store</b>	0.99 (0.83, 1.17)	0.91 (0.74, 1.11)	0.99 (0.85, 1.14)	1.21 (0.96, 1.53)
<b>Fast Food Chain</b>	1.03 (0.89, 1.20)	0.96 (0.78, 1.19)	0.99 (0.87, 1.14)	1.09 (0.88, 1.34)
<b>Alcohol Establishment</b>	0.96 (0.83, 1.11)	1.07 (0.87, 1.32)	1.01 (0.88, 1.15)	1.01 (0.82, 1.25)
<b>Physical Activity Facility</b>	1.02 (0.90, 1.17)	1.04 (0.82, 1.30)	0.99 (0.87, 1.13)	0.96 (0.89, 1.14)

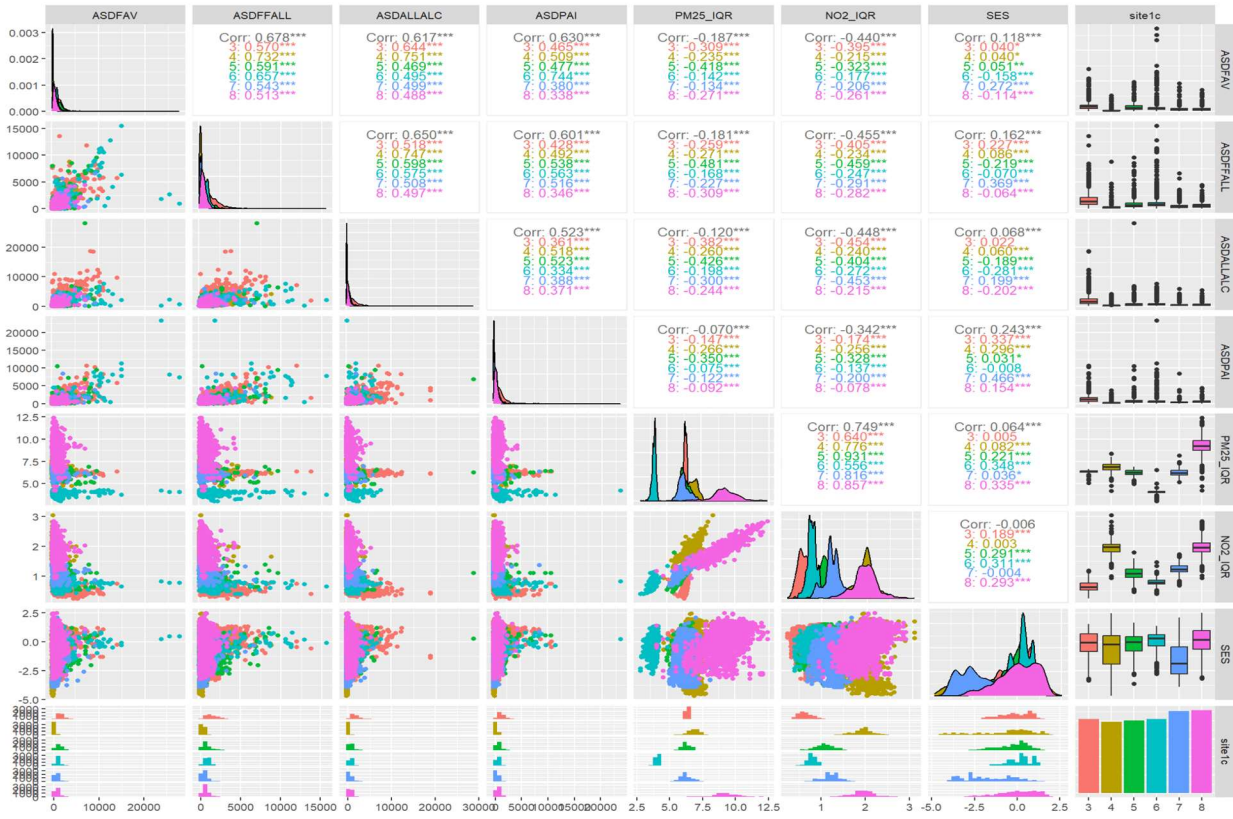
### Supplementary Table 3. Posterior inclusion probabilities

In BKMR, posterior inclusion probabilities (PIPs) represent the likelihood that a variable is important in the model in explaining the outcome. Higher values indicate stronger evidence of the variable's importance. Bolded cells indicate PIPs greater than 0.5.

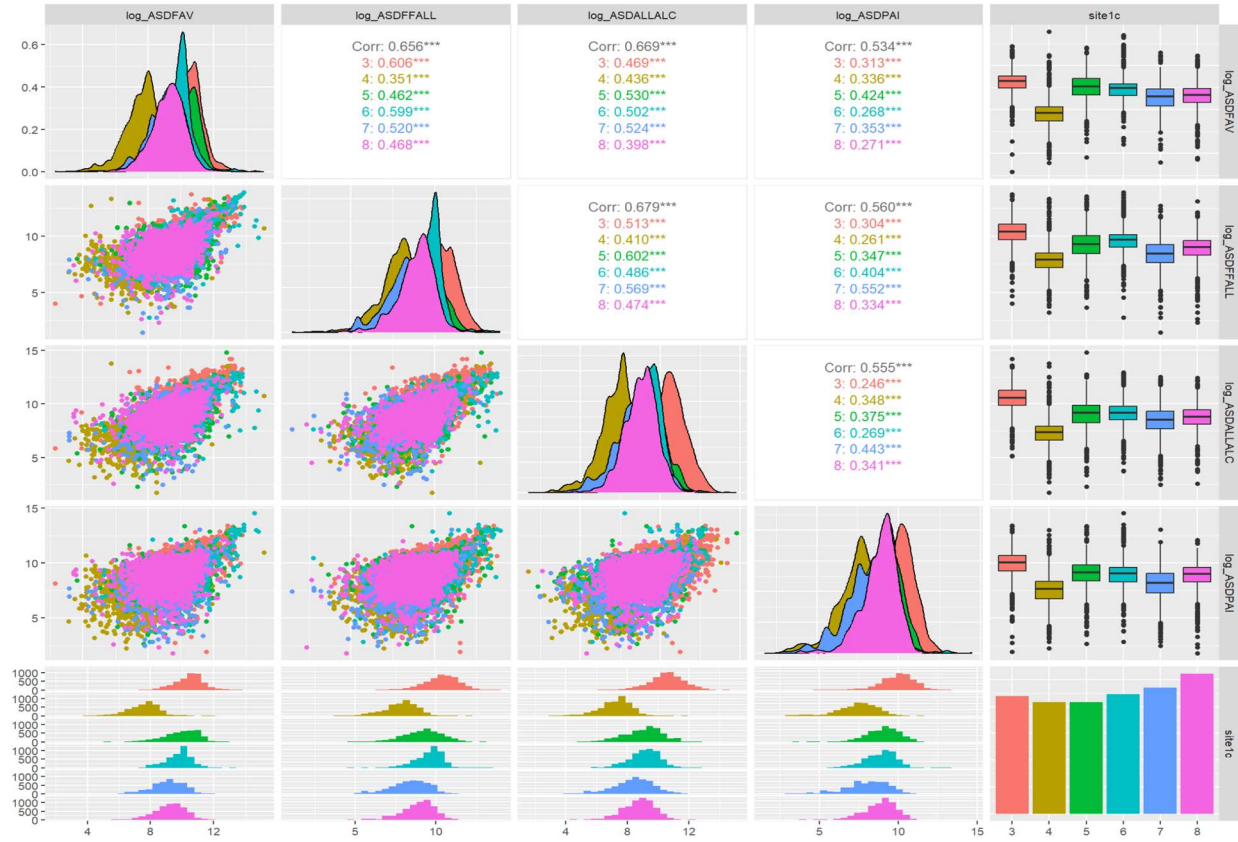
	<i>e4</i> Carriers	Non- <i>e4</i> , High	Non- <i>e4</i> , Mid	Non- <i>e4</i> , Low
<b>SES</b>	0.0000	0.0024	0.0000	0.0016
<b>PM2.5</b>	0.0030	0.0036	0.0608	0.0428
<b>NO2</b>	<b>0.4054</b>	0.0236	<b>0.8018</b>	0.0926
<b>ASDFAV</b>	0.2182	0.0052	0.0010	0.0014
<b>ASDFFALL</b>	0.0812	0.0018	0.0000	0.0000
<b>ASDALC</b>	0.0832	0.0000	0.0082	0.0000
<b>ASDPAI</b>	0.0036	0.0042	0.0000	0.0040

# Supplementary Figure 1. Distribution of contextual determinants

Distributions differ by site. The built environment variables are right skewed.



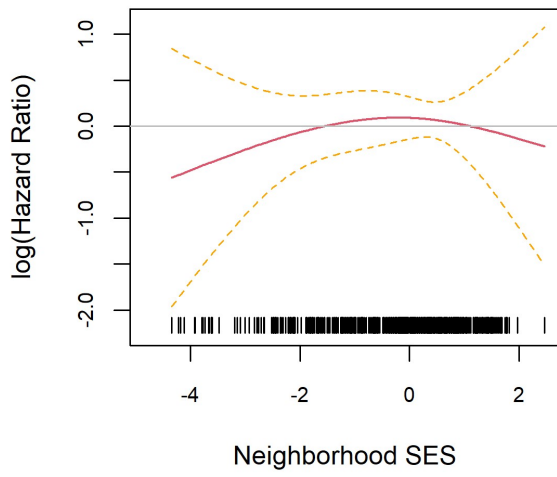
## Supplementary Figure 2. Distribution of built environment variables after log<sub>2</sub> transformation



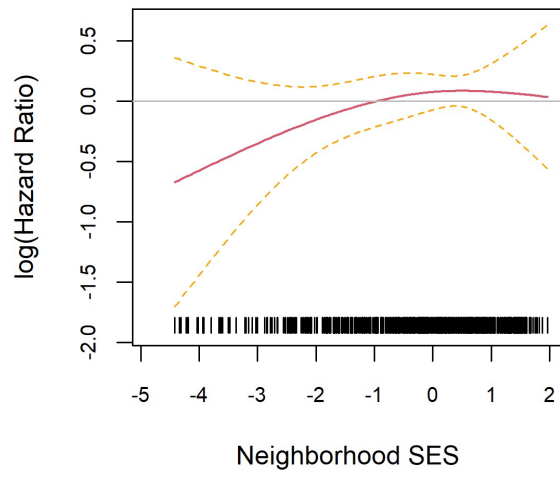
### **Supplementary Figure 3. Nonlinear interactions between Neighborhood SES and genetic risk on incident dementia risk**

The MESA participants were stratified by genetic risk into four groups. The *e4* carriers are considered the highest genetic risk group. The non-*e4* carriers are then split into three groups, with the high and low groups being the top and bottom quartile of polygenic risk scores among the non-*e4* carriers. Neighborhood SES was modeled using a two degree of freedom cubic spline. ANOVA tests were used to assess the statistical significance of differences in association between SES and dementia across the four genetic risk groups. The associations between SES and incident dementia were not significantly different across genetic risk groups (ANOVA:  $P > 0.5$ ).

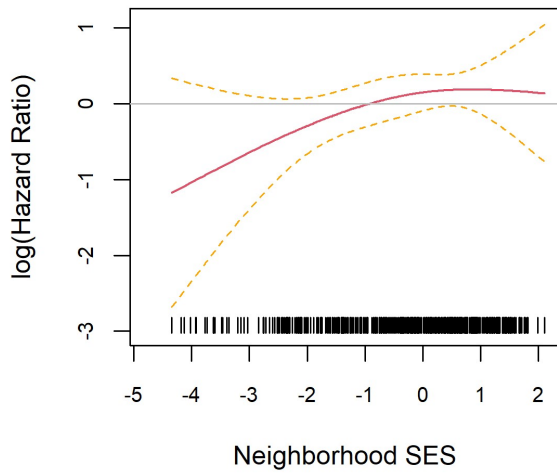
**non-e4, Low, n = 1042**



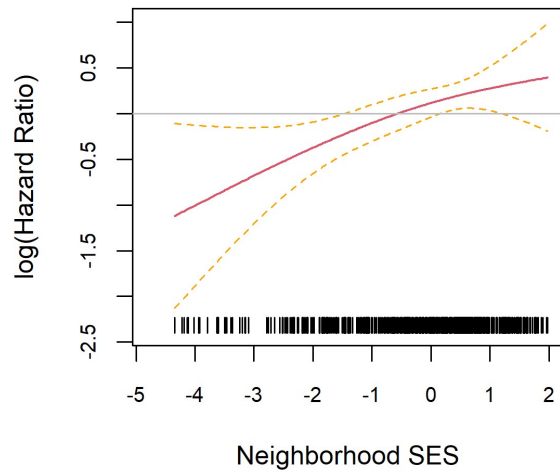
**non-e4, Intermediate, n = 2083**



**non-e4, High, n = 1042**

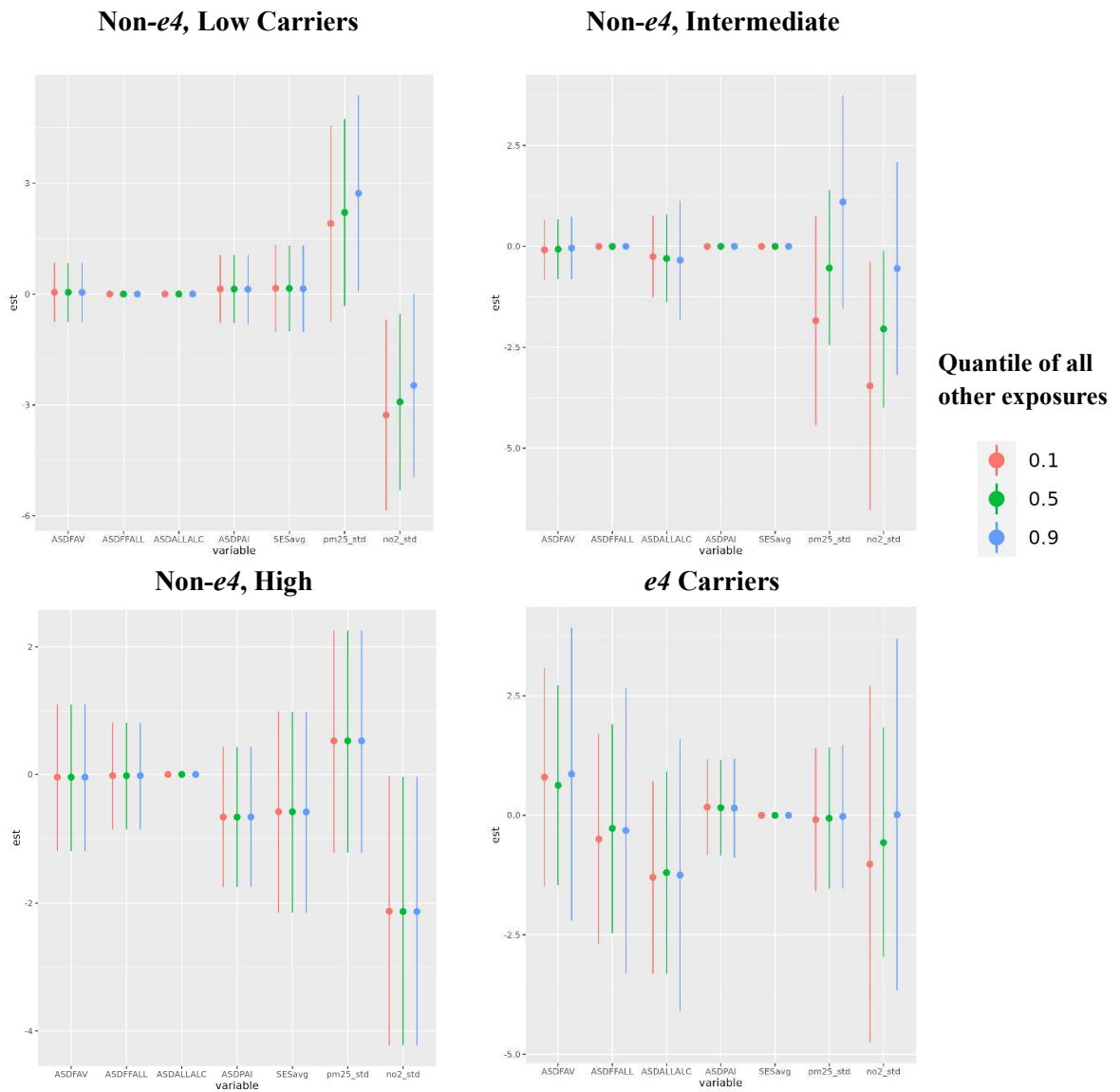


**e4 Carriers, n = 1546**



### Supplementary Figure 4. The effect of single exposures on the overall mixture

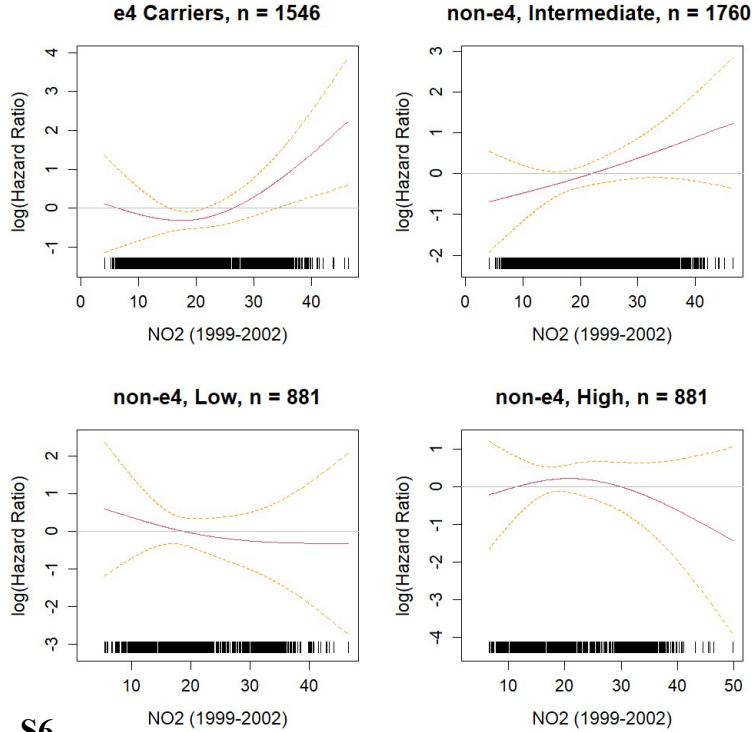
BKMR also allows for the investigation of the contribution of individual exposures to the mixture. Within the group of *e4* carriers, no individual exposure is significantly associated. Among the non-*e4* carriers with low and intermediate polygenic risk, only NO<sub>2</sub> contributes to variation in cognition scores. Higher values of NO<sub>2</sub> are associated with lower values of the mixture. As the other non-NO<sub>2</sub> exposures increase in value from their 10<sup>th</sup> to 90<sup>th</sup> percentile, the effect of NO<sub>2</sub> decreases.



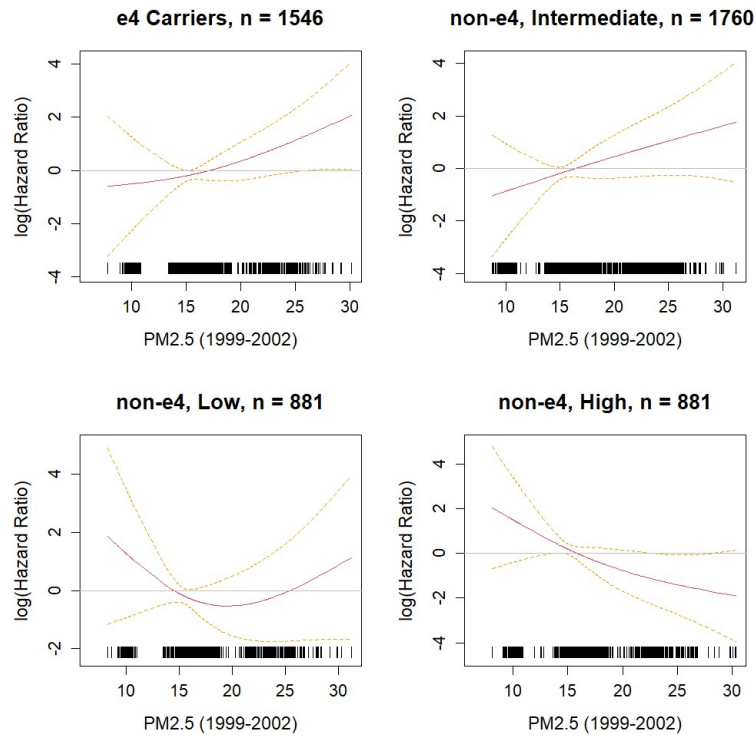
## Supplementary Figures 5-10. Nonlinear interaction plots

S5) NO<sub>2</sub>, S6) PM<sub>2.5</sub>, S7) Favorable Food, S8) Fast Food, S9) Alcohol, S10) Physical Activity

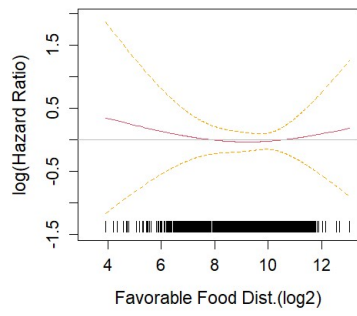
**S5**



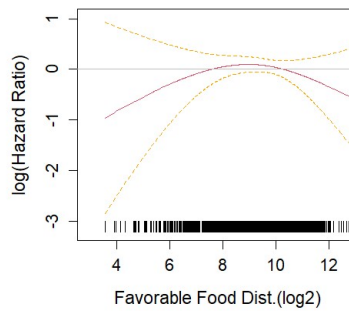
**S6**



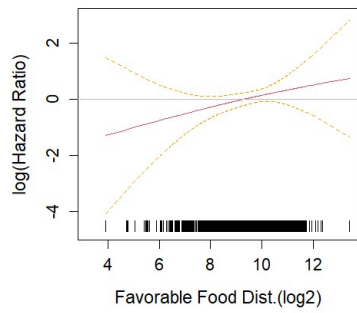
**S7** e4 Carriers, n = 1546



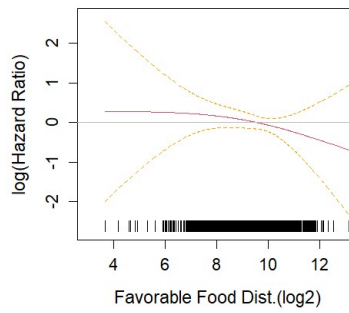
non-e4, Intermediate, n = 1760



non-e4, Low, n = 881

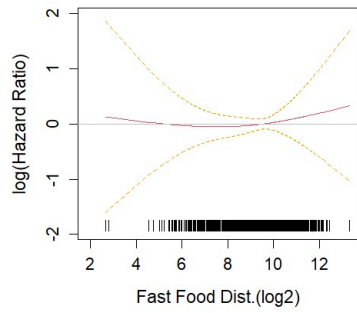


non-e4, High, n = 881

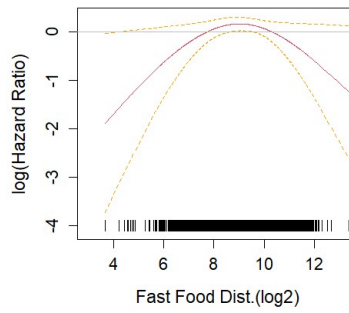


**S8**

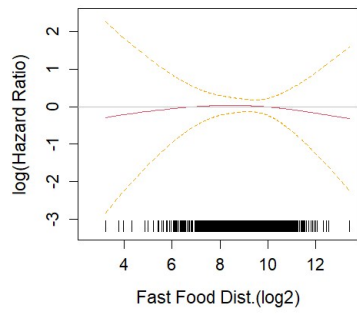
e4 Carriers, n = 1546



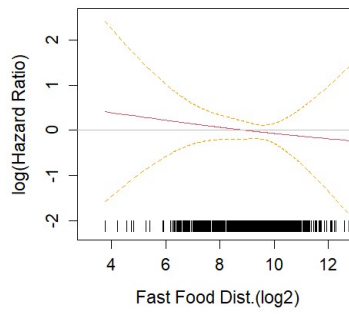
non-e4, Intermediate, n = 1760



non-e4, Low, n = 881

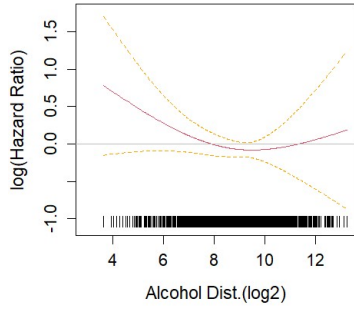


non-e4, High, n = 881

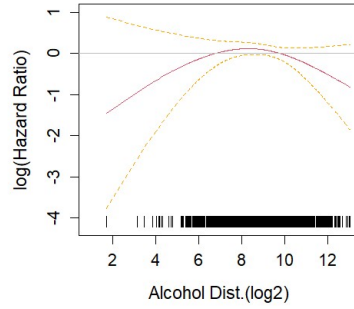


**S9**

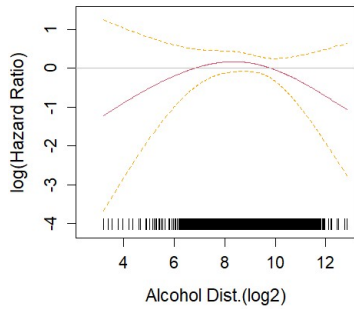
**e4 Carriers, n = 1546**



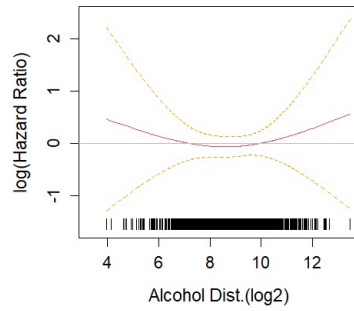
**non-e4, Intermediate, n = 1760**



**non-e4, Low, n = 881**

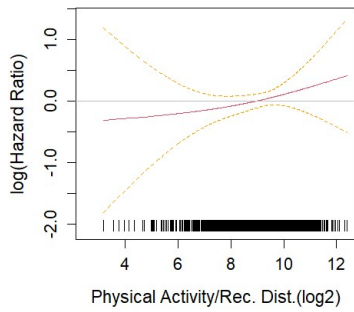


**non-e4, High, n = 881**

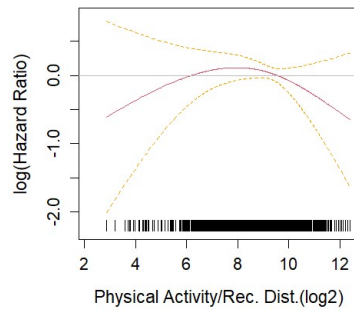


**S10**

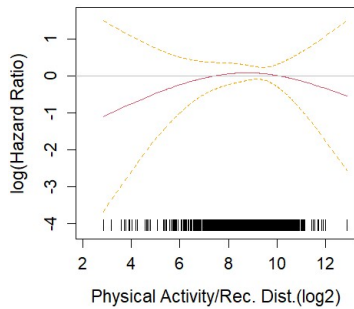
**e4 Carriers, n = 1546**



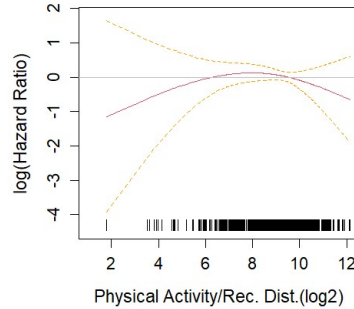
**non-e4, Intermediate, n = 1760**



**non-e4, Low, n = 881**



**non-e4, High, n = 881**



End of material.