

1 **A fast and robust strategy to remove variant level artifacts in Alzheimer’s Disease Sequencing**
2 **Project data**

3
4 **Authors:** Michael E. Belloy¹, PhD, Yann Le Guen^{1,2}, PhD, Sarah J. Eger¹, BA, Valerio Napolioni³, PhD,
5 Michael D. Greicius¹, MD, MPH, Zihuai He^{1,4}, PhD.

6
7 ¹Department of Neurology and Neurological Sciences, Stanford University, Stanford, CA, 94304, USA

8 ²Institut du Cerveau - Paris Brain Institute - ICM, Paris, 75013, France

9 ³School of Biosciences and Veterinary Medicine, University of Camerino, Camerino, 62032, Italy

10 ⁴Quantitative Sciences Unit, Department of Medicine, Stanford University, Stanford, CA, 94304, USA

11

12 **Corresponding Author**

13 Michael E. Belloy

14 Department of Neurology and Neurological Sciences – Greicius lab

15 Stanford University

16 290 Jane Stanford Way, Stanford, CA, USA

17 Tel: 650 498 4624

18 Email: mbelloy@stanford.edu

19

20 **Word/Table count**

21	Abstract	Word count: 259
22	Author Summary	Word count:
23	Main text	Word count: 3774
24	Introduction	Word count: 429
25	Materials & Methods	Word count: 1314
26	Results	Word count: 1141
27	Discussion	Word count: 890
28	References	Count: 30
29	Figures	Count: 4
30	Tables	Count: 2

31 **Abstract**

32 Whole-exome sequencing (WES) and whole-genome sequencing (WGS) are expected to be critical to
33 further elucidate the missing genetic heritability of Alzheimer’s disease (AD) risk by identifying rare
34 coding and/or noncoding variants that contribute to AD pathogenesis. In the United States, the
35 Alzheimer’s Disease Sequencing Project (ADSP) has taken a leading role in sequencing AD-related
36 samples at scale, with the resultant data being made publicly available to researchers to generate new
37 insights into the genetic etiology of AD. In order to achieve sufficient power, the ADSP has adapted a
38 study design where subsets of larger AD cohorts are collected and sequenced across multiple centers,
39 using a variety of sequencing kits. This approach may lead to variable variant quality across sequencing
40 centers and/or kits. Here, we performed exome-wide and genome-wide association analyses on AD risk
41 using the latest ADSP WES and WGS data releases. We observed that many variants displayed large
42 variation in allele frequencies across sequencing centers/kits and contributed to spurious association
43 signals with AD risk. We also observed that sequencing kit/center adjustment in association models
44 could not fully account for these spurious signals. To address this issue, we designed and implemented
45 novel filters that aim to capture and remove these center/kit-specific artifactual variants. We conclude
46 by deriving a novel, fast, and robust approach to filter variants that represent sequencing center- or kit-
47 related artifacts underlying spurious associations with AD risk in ADSP WES and WGS data. This
48 approach will be important to support future robust genetic association studies on ADSP data, as well as
49 other studies with similar designs.

50 **Author Summary**

51 Next generation sequencing data represents a highly valuable resource to uncover rare coding and/or
52 noncoding genetic variants that contribute to Alzheimer’s disease risk. In order to achieve large sample
53 sizes that are required for such data, the Alzheimer’s Disease Sequencing Project (ADSP) has taken the
54 leading role in sequencing Alzheimer’s disease related samples at scale in the United States. The ADSP’s
55 study design however leads to variable variant quality across the involved sequencing centers,
56 necessitating a quality control approach that ensures robust genetic association analyses. Here, we
57 present and validate a rigorous quality control pipeline, where we specifically developed a new strategy
58 to handle inter-center variant quality issues in the ADSP. In doing so, we provide a first glance into
59 exome- and genome-wide associations with Alzheimer’s disease risk using the latest releases of ADSP
60 data (respectively 20.5k and 16.9k individuals). In sum, our pipeline is important to support future
61 robust genetic association studies on ADSP data, as well as other studies with similar design. This in turn
62 will contribute to accelerating Alzheimer’s disease gene discovery and gene-driven therapy
63 development.

64 Introduction

65 Late-onset Alzheimer's disease (AD) is marked by a strong genetic component, with heritability
66 estimates ranging from 59% to 79%^{1,2}. Largely supported by single nucleotide polymorphism (SNP)
67 genotyping arrays and variant imputation, large-scale meta-analyses of genome-wide association
68 studies have so far implicated over 50 loci relevant to AD in subjects of European ancestry²⁻⁶. Despite
69 these important advances, most risk variants identified so far have common allele frequencies and it's
70 estimated that only about half of AD's genetic heritability has been captured, such that much of AD's
71 genetic component remains to be identified². In response to this observation, there has been a shift to
72 start using whole-exome sequencing (WES) or whole-genome sequencing (WGS) to help capture rare
73 and/or coding variants that contribute to AD risk, which has led to several recent initial successes⁷⁻¹⁵.

74 In the United States, the Alzheimer's Disease Sequencing Project (ADSP) has taken a leading role in
75 sequencing of AD-related samples at scale, with resultant data being made publicly available to
76 researchers to generate new insights into the genetic etiology of AD. In order to achieve sufficient
77 power to support analyses of sequencing data and rare variants, the ADSP has adapted a study design
78 where subsets of larger AD cohorts are collected and sequenced across multiple centers, using a variety
79 of sequencing kits¹⁶⁻¹⁸. This in turn can lead to "center" or "kit" effects that traditionally are accounted
80 for by using center/kit covariate adjustment. However, a prior study using a prior version of the ADSP
81 WES discovery phase observed that center/kit covariate adjustment could not account for variable
82 variant qualities across centers and kits, which in turn may lead to spurious associations or impact the
83 identification of AD-associated risk variants¹⁹.

84 Since then, the ADSP has further expanded its efforts and as of 2021 provides WES and WGS data on
85 respectively 20.5k and 16.9k individuals across diverse ancestries¹⁸. In our exploratory analyses of these
86 data, we observed many variants that displayed large variation in allele frequencies across centers/kits
87 and contributed to spurious association signals with AD risk. Similar to the prior study¹⁹, we also
88 observed that kit/center adjustment could not fully account for these signals. Thus, in the current study,
89 we design and implement novel filters that aim to capture and remove these center/kit-specific
90 artifactual variants. We additionally test filters containing putatively artifactual variants identified in the
91 gnomAD reference database²⁰. The filters are designed such that they can be implemented post-hoc to
92 association analyses, leaving flexibility to researchers to either run full sample analyses with robust
93 variant quality control, or, to identify variants that require targeted analyses.

94 **Methods**

95 **Ascertainment of Genotype and Phenotype Data**

96 Genotype data for subjects with AD-related clinical outcome measures were available from the
97 Alzheimer’s Disease Sequencing project (ADSP) whole exome sequencing (WES) and whole genome
98 sequencing (WGS) data. Notably, the ADSP performed targeted sequencing of samples in case-control
99 (majority), family-based, population-based, and longitudinal cohorts, performing sequencing across
100 multiple sequencing centers and using various sequencing kits (**Table S1-2**). Ascertainment of
101 genotype/phenotype data for these samples is described in detail elsewhere¹⁸. In addition to the ADSP
102 samples, we also had access to several publicly available SNP microarray and WGS datasets (**Table S1**),
103 largely comprising data from the Alzheimer’s Disease Genetics Consortium (ADGC). The latter have a
104 large degree of sample overlap with ADSP. In order to ensure the most up-to-date and parsimonious
105 phenotypes, we performed a cross-sample genotype/phenotype harmonization, which is summarized in
106 **Supplementary Methods**.

107 Participants or their caregivers provided written informed consents in the original studies. The
108 current study protocol was granted an exemption by the Stanford Institutional Review Board because
109 the analyses were carried out on “de-identified, off-the-shelf” data.

110 **Genetic Data Quality Control and Processing**

111 The ADSP WES and WGS data (NG00067.v5) were joint called by the ADSP following the SNP/Indel
112 Variant Calling Pipeline and data management tool used for analysis of whole genome and exome
113 sequencing (WGS/WES) for the Alzheimer’s Disease Sequencing Project (VCPA)²¹. The WES data was
114 currently only released for bi-allelic variants, which the ADSP has quality controlled. The WGS data was
115 released for bi-allelic and multi-allelic variants separately, which the ADSP had not yet quality controlled.
116 The current analyses of ADSP WGS were restricted to bi-allelic variants, to which we applied the Variant
117 Quality Score Recalibration (VSQR) quality control filter (“PASS” variants; GATK v4.1)²². The WES/WGS
118 data were available in genome build hg38, which we annotated using dbSNP153 variant identifiers.

119 Genetic data underwent standard quality control (QC). Detailed descriptions of all processing
120 procedures and sequential sample filtering steps are in **Supplementary Methods** and **Table S3-4**. For the
121 purpose of the presented genetic association analyses, only non-Hispanic subjects of European ancestry
122 were considered to focus on the largest ancestry population (SNPweights v2.1; **Figure S1**)²³. Principal

123 component analysis of genotyped SNPs provided principal components (PCs) capturing population
124 substructure (PC-AiR, **Figure S2**)²⁴. In both the WES and WGS data respectively, variants with a
125 genotyping rate less than 95%, deviating from Hardy Weinberg Equilibrium (HWE) in the full sample or in
126 controls ($p < 10^{-6}$), and a minor allele count less than 10, were excluded. After this standard quality
127 control, the total number of remaining variants was 224,270 for ADSP WES and 14,772,936 for ADSP
128 WGS.

129 **Primary filters to remove sequencing center/kit-related variant level artifacts**

130 We designed filters to assess whether there were significant deviations in genotype distributions for
131 any given variant across sequencing centers and kits respectively. To avoid bias from frequency
132 differences across cases and controls, we only assessed genotypes in control individuals.

133 The primary filters made use of the fast Fisher exact test as implemented by Plink (v.1.9; command --
134 fisher)²⁵. However, this test can currently only be implemented by comparing two groups at a time (e.g.
135 two genotyping centers) while we observed variant issues across multiple groups. We therefore
136 compared every individual sequencing center/kit to all others and combined the P-values from the
137 multiple tests through the Cauchy combination test²⁵. Variants with a combined P-value lower than the
138 heuristic threshold of 10^{-5} were flagged to be filtered.

139 We additionally tested two other types of sequencing center/kit-based variant filters. On one hand,
140 we performed chi square tests (R v.3.6.0) that respectively considered all sequencing centers or kits at
141 once. Variants with a P-value lower than the heuristic threshold of 10^{-5} were flagged to be filtered. On
142 the other hand, we performed Fisher tests with Monte Carlo (MC) simulation of P-values (R v.3.6.0) that
143 respectively considered all sequencing centers or kits. The MC approach was chosen to allow feasible
144 run times. Variants with a P-value lower than the heuristic threshold of 10^{-3} were flagged to be filtered
145 (this threshold reflects that the P-values from MC simulation are less small than those obtained for the
146 other tests).

147 The three filters were compared in terms of speed by calculating the time needed to derive the
148 respective variant filters on a 1MB genetic region of chromosome 1 in ADSP WGS. Computing time was
149 evaluated on a single CPU from an 80-core Xeon Gold 6138T processor @ 2.00GHz.

150 **Filters from the Genome Aggregation Database (gnomAD)**

151 In addition to the filters proposed above, we used the gnomAD data base (v3.1.1) reference to
152 identify potential variant artifacts²⁰. Specifically, we created filters for variants that have: (1) a “non-
153 PASS” flag in gnomAD, corresponding to those that did not pass gnomAD sample quality control filters
154 and may thus be more prone to sequencing issues; (2) a “LCR” flag in gnomAD, corresponding to those
155 located in a Low Complexity Region and may thus be more prone to low coverage, read misalignment,
156 and subsequent genotype issues; (3) a differential frequency of more than 10% between our current
157 samples and non-Finish European (nfe) participants in gnomAD, which may indicate an issue with those
158 variants in our samples. The three gnomAD filters were evaluated with the goal of supporting the
159 primary ADSP WES/WGS center/kit-based variant filters.

160 **Filters for discordant variants across duplicate samples**

161 A final set of filters was designed to flag variants that are discordant across duplicate samples.
162 Notably, the ADSP WES and WGS data both respectively contain a few hundred duplicate samples,
163 generally covering multiple sequencing centers and/or kits. Discordant variants across such duplicates
164 therefore provide a reference of artifactual variants that should be removed and are largely reflecting
165 center/kit-related genotyping issues. We evaluated these filters with the primary goal of comparing
166 them with the primary ADSP WES/WGS center/kit-based variant filters as well as the gnomAD-based
167 variant filters. In a secondary goal, we also assessed to what extent these duplicate discordant variant
168 filters themselves could handle center/kit-related variant issues that drove observations of spurious
169 association signals.

170 **Statistical analyses, Variant Annotation, and Visualization**

171 Exome-wide and genome-wide association studies on AD case-control status were conducted
172 respectively on ADSP WES and WGS, using LMM-BOLT (v.2.3.5). LMM-BOLT employs a Bayesian mixture
173 model that allows the inclusion of related individuals by adjusting for the genetic relationship matrix
174 (GRM)²⁶, thereby maximizing sample size and power. Given the current minor allele count thresholds,
175 the approximate fifty-fifty ratio of cases to controls, and sample sizes exceeding 5,000 participants for
176 both ADSP WES and WGS, the resultant test statistics are expected to be well-calibrated²⁶. After
177 analyses, association statistics were transformed back to a logistic scale taking into account the case
178 fraction²⁶. Per convention, variants were considered at suggestive ($P \leq 10^{-5}$) or genome-wide ($P \leq 5 \times 10^{-8}$)
179 significance.

180 Case-control association analyses considered two models. Model-1 included covariates for sex,
181 *APOE**4 dosage, *APOE**2 dosage, and the first 5 genetic PCs. We did not adjust for age as we previously
182 showed that this can lead to significant power loss when the age of cases is younger than for controls¹⁵,
183 which is true for ADSP given their initial design to prioritize old controls and young cases (**Table 1 &**
184 **Table S5-6**). Model-2 was the same as Model-1 but additionally included covariates for sequencing
185 center and kit.

186 The *APOE* locus (1Mb region centered on *APOE*) was removed from all summary statistics.
187 Independent loci were determined by sliding window when no variants with $P \leq 10^{-5}$ were observed
188 within 200Kb from one another. Manhattan plots provide RefSeq curated gene annotations for the gene
189 closest (<500Kb) to the top significant variant per locus. Only variants with $P \leq 10^{-6}$ were annotated to
190 improve visualization. Suggestive significance levels were indicated by gray dotted lines and green dots
191 for variants. Genome-wide significance levels were indicated by black solid lines and red dots for
192 variants. Variant densities were indicated at the bottom of Manhattan plots (dark green = low density,
193 yellow=medium density, red = high density). Plots were generated using the R package CMplot²⁷.

194 Results

195 Sample demographics are provided in **Table 1**, with per center/kit demographics in **Table S5-6**. In
196 initial exome and genome-wide analyses using model-1, we observed many spurious associations
197 ($P \leq 1e-5$). We identified that variants underlying these spurious signals displayed increased variation in
198 allele frequency across sequencing centers/kits for the full frequency range (**Figure 1A-B**). We also
199 observed that such variants could not consistently be accounted for by adjustment for sequencing
200 center/kit in model-2; A specific example of such a variant is provided in **Figure 1C**.

201 Based on these observations, three versions of filters were designed and evaluated for their capacity
202 to capture putative center/kit-related variant artifacts. In assessing computing time, the filter using the
203 Fisher exact test implemented in Plink followed by Cauchy combination of P-values implemented in R
204 proved to be the fastest, taking 32 seconds to be constructed using a single CPU for a 1Mb region in
205 ADSP WGS (5,402 variants). Comparatively, constructing the chi square test filter implemented in R took
206 93 seconds, while the Fisher test with MC filter implemented in R took 128 seconds. Given the faster
207 speed, as well as the expected higher robustness provided by an exact test, we present the filter using
208 the Fisher exact test implemented in Plink as the primary filter, while the other two represent
209 supporting analyses. Throughout the remainder of the manuscript we will use the term “filtered” to
210 describe variants that were removed by filters and the term “non-filtered” to describe variants that
211 were not removed by filters.

212 The Fisher exact center/kit-based variant filters showed they strongly reduced the number of
213 spurious associations observed with model-1 in ADSP WES (**Figure 2A & 2C**) and WGS (**Figure 3A & 3C**).
214 When further adjusting for sequencing center/kit in model-2, spurious associations appeared essentially
215 absent in ADSP WES (**Figure 2D**) and WGS (**Figure 3D**). Notably, the spurious associations did not appear
216 to be driven by inflation, as for instance the genomic control factor (λ) was consistent prior to and after
217 applying variant filters in ADSP WGS for the respective models (**Figure 3**). The slightly larger λ for ADSP
218 WES in model-1 prior to applying the variant filters (**Figure 2A**) indicated that the large number of
219 spurious variants with regard to the relatively small total set of variants was likely driving some modest
220 inflation. Consistent observations were made for the other two center/kit-based variant filters (**Figure**
221 **S3-6**). When intersecting variants identified across these three sets of filters, the filter derived from the
222 fisher exact test implemented in Plink overlapped strongly (>96%) with the other two filters that in turn

223 showed less overlap (**Figure S7**). This was consistent with the Fisher exact test being the most
224 conservative and robust.

225 Closer inspection of the center/kit-based variant filters showed that non-filtered variants displayed
226 fairly concordant P-values across model-1 and model-2, whereas filtered variants showed many
227 discrepancies (**Figure 4A-B & 4D-E**). This was consistent with the filtered variants driving spurious
228 associations. Additionally, it was apparent that filters removed variants across the full frequency range
229 (**Figure 4C & 4F**) consistent with the increased MAF variation across all frequency ranges for variants
230 underlying spurious association signals (**Figure 1A-B**).

231 We then assessed to what extent the gnomAD-based filters could remove the observed spurious
232 associations. Visual assessment of Manhattan plots showed that the gnomAD-based filters could only
233 account for a part of the spurious associations (**Figure S8-9**). Similarly, closer inspection of the gnomAD-
234 based filters showed that they mainly removed variants with frequencies <1% (**Figure S10**). P-values
235 across model-1 and model-2 further showed many discrepancies both for non-filtered and filtered
236 variants (although fewer for non-filtered variants). In sum, the gnomAD-based filters could remove some
237 spurious signals, but were less effective than the center/kit-based variant filters.

238 We further assessed to what extent the duplicate discordant variants filters could remove the
239 observed spurious associations. Manhattan plots showed that the duplicate discordant variant filters
240 could account for many of the spurious associations, but several remained when using model-1, while
241 when using model-2 the Manhattan plots looked similar to those using the center-kit-based variant
242 filters (**Figure S11-12**). Closer inspection of the duplicate discordant variant filters similarly showed they
243 mainly removed variants with frequencies >10% and did not remove a set of variants that lose
244 suggestive significance when going from model-1 to model-2 (**Figure S13**). An illustrative example of
245 such a variant is provided in **Table S7**, confirming these variants represent genotyping issues that more
246 ideally should be removed from the data. In sum, the duplicate discordant filters could remove many
247 spurious signals, but were less effective than the center/kit-based variant filters, yet more effective than
248 the gnomAD-based variant filters.

249 We also sought to understand the overlap between the different proposed filters. The three
250 gnomAD-based variant filters appeared to show little overlap with one another (**Figure S14**) and
251 overlapped with less than 20% of the variants in the center/kit-based variant filters (**Figure S15**).
252 Further, in ADSP WES and WGS, respectively 32% and 14% of duplicate discordant variants overlapped

253 center/kit-based variant filters, while vice versa 31% and 15% of center/kit-based filtered variants
254 overlapped duplicate discordant variants (**Figure S16**). In the same comparison, respectively 28% and
255 49% of duplicate discordant variants overlapped gnomAD-based variant filters, while vice versa 53% and
256 17% of gnomAD-based filtered variants overlapped duplicate discordant variants (**Figure S17**). In sum,
257 this confirmed that all three types of filters captured overlapping as well as unique variants. Notably, the
258 center/kit- and gnomAD-based variant filters could capture a subset of reference artifactual variants
259 present in the duplicate discordant variant filters, but identified many additional signals that
260 represented likely artifactual variants and that contributed to spurious association signals.

261 Then, we sought to assess whether the use of these different types of variant filters could omit the
262 need for adjusting for sequencing center/kit as implemented in model-2, which may be desirable for
263 certain studies or research questions. We thus inspected all variants that passed suggestive significance
264 in either model-1 or model-2 in ADSP WES (**Table 2**) and WGS (**Table S8**) after applying the center/kit-
265 based filters (which we showed removed the most spurious signals). We observed that many variants
266 that lose suggestive significance after center/kit adjustment in model-2 have fairly small (above
267 threshold) P-values in the center/kit-based Fisher exact tests and/or are covered in the gnomAD-based
268 and duplicate discordant variant filters. Similarly, assessing Manhattan plots and variant metrics
269 suggested that the gnomAD-based and/or duplicate discordant variant filters removed few additional
270 variants underlying spurious signals (**Figure S18-23**). This suggests there may be added value in using
271 model-2 and/or applying the gnomAD-based filters to reduce spurious signals. Obviously, adding the
272 duplicate discordant variant filters will inherently remove artifactual signals and help reduce spurious
273 signals.

274 Lastly, as a robustness check, we compared association statistics from the current ADSP WES analyses
275 to variants that we identified in a prior study using a prior version of the ADSP WES data and observed
276 highly concordant findings (**Table S9**)¹⁵.

277 Discussion

278 We present a novel, fast, and robust approach to filter variants that represent sequencing center- or
279 kit-related artifacts underlying spurious associations with AD risk in ADSP WES and WGS data, that
280 cannot fully be accounted for by center/kit covariate adjustment. In addition, we show that filters
281 comprising variants that may be prone to artifacts as identified by gnomAD were less efficient in
282 removing spurious signals, but may still have added value on top of the center/kit-based filters.
283 Similarly, filters containing variants that were discordant across duplicate samples could remove many,
284 but not all, spurious signals, and added onto the center/kit-based filters. In sum, the presented filters
285 are important to support future robust studies on ADSP data. In addition, these filters allow flexibility
286 given that they can be applied in post-hoc quality control. Researchers may thus inspect filtered variants
287 in targeted analyses in subsets of the ADSP data where no artifactual genotype enrichment is observed
288 (e.g. excluding a single sequencing center/kit that showed an artifactual increase in genotype counts
289 compared to the others).

290 Certain study designs or research questions may benefit from not adjusting by sequencing center/kit
291 (i.e. cohort adjustment). For example, a study that considers specific strata and/or low frequency
292 variants may observe some co-linearity between variant genotype observations and sequencing
293 centers/kits. However, this does not necessarily indicate artifactual variants and may be driven by
294 chance or variable cohort study designs across samples sequenced by different centers. We observed
295 that the presented center/kit-based variant filters could handle nearly all spurious associations when
296 not adjusting for sequencing center/kit in model-1. Inspecting the remaining signals passing suggestive
297 significance, it was apparent that the gnomAD-based and duplicate discordant variant filters could
298 remove a few additional spurious signals. Similarly, the P-values from the Fisher exact tests across
299 sequencing centers/kits was fairly small for several variants that passed suggestive significance in their
300 association with AD risk in model-1 but lost suggestive significance upon center/kit adjustment in
301 model-2. In sum, we suggest that model-2 with application of center/kit-based, gnomAD-based, and
302 duplicate discordant variant filters is the most conservative approach, but model-1 using only center/kit-
303 based and duplicate discordant variant filters may reasonably be implemented, contingent on post-hoc
304 assessment of the association signals' robustness.

305 The center/kit-based filtering approach will further be valuable beyond the currently presented
306 exome- and genome-wide univariate AD risk association analyses in European ancestry samples.

307 Notably, removal of artifactual variants may lead to improved association statistics in gene-based
308 testing, which is particularly relevant for WES/WGS data⁷. The filter approach can also be applied to
309 non-European samples available in ADSP WES/WGS. Lastly, the approach to check for variant artifacts by
310 comparing genotype distributions across sequencing centers/kits may also be used in other studies with
311 a similar design as the current ADSP data. Notably, pre-processing of UK Biobank SNP array data has
312 already implemented a similar type of filter as the one we described here in order to remove variants
313 that may represent batch or array effects²⁸. In turn, the approach described here and applied to
314 WES/WGS data could also be applied to the large amount of SNP array data sets used in large-scale
315 genetic studies of AD³.

316 The current study is the first to report exome- and genome-wide AD risk association findings for the
317 newly released ADSP 20.5k WES and 16.9k WGS data. After quality control and filter implementation,
318 we observed few signals passing the genome-wide significance threshold. In the ADSP WES data, *TREM2*
319 and *ABCA7*—well-established AD risk genes^{2,6}—were observed with variants respectively at genome-
320 wide and suggestive significance, consistent with observations for similar models in prior studies on the
321 prior ADSP WES discovery phase data^{7,15}. Despite only observing 4 variants in ADSP WES that passed
322 suggestive significance in model-2, our findings were overall highly consistent with prior work¹⁵. We also
323 observed that certain variants identified previously were not present in our current summary statistics
324 (**Table S9**), reflecting differences in joint calling, quality control, and the fact that currently only bi-allelic
325 data were available for the new ADSP WES data. Notably, the common protective variant on *ABCA7*
326 identified here has not been previously reported (and we confirm it appears to not have been
327 successfully joint called in the prior ADSP WES data; dbGaP accession ID: phs000572). In the ADSP WGS
328 data, in addition to several suggestive hits, *BIN1*—a well-established AD risk gene^{2,6}—and *CNTN4* were
329 identified with variants at genome-wide significance. The common protective variant on *CNTN4* appears
330 novel and may be of relevance to AD pathogenesis given that Contactin 4 (CNTN4) is a binding partner of
331 Amyloid Precursor Protein (APP) and CNTN4/APP interaction may play a role in promoting target-specific
332 axon arborization^{29,30}. Overall, these initial findings appear promising but suggest that the current ADSP
333 WES/WGS data may still suffer from power limitations limiting discovery of novel risk variants. As such,
334 gene-based testing, analyses on available non-European ancestry samples, and novel methodological
335 approaches to gain additional power^{12,15}, will all be crucial to support future advances into disentangling
336 the missing heritability of AD using ADSP samples and other complimentary large-scale sequencing data.

337 **Conclusion**

338 We present a novel, fast, and robust approach to filter variants that represent sequencing center- or
339 kit-related artifacts underlying spurious associations with AD risk in ADSP WES and WGS data. This
340 approach will be important to support future robust studies on ADSP data, as well as other studies with
341 similar designs.

342 **Contributions**

343 M.E.B. performed data processing, performed data analyses, designed analyses, designed study, wrote
344 paper, and obtained funding. Y.L.G. performed data processing and designed analyses. S.J.E. performed
345 data processing. V.N. performed data processing and supervised work. M.D.G supervised analyses,
346 supervised work, and obtained funding. Z.H designed study, designed analyses, supervised analyses,
347 supervised work, wrote paper, and obtained funding.

348 **Declaration of interests**

349 The authors declare no competing interests.

350 **Data sharing statement**

351 All data used in the analyses are available upon application to:

- 352 - dbGaP (<https://www.ncbi.nlm.nih.gov/gap/>)
- 353 - NIAGADS (<https://www.niagads.org/>)
- 354 - LONI (<https://ida.loni.usc.edu/>)
- 355 - Synapse (<https://www.synapse.org/>)
- 356 - Rush (<https://www.radc.rush.edu/>)
- 357 - NACC (<https://nacccdata.org/>)

358 The specific data repository and identifier for each cohort is indicated in **Table S1-2** of the supplement.

359 Summary statistics from the current study will be available at <https://www.niagads.org/home/>

360 Acknowledgements

361 Funding for this study was provided by the The Iqbal Farrukh & Asad Jamal Fund, the NIH (AG060747
362 and AG047366, granted to M.D.G, AG066206 and AG066515 granted to Z.H), and the Alzheimer's
363 Association (AARF-20-683984, granted to M.E.B), the European Union's Horizon 2020 research and
364 innovation program under the Marie Skłodowska-Curie (grant agreement No. 890650, granted to Y.L.G).

365 Biological samples used in this study were stored at study investigators' institutions and at the National
366 Cell Repository for Alzheimer's Disease (NCRAD) at Indiana University, which receives government
367 support under a cooperative agreement grant (U24 AG21886) awarded by the National Institute on
368 Aging (NIA). We thank contributors who collected samples used in this study, as well as patients and
369 their families, whose help and participation made this work possible. Phenotypic data were provided by
370 principal investigators, the NIA funded Alzheimer's Disease Centers (ADCs), the National Alzheimer's
371 Coordinating Center (NACC, U01AG016976), and the National Institute on Aging Genetics of Alzheimer's
372 Disease Data Storage Site (NIAGADS, U24AG041689) at the University of Pennsylvania, funded by NIA.
373 Contributors to the Genetic Analysis Data included Study Investigators on projects that were individually
374 funded by NIA, and other NIH institutes, and by private U.S. organizations, or foreign governmental or
375 nongovernmental organizations.

376 Data for this study were prepared, archived, and distributed by the National Institute on Aging
377 Alzheimer's Disease Data Storage Site (NIAGADS) at the University of Pennsylvania (U24-AG041689-01);
378 Alzheimer's Disease Genetics Consortium (ADGC), U01 AG032984, RC2 AG036528; NACC, U01
379 AG016976; NIA-LOAD (Columbia University), U24 AG026395, U24 AG026390, R01AG041797; Banner Sun
380 Health Research Institute P30 AG019610; Boston University, P30 AG013846, U01 AG10483, R01
381 CA129769, R01 MH080295, R01 AG017173, R01 AG025259, R01 AG048927, R01AG33193, R01
382 AG009029; Columbia University, P50 AG008702, R37 AG015473, R01 AG037212, R01 AG028786; Duke
383 University, P30 AG028377, AG05128; Emory University, AG025688; Group Health Research Institute,
384 UO1 AG006781, UO1 HG004610, UO1 HG006375, UO1 HG008657; Indiana University, P30 AG10133, R01
385 AG009956, RC2 AG036650; Johns Hopkins University, P50 AG005146, R01 AG020688; Massachusetts
386 General Hospital, P50 AG005134; Mayo Clinic, P50 AG016574, R01 AG032990, KL2 RR024151; Mount
387 Sinai School of Medicine, P50 AG005138, P01 AG002219; New York University, P30 AG08051, UL1
388 RR029893, 5R01AG012101, 5R01AG022374, 5R01AG013616, 1RC2AG036502, 1R01AG035137; North
389 Carolina A&T University, P20 MD000546, R01 AG28786-01A1; Northwestern University, P30 AG013854;
390 Oregon Health & Science University, P30 AG008017, R01 AG026916; Rush University, P30 AG010161,

391 R01 AG019085, R01 AG15819, R01 AG17917, R01 AG030146, R01 AG01101, RC2 AG036650, R01
392 AG22018; TGEN, R01 NS059873; University of Alabama at Birmingham, P50 AG016582, UL1RR02777;
393 University of Arizona, R01 AG031581; University of California, Davis, P30 AG010129; University of
394 California, Irvine, P50 AG016573, P50 AG016575, P50 AG016576, P50 AG016577; University of
395 California, Los Angeles, P50 AG016570; University of California, San Diego, P50 AG005131; University of
396 California, San Francisco, P50 AG023501, P01 AG019724; University of Kentucky, P30 AG028383,
397 AG05144; University of Michigan, P30 AG053760 and AG063760; University of Pennsylvania, P30
398 AG010124; University of Pittsburgh, P50 AG005133, AG030653, AG041718, AG07562, AG02365;
399 University of Southern California, P50 AG005142; University of Texas Southwestern, P30 AG012300;
400 University of Miami, R01 AG027944, AG010491, AG027944, AG021547, AG019757; University of
401 Washington, P50 AG005136, R01 AG042437; University of Wisconsin, P50 AG033514; Vanderbilt
402 University, R01 AG019085; and Washington University, P50 AG005681, P01 AG03991, P01 AG026276.
403 The Kathleen Price Bryan Brain Bank at Duke University Medical Center is funded by NINDS grant #
404 NS39764, NIMH MH60451 and by Glaxo Smith Kline. Genotyping of the TGEN2 cohort was supported by
405 Kronos Science. The TGen series was also funded by NIA grant AG041232, The Banner Alzheimer's
406 Foundation, The Johnnie B. Byrd Sr. Alzheimer's Institute, the Medical Research Council, and the state of
407 Arizona and also includes samples from the following sites: Newcastle Brain Tissue Resource (funding via
408 the Medical Research Council, local NHS trusts and Newcastle University), MRC London Brain Bank for
409 Neurodegenerative Diseases (funding via the Medical Research Council), South West Dementia Brain
410 Bank (funding via numerous sources including the Higher Education Funding Council for England
411 (HEFCE), Alzheimer's Research Trust (ART), BRACE as well as North Bristol NHS Trust Research and
412 Innovation 58 Department and DeNDROn), The Netherlands Brain Bank (funding via numerous sources
413 including Stichting MS Research, Brain Net Europe, Hersenstichting Nederland Breinbrekend Werk,
414 International Parkinson Fonds, Internationale Stichting Alzheimer Onderzoek), Institut de
415 Neuropatologia, Servei Anatomia Patologica, Universitat de Barcelona.
416 The NACC database is funded by NIA/NIH Grant U01 AG016976. NACC data are contributed by the NIA-
417 funded ADCs: P30 AG019610 (PI Eric Reiman, MD), P30 AG013846 (PI Neil Kowall, MD), P30 AG062428-
418 01 (PI James Leverenz, MD) P50 AG008702 (PI Scott Small, MD), P50 AG025688 (PI Allan Levey, MD,
419 PhD), P50 AG047266 (PI Todd Golde, MD, PhD), P30 AG010133 (PI Andrew Saykin, PsyD), P50 AG005146
420 (PI Marilyn Albert, PhD), P30 AG062421-01 (PI Bradley Hyman, MD, PhD), P30 AG062422-01 (PI Ronald
421 Petersen, MD, PhD), P50 AG005138 (PI Mary Sano, PhD), P30 AG008051 (PI Thomas Wisniewski, MD),
422 P30 AG013854 (PI Robert Vassar, PhD), P30 AG008017 (PI Jeffrey Kaye, MD), P30 AG010161 (PI David

423 Bennett, MD), P50 AG047366 (PI Victor Henderson, MD, MS), P30 AG010129 (PI Charles DeCarli, MD),
424 P50 AG016573 (PI Frank LaFerla, PhD), P30 AG062429-01(PI James Brewer, MD, PhD), P50 AG023501 (PI
425 Bruce Miller, MD), P30 AG035982 (PI Russell Swerdlow, MD), P30 AG028383 (PI Linda Van Eldik, PhD),
426 P30 AG053760 (PI Henry Paulson, MD, PhD), P30 AG010124 (PI John Trojanowski, MD, PhD), P50
427 AG005133 (PI Oscar Lopez, MD), P50 AG005142 (PI Helena Chui, MD), P30 AG012300 (PI Roger
428 Rosenberg, MD), P30 AG049638 (PI Suzanne Craft, PhD), P50 AG005136 (PI Thomas Grabowski, MD),
429 P30 AG062715-01 (PI Sanjay Asthana, MD, FRCP), P50 AG005681 (PI John Morris, MD), P50 AG047270
430 (PI Stephen Strittmatter, MD, PhD).

431 The genotypic and associated phenotypic data used in the study “Multi-Site Collaborative Study for
432 Genotype-Phenotype Associations in Alzheimer’s Disease (GenADA)” were provided by the
433 GlaxoSmithKline, R&D Limited.

434 ROSMAP study data were provided by the Rush Alzheimer’s Disease Center, Rush University Medical
435 Center, Chicago. Data collection was supported through funding by NIA grants P30AG10161,
436 R01AG15819, R01AG17917, R01AG30146, R01AG36836, U01AG32984, U01AG46152, the Illinois
437 Department of Public Health, and the Translational Genomics Research Institute.

438 The AddNeuroMed data are from a public-private partnership supported by EFPIA companies and SMEs
439 as part of InnoMed (Innovative Medicines in Europe), an Integrated Project funded by the European
440 Union of the Sixth Framework program priority FP6-2004-LIFESCIHEALTH-5. Clinical leads responsible for
441 data collection are Iwona Kloszewska (Lodz), Simon Lovestone (London), Patrizia Mecocci (Perugia),
442 Hilikka Soininen (Kuopio), Magda Tsolaki (Thessaloniki), and Bruno Vellas (Toulouse), imaging leads are
443 Andy Simmons (London), Lars-Olad Wahlund (Stockholm) and Christian Spenger (Zurich) and
444 bioinformatics leads are Richard Dobson (London) and Stephen Newhouse (London).

445 Data collection and sharing for this project was funded by the Alzheimer's Disease Neuroimaging
446 Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI (Department of
447 Defense award number W81XWH-12-2-0012). ADNI is funded by the National Institute on Aging, the
448 National Institute of Biomedical Imaging and Bioengineering and through generous contributions from
449 the following: AbbVie. Alzheimer’s Association; Alzheimer’s Drug Discovery Foundation; Araclon Biotech;
450 BioClinica. Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir. Inc.; Cogstate; Eisai Inc.; Elan
451 Pharmaceuticals. Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated
452 company Genentech. Inc.; Fujirebio; GE HealthControlsare; IXICO Ltd.; Janssen Alzheimer Immunotherapy
453 Research & Development. LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.;
454 Lumosity; Lundbeck; Merck & Co. Inc.; Meso Scale Diagnostics. LLC.; NeuroRx Research; Neurotrack

455 Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda
456 Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes of Health Research is
457 providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by
458 the Foundation for the National Institutes of Health. The grantee organization is the Northern California
459 Institute for Research and Education, and the study is coordinated by the Alzheimer's Therapeutic
460 Research Institute at the University of Southern California. ADNI data are disseminated by the
461 Laboratory for Neuro Imaging at the University of Southern California.

462 The Alzheimer's Disease Sequencing Project (ADSP) is comprised of two Alzheimer's Disease (AD)
463 genetics consortia and three National Human Genome Research Institute (NHGRI) funded Large Scale
464 Sequencing and Analysis Centers (LSAC). The two AD genetics consortia are the Alzheimer's Disease
465 Genetics Consortium (ADGC) funded by NIA (U01 AG032984), and the Cohorts for Heart and Aging
466 Research in Genomic Epidemiology (CHARGE) funded by NIA (R01 AG033193), the National Heart, Lung,
467 and Blood Institute (NHLBI), other National Institute of Health (NIH) institutes and other foreign
468 governmental and non-governmental organizations. The Discovery Phase analysis of sequence data is
469 supported through UF1AG047133 (to Drs. Schellenberg, Farrer, Pericak-Vance, Mayeux, and Haines);
470 U01AG049505 to Dr. Seshadri; U01AG049506 to Dr. Boerwinkle; U01AG049507 to Dr. Wijsman; and
471 U01AG049508 to Dr. Goate and the Discovery Extension Phase analysis is supported through
472 U01AG052411 to Dr. Goate, U01AG052410 to Dr. Pericak-Vance and U01 AG052409 to Drs. Seshadri and
473 Fornage.

474 The ADGC cohorts included in ADSP include: Adult Changes in Thought (ACT) (UO1 AG006781, UO1
475 HG004610, UO1 HG006375, UO1 HG008657), the Alzheimer's Disease Centers (ADC) (P30 AG019610,
476 P30 AG013846, P50 AG008702, P50 AG025688, P50 AG047266, P30 AG010133, P50 AG005146, P50
477 AG005134, P50 AG016574, P50 AG005138, P30 AG008051, P30 AG013854, P30 AG008017, P30
478 AG010161, P50 AG047366, P30 AG010129, P50 AG016573, P50 AG016570, P50 AG005131, P50
479 AG023501, P30 AG035982, P30 AG028383, P30 AG010124, P50 AG005133, P50 AG005142, P30
480 AG012300, P50 AG005136, P50 AG033514, P50 AG005681, and P50 AG047270), the Chicago Health and
481 Aging Project (CHAP) (R01 AG11101, RC4 AG039085, K23 AG030944), Indianapolis Ibadan (R01
482 AG009956, P30 AG010133), the Memory and Aging Project (MAP) (R01 AG17917), Mayo Clinic (MAYO)
483 (R01 AG032990, U01 AG046139, R01 NS080820, RF1 AG051504, P50 AG016574), Mayo Parkinson's
484 Disease controls (NS039764, NS071674, 5RC2HG005605), University of Miami (R01 AG027944, R01
485 AG028786, R01 AG019085, IIRG09133827, A2011048), the Multi-Institutional Research in Alzheimer's
486 Genetic Epidemiology Study (MIRAGE) (R01 AG09029, R01 AG025259), the National Cell Repository for

487 Alzheimer's Disease (NCRAD) (U24 AG21886), the National Institute on Aging Late Onset Alzheimer's
488 Disease Family Study (NIA- LOAD) (R01 AG041797), the Religious Orders Study (ROS) (P30 AG10161, R01
489 AG15819), the Texas Alzheimer's Research and Care Consortium (TARCC) (funded by the Darrell K Royal
490 Texas Alzheimer's Initiative), Vanderbilt University/Case Western Reserve University (VAN/CWRU) (R01
491 AG019757, R01 AG021547, R01 AG027944, R01 AG028786, P01 NS026630, and Alzheimer's Association),
492 the Washington Heights-Inwood Columbia Aging Project (WHICAP) (RF1 AG054023), the University of
493 Washington Families (VA Research Merit Grant, NIA: P50AG005136, R01AG041797, NINDS:
494 R01NS069719), the Columbia University Hispanic Estudio Familiar de Influencia Genetica de Alzheimer
495 (EFIGA) (RF1 AG015473), the University of Toronto (UT) (funded by Wellcome Trust, Medical Research
496 Council, Canadian Institutes of Health Research), and Genetic Differences (GD) (R01 AG007584). The
497 CHARGE cohorts are supported in part by National Heart, Lung, and Blood Institute (NHLBI)
498 infrastructure grant HL105756 (Psaty), RC2HL102419 (Boerwinkle) and the neurology working group is
499 supported by the National Institute on Aging (NIA) R01 grant AG033193.

500 The CHARGE cohorts participating in the ADSP include the following: Austrian Stroke Prevention Study
501 (ASPS), ASPS-Family study, and the Prospective Dementia Registry-Austria (ASPS/PRODEM-Aus), the
502 Atherosclerosis Risk in Communities (ARIC) Study, the Cardiovascular Health Study (CHS), the Erasmus
503 Rucphen Family Study (ERF), the Framingham Heart Study (FHS), and the Rotterdam Study (RS). ASPS is
504 funded by the Austrian Science Fond (FWF) grant number P20545-P05 and P13180 and the Medical
505 University of Graz. The ASPS-Fam is funded by the Austrian Science Fund (FWF) project I904), the EU
506 Joint Programme - Neurodegenerative Disease Research (JPND) in frame of the BRIDGET project
507 (Austria, Ministry of Science) and the Medical University of Graz and the Steiermärkische
508 Krankenanstalten Gesellschaft. PRODEM-Austria is supported by the Austrian Research Promotion
509 agency (FFG) (Project No. 827462) and by the Austrian National Bank (Anniversary Fund, project 15435.
510 ARIC research is carried out as a collaborative study supported by NHLBI contracts
511 (HHSN268201100005C, HHSN268201100006C, HHSN268201100007C, HHSN268201100008C,
512 HHSN268201100009C, HHSN268201100010C, HHSN268201100011C, and HHSN268201100012C).
513 Neurocognitive data in ARIC is collected by U01 2U01HL096812, 2U01HL096814, 2U01HL096899,
514 2U01HL096902, 2U01HL096917 from the NIH (NHLBI, NINDS, NIA and NIDCD), and with previous brain
515 MRI examinations funded by R01-HL70825 from the NHLBI. CHS research was supported by contracts
516 HHSN268201200036C, HHSN268200800007C, N01HC55222, N01HC85079, N01HC85080, N01HC85081,
517 N01HC85082, N01HC85083, N01HC85086, and grants U01HL080295 and U01HL130114 from the NHLBI
518 with additional contribution from the National Institute of Neurological Disorders and Stroke (NINDS).

519 Additional support was provided by R01AG023629, R01AG15928, and R01AG20098 from the NIA. FHS
520 research is supported by NHLBI contracts N01-HC-25195 and HHSN268201500001I. This study was also
521 supported by additional grants from the NIA (R01s AG054076, AG049607 and AG033040 and NINDS
522 (R01 NS017950). The ERF study as a part of EUROSPAN (European Special Populations Research
523 Network) was supported by European Commission FP6 STRP grant number 018947 (LSHG-CT-2006-
524 01947) and also received funding from the European Community's Seventh Framework Programme
525 (FP7/2007-2013)/grant agreement HEALTH-F4- 2007-201413 by the European Commission under the
526 programme "Quality of Life and Management of the Living Resources" of 5th Framework Programme
527 (no. QL2-CT-2002- 01254). High-throughput analysis of the ERF data was supported by a joint grant
528 from the Netherlands Organization for Scientific Research and the Russian Foundation for Basic
529 Research (NWO-RFBR 047.017.043). The Rotterdam Study is funded by Erasmus Medical Center and
530 Erasmus University, Rotterdam, the Netherlands Organization for Health Research and Development
531 (ZonMw), the Research Institute for Diseases in the Elderly (RIDE), the Ministry of Education, Culture and
532 Science, the Ministry for Health, Welfare and Sports, the European Commission (DG XII), and the
533 municipality of Rotterdam. Genetic data sets are also supported by the Netherlands Organization of
534 Scientific Research NWO Investments (175.010.2005.011, 911-03-012), the Genetic Laboratory of the
535 Department of Internal Medicine, Erasmus MC, the Research Institute for Diseases in the Elderly (014-
536 93-015; RIDE2), and the Netherlands Genomics Initiative (NGI)/Netherlands Organization for Scientific
537 Research (NWO) Netherlands Consortium for Healthy Aging (NCHA), project 050-060-810. All studies are
538 grateful to their participants, faculty and staff. The content of these manuscripts is solely the
539 responsibility of the authors and does not necessarily represent the official views of the National
540 Institutes of Health or the U.S. Department of Health and Human Services.

541 The FUS cohorts include: the Alzheimer's Disease Centers (ADC) (P30 AG019610, P30 AG013846, P50
542 AG008702, P50 AG025688, P50 AG047266, P30 AG010133, P50 AG005146, P50 AG005134, P50
543 AG016574, P50 AG005138, P30 AG008051, P30 AG013854, P30 AG008017, P30 AG010161, P50
544 AG047366, P30 AG010129, P50 AG016573, P50 AG016570, P50 AG005131, P50 AG023501, P30
545 AG035982, P30 AG028383, P30 AG010124, P50 AG005133, P50 AG005142, P30 AG012300, P50
546 AG005136, P50 AG033514, P50 AG005681, and P50 AG047270), Alzheimer's Disease Neuroimaging
547 Initiative (ADNI) (U19AG024904), Amish Protective Variant Study (RF1AG058066), Cache County Study
548 (R01AG11380, R01AG031272, R01AG21136, RF1AG054052), Case Western Reserve University Brain
549 Bank (CWRUBB) (P50AG008012), Case Western Reserve University Rapid Decline (CWRURD)
550 (RF1AG058267, NU38CK000480), CubanAmerican Alzheimer's Disease Initiative (CuAADI)

551 (3U01AG052410), Estudio Familiar de Influencia Genetica en Alzheimer (EFIGA) (5R37AG015473,
552 RF1AG015473, R56AG051876), Genetic and Environmental Risk Factors for Alzheimer Disease Among
553 African Americans Study (GenerAAtions) (2R01AG09029, R01AG025259, 2R01AG048927), Gwangju
554 Alzheimer and Related Dementias Study (GARD) (U01AG062602), Hussman Institute for Human
555 Genomics Brain Bank (HIHGBB) (R01AG027944, Alzheimer's Association "Identification of Rare Variants
556 in Alzheimer Disease"), Ibadan Study of Aging (IBADAN) (5R01AG009956), Mexican Health and Aging
557 Study (MHAS) (R01AG018016), Multi-Institutional Research in Alzheimer's Genetic Epidemiology
558 (MIRAGE) (2R01AG09029, R01AG025259, 2R01AG048927), Northern Manhattan Study (NOMAS)
559 (R01NS29993), Peru Alzheimer's Disease Initiative (PeADI) (RF1AG054074), Puerto Rican 1066 (PR1066)
560 (Wellcome Trust (GR066133/GR080002), European Research Council (340755)), Puerto Rican Alzheimer
561 Disease Initiative (PRADI) (RF1AG054074), Reasons for Geographic and Racial Differences in Stroke
562 (REGARDS) (U01NS041588), Research in African American Alzheimer Disease Initiative (REAAADI)
563 (U01AG052410), Rush Alzheimer's Disease Center (ROSMAP) (P30AG10161, R01AG15819,
564 R01AG17919), University of Miami Brain Endowment Bank (MBB), and University of Miami/Case
565 Western/North Carolina A&T African American (UM/CASE/NCAT) (U01AG052410, R01AG028786).
566 The four LSACs are: the Human Genome Sequencing Center at the Baylor College of Medicine (U54
567 HG003273), the Broad Institute Genome Center (U54HG003067), The American Genome Center at the
568 Uniformed Services University of the Health Sciences (U01AG057659), and the Washington University
569 Genome Institute (U54HG003079).

570 References

- 571 1 Sierksma A, Escott-Price V, De Strooper B. Translating genetic risk of Alzheimer's disease into
572 mechanistic insight and drug targets. *Science* 2020; **370**: 61–6.
- 573 2 Sims R, Hill M, Williams J. The multiplex model of the genetics of Alzheimer's disease. *Nat*
574 *Neurosci* 2020; **23**: 311–322.
- 575 3 Kunkle BW, Grenier-Boley B, Sims R, *et al.* Genetic meta-analysis of diagnosed Alzheimer's
576 disease identifies new risk loci and implicates A β , tau, immunity and lipid processing. *Nat Genet*
577 2019; **51**: 414–30.
- 578 4 Jansen IE, Savage JE, Watanabe K, *et al.* Genome-wide meta-analysis identifies new loci and
579 functional pathways influencing Alzheimer's disease risk. *Nat Genet* 2019; **51**: 404–13.
- 580 5 de Rojas I, Moreno-Grau S, Tesi N, *et al.* Common variants in Alzheimer's disease and risk
581 stratification by polygenic risk scores. *Nat Commun* 2021; **12**. DOI:10.1038/s41467-021-22491-8.
- 582 6 Andrews SJ, Fulton-Howard B, Goate A. Interpretation of risk loci from genome-wide association
583 studies of Alzheimer's disease. *Lancet Neurol* 2020; **19**: 326–35.
- 584 7 Bis JC, Jian X, Chen BWK, *et al.* Whole exome sequencing study identifies novel rare and common
585 Alzheimer's-Associated variants involved in immune response and transcriptional regulation. *Mol*
586 *Psychiatry* 2020; **25**: 1859–75.
- 587 8 Patel D, Mez J, Vardarajan BN, *et al.* Association of Rare Coding Mutations With Alzheimer
588 Disease and Other Dementias Among Adults of European Ancestry. *JAMA Netw open* 2019; **2**:
589 e191350.
- 590 9 Ma Y, Jun GR, Zhang X, *et al.* Analysis of Whole-Exome Sequencing Data for Alzheimer Disease
591 Stratified by APOE Genotype. *JAMA Neurol* 2019; **76**: 1099–108.
- 592 10 Blue EE, Thornton TA, Kooperberg C, *et al.* Non-coding variants in MYH11, FZD3, and SORCS3 are
593 associated with dementia in women. *Alzheimer's Dement* 2021; **17**: 215–25.
- 594 11 Park JH, Park I, Youm EM, *et al.* Novel Alzheimer's disease risk variants identified based on whole-
595 genome sequencing of APOE ϵ 4 carriers. *Transl Psychiatry* 2021; **11**. DOI:10.1038/s41398-021-
596 01412-9.
- 597 12 He Z, Liu L, Wang C, *et al.* Identification of putative causal loci in whole-genome sequencing data
598 via knockoff statistics. *Nat Commun* 2021; **12**: 3152.
- 599 13 He L, Loika Y, Park Y, Bennett DA, Kellis M, Kulminski AM. Exome-wide age-of-onset analysis
600 reveals exonic variants in ERN1, TACR3 and SPPL2C associated with Alzheimer's disease. *Transl*
601 *Psychiatry* 2021; **11**: 146.
- 602 14 Prokopenko D, Morgan SL, Mullin K, *et al.* Whole-genome sequencing reveals new Alzheimer's
603 disease – associated rare variants in loci related to synaptic function and neuronal development.
604 *Alzheimer's Dement J Alzheimer's Assoc* 2021; **17**: 1509–27.
- 605 15 Le Guen Y, Belloy ME, Napolioni V, *et al.* A novel age-informed approach for genetic association

- 606 analysis in Alzheimer's disease. *Alzheimers Res Ther* 2021; **13**.
607 <http://medrxiv.org/content/early/2021/01/06/2021.01.05.21249292.abstract>.
- 608 16 Beecham GW, Bis JC, Martin ER, *et al*. The Alzheimer's disease sequencing project: Study design
609 and sample selection. *Neurol Genet* 2017; **3**: e194.
- 610 17 Crane PK, Foroud T, Montine TJ, Larson EB. Alzheimer's Disease Sequencing Project Discovery
611 and Replication criteria for cases and controls: data from a community-based prospective cohort
612 study with autopsy follow-up. *Alzheimers Dement* 2017; **13**: 1410–3.
- 613 18 NIAGADS. NG00067 – ADSP Umbrella. 2021. <https://dss.niagads.org/datasets/ng00067/>.
- 614 19 Wickland DP, Ren Y, Sinnwell JP, *et al*. Impact of variant-level batch effects on identification of
615 genetic risk factors in large sequencing studies. *PLoS One* 2021; **16**: e0249305.
- 616 20 Karczewski KJ, Francioli LC, Tiao G, *et al*. The mutational constraint spectrum quantified from
617 variation in 141,456 humans. *Nature* 2020; **581**: 434–43.
- 618 21 Leung YY, Valladares O, Chou YF, *et al*. VCPA: Genomic variant calling pipeline and data
619 management tool for Alzheimer's Disease Sequencing Project. *Bioinformatics* 2019; **35**: 1768–70.
- 620 22 GATK team. GATK Best Practices Workflows. [https://gatk.broadinstitute.org/hc/en-](https://gatk.broadinstitute.org/hc/en-us/articles/360035894751)
621 [us/articles/360035894751](https://gatk.broadinstitute.org/hc/en-us/articles/360035894751) (accessed Feb 1, 2021).
- 622 23 Chen CY, Pollack S, Hunter DJ, Hirschhorn JN, Kraft P, Price AL. Improved ancestry inference using
623 weights from external reference panels. *Bioinformatics* 2013; **29**: 1399–406.
- 624 24 Conomos MP, Miller MB, Thornton TA. Robust inference of population structure for ancestry
625 prediction and correction of stratification in the presence of relatedness. *Genet Epidemiol* 2015;
626 **39**: 276–93.
- 627 25 Liu Y, Xie J. Cauchy Combination Test: A Powerful Test With Analytic p-Value Calculation Under
628 Arbitrary Dependency Structures. *J Am Stat Assoc* 2020; **115**: 393–402.
- 629 26 Sun Y, Wu S, Bu G, *et al*. Glial Fibrillary Acidic Protein – Apolipoprotein E (apoE) Transgenic Mice:
630 Astrocyte-Specific Expression and Differing Biological Effects of Astrocyte-Secreted apoE3 and
631 apoE4 Lipoproteins. *J Neurosci* 1998; **18**: 3261–72.
- 632 27 Yizhar O, Fenno L, Zhang F, Hegemann P, Diesseroth K. Microbial opsins: A family of single-
633 component tools for optical control of neural activity. *Cold Spring Harb Protoc* 2011; **6**.
634 DOI:10.1101/pdb.top102.
- 635 28 Bycroft C, Freeman C, Petkova D, *et al*. The UK Biobank resource with deep phenotyping and
636 genomic data. *Nature* 2018; **562**: 203–9.
- 637 29 Osterfield M, Egelund R, Young LM, Flanagan JG. Interaction of amyloid precursor protein with
638 contactins and NgCAM in the retinotectal system. *Dev Dis* 2008; **135**: 1189–99.
- 639 30 Osterhout JA, Stafford BK, Yoshihara Y, *et al*. Functional Development of the Accessory Optic
640 Article Contactin-4 Mediates Axon-Target Specificity and Functional Development. *Neuron* 2015;
641 **86**: 985–99.
- 642

643 **Table 1. Sample demographics.** Samples were restricted to those passing genetic/phenotypic quality
 644 control, being non-Hispanic, and being of European ancestry.

Samples		Diagnosis		Sex	Age	APOE status	
Name	Participants after QC (N)	Type	(N)	Female (N (%))	Age (Mean (SD))	APOE *4-pos	APOE *2-pos
ADSP WES	11573	CN	5418	3152 (58.2 %)	85.4 (6.5)	926 (17.1 %)	1057 (19.5 %)
		AD	6155	3619 (58.8 %)	75.4 (8.6)	2938 (47.7 %)	493 (8.0 %)
ADSP WGS	6533	CN	2949	1791 (60.7 %)	81.6 (6.6)	1075 (36.4 %)	204 (6.9 %)
		AD	3584	2051 (57.2 %)	76.7 (8.3)	2078 (58.0 %)	177 (4.9 %)

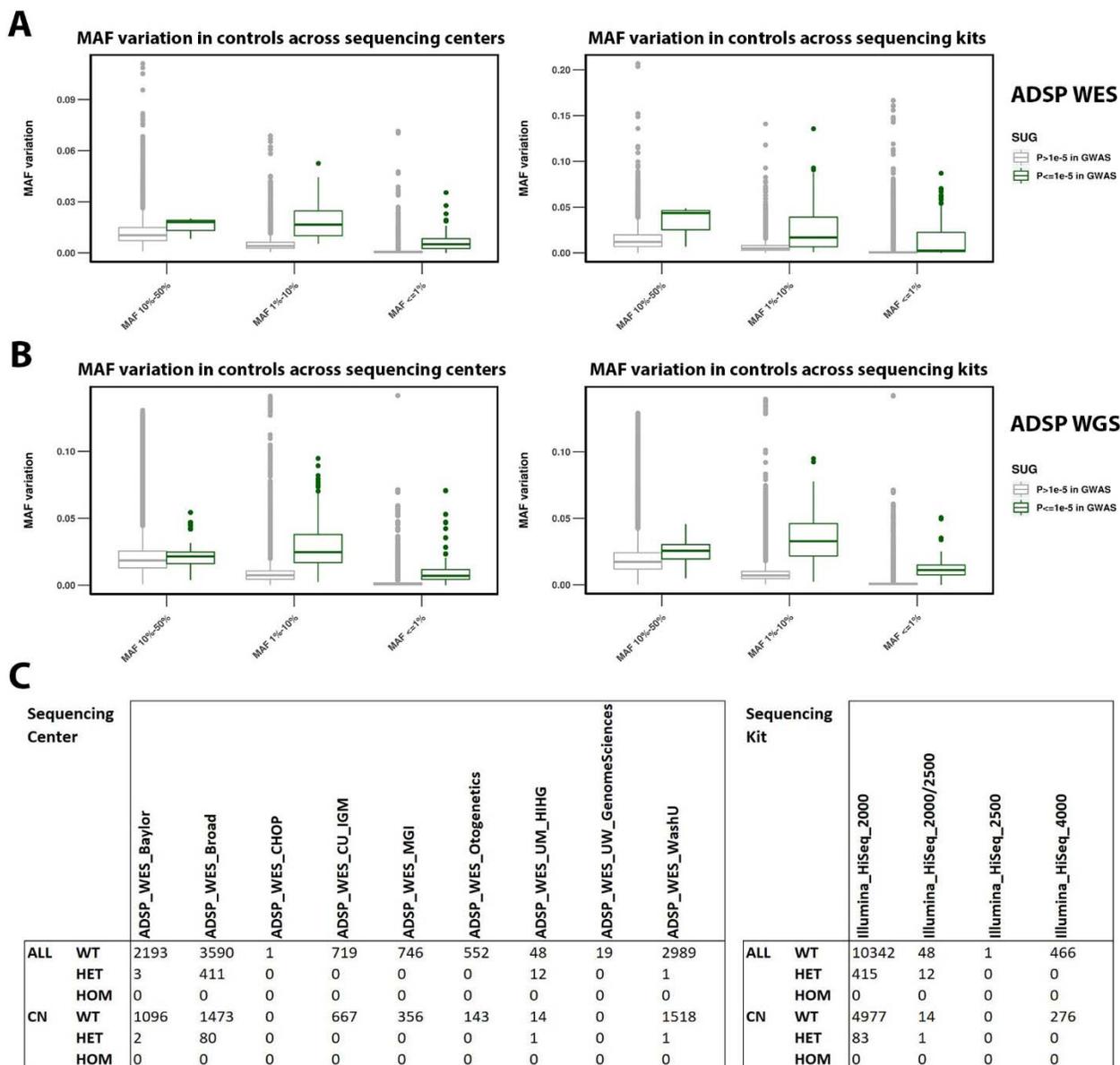
645

646 *Abbreviations: QC, quality control; CN, cognitively normal; AD, Alzheimer’s disease; SD, standard*
 647 *deviation*

648 **Table 2. ADSP WES variants passing suggestive significance after applying centers/kit-based filters.** Variants shown passed suggestive
649 significance in either model-1 or model-2. Note that many variants that lose suggestive significance after center/kit adjustment in model-2 have
650 fairly small P-values (but above threshold) in the center/kit Fisher tests and/or have a non-PASS flag in gnomAD or are flagged by the duplicate
651 discordant variant filter. This suggests there is added value in using model 2 and/or applying the gnomAD and duplicate discordant variant filters
652 to reduce spurious signals, or, that model 1 without gnomAD filters can be used contingent on post-hoc assessment of the association signal's
653 robustness.

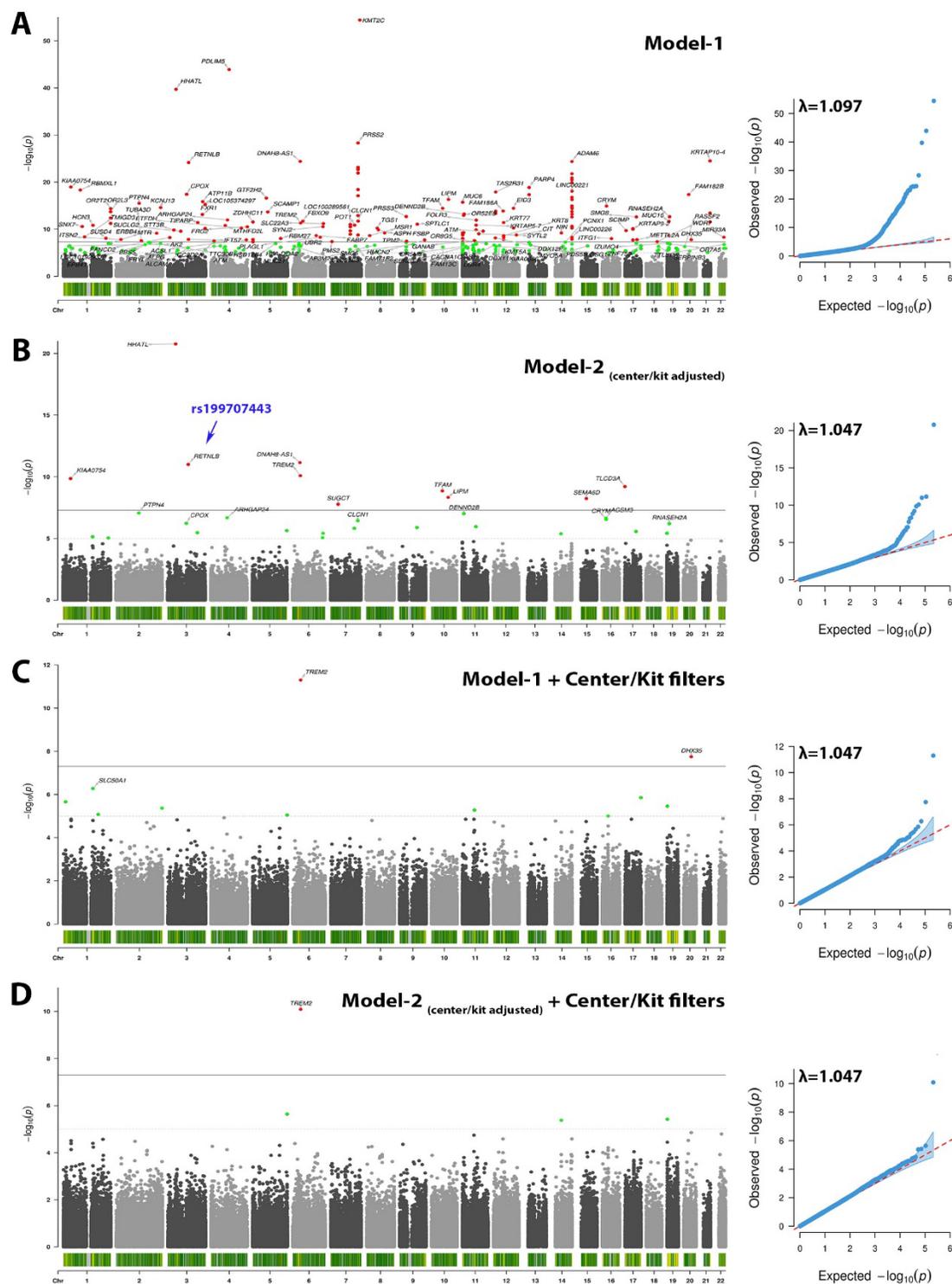
Variant info							Model-1				Model-2				Filters			
GENE	CHR	BP	dbSNP153 ID	effect allele	other allele	effect allele frq.	OR	95% CI (lb)	95% CI (ub)	P	OR	95% CI (lb)	95% CI (ub)	P	Center Fisher P	Kit Fisher P	gnomAD filter	Duplicate check
DRAXIN	1	11709400	rs769650621	C	A	0.45%	2.44	1.69	3.53	2.2E-06	1.83	1.28	2.63	1.0E-03	2.2E-04	0.97	non-PASS	discordant
SLC50A1	1	155136277	rs765315443	C	T	0.37%	2.82	1.88	4.22	5.3E-07	2.06	1.38	3.05	3.6E-04	2.3E-05	0.99	PASS	ok
LAMC1-AS1	1	183135182	rs1385675950	A	C	0.18%	3.69	2.08	6.56	8.4E-06	2.70	1.54	4.72	5.2E-04	0.10	0.99	PASS	discordant
LOC150935	2	239780397	rs1355381797	C	A	0.15%	4.39	2.34	8.25	4.3E-06	3.17	1.71	5.86	2.4E-04	0.99	0.99	PASS	ok
RASGEF1C	5	180127602	rs57288534	T	C	18.67%	0.87	0.81	0.92	9.0E-06	0.86	0.81	0.92	2.3E-06	0.76	0.34	PASS	ok
TREM2	6	41161514	rs75932628	T	C	0.69%	2.82	2.10	3.79	5.0E-12	2.58	1.94	3.44	8.2E-11	0.02	0.09	PASS	ok
HNRNPUL2-BSCL2	11	62724366	rs772898628	C	A	0.24%	3.14	1.92	5.13	5.3E-06	2.21	1.37	3.57	1.2E-03	6.1E-03	0.98	non-PASS	discordant
CDKL1	14	50390164	rs61981931	T	C	4.82%	0.79	0.71	0.89	5.6E-05	0.77	0.69	0.86	4.2E-06	1.1E-05	0.02	PASS	ok
C16orf92	16	30025807	rs11544328	C	A	46.34%	0.89	0.85	0.94	1.0E-05	0.91	0.87	0.95	8.0E-05	0.83	1.7E-03	PASS	discordant
ZNF750	17	82831739	rs751362098	G	A	0.39%	2.61	1.77	3.85	1.4E-06	2.02	1.38	2.96	2.8E-04	1.5E-04	0.99	non-PASS	discordant
ABCA7	19	1042810	rs3764645	G	A	46.71%	0.89	0.85	0.94	3.5E-06	0.89	0.85	0.94	3.8E-06	0.22	3.9E-04	PASS	ok
DHX35	20	39018815	rs779184241	A	G	0.32%	3.49	2.26	5.38	1.8E-08	2.56	1.67	3.91	1.4E-05	2.5E-04	0.99	PASS	discordant

654

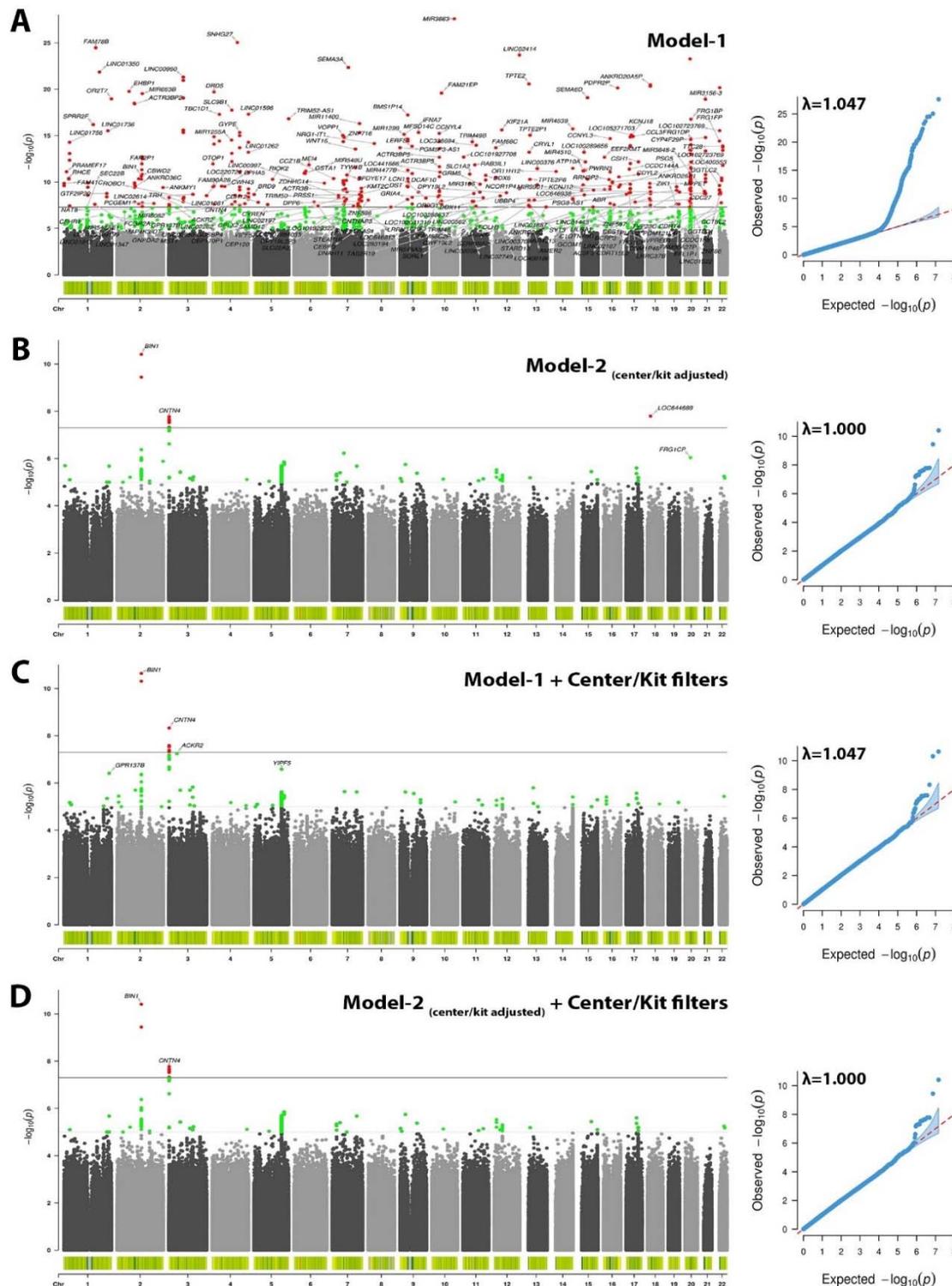


655

656 **Figure 1. Variant artifacts across different sequencing centers/kits drive spurious associations in ADSP**
 657 **WES and WGS data.** In initial exome-wide and genome-wide association studies of ADSP WES and WGS,
 658 we observed many spurious associations ($P \leq 1e-5$) using model-1 (i.e. not adjusting for sequencing
 659 center/kit; cf. **Figure 2A & 3A**). Upon inspection of these signals, it was notable that these variants
 660 displayed large variation in genotype counts across sequencing centers/kits. The MAF variation in
 661 control subjects for all analyzed variants is visualized in **A**) for ADSP WES, and in **B**) for ADSP WGS. **C**) A
 662 specific example of a variant showing spurious association is provided. This variant, rs199707443, has a
 663 MAF of 0.003% in non-Finnish Europeans in gnomAD v3.1.1, contrasting the 411 heterozygote counts in
 664 the Broad sequencing center. Notably, this particular variant still showed genome-wide significant
 665 association with AD risk even after sequencing center/kit adjustment (cf. **Figure 2B**). *Abbreviations:*
 666 *MAF, minor allele frequency; CN, cognitively normal; WT, Wild type; HET, heterozygote; HOM,*
 667 *homozygote.*



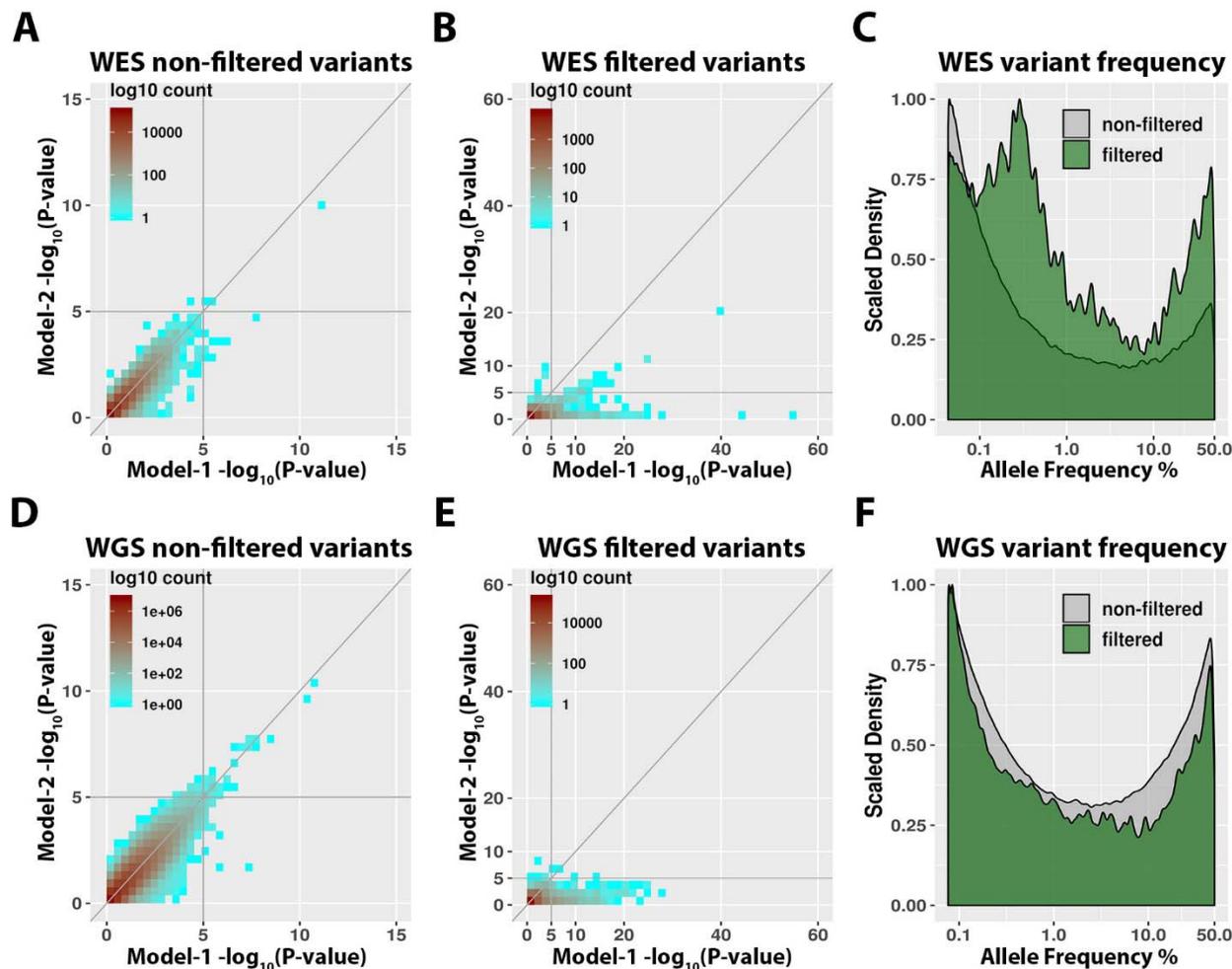
668
 669 **Figure 2. The proposed center/kit-based variant filters remove spurious associations in ADSP WES.**
 670 Figure shows Manhattan (left) and quantile-quantile (right) plots. **A)** Model-1 indicates many spurious
 671 hits. **B)** Model-2 shows that adjustment for center/kit can reduce many, but not all, spurious hits. The
 672 variant described in **Figure 1C** is highlighted by the blue arrow. **C)** Filters remove most spurious hits. **D)**
 673 Further adjustment for center/kit removes few additional spurious hits.



674

675 **Figure 3. The proposed center/kit-based variant filters remove spurious associations in ADSP WGS.**
 676 Figure shows Manhattan (left) and quantile-quantile (right) plots. **A)** Model-1 indicates many spurious
 677 hits. **B)** Model-2 shows that adjustment for center/kit can reduce many, but not all, spurious hits. **C)**

678 Filters remove most spurious hits. **D)** Further adjustment for center/kit removes few additional spurious
679 hits.



680
681 **Figure 4. Metrics of variants removed by the proposed center/kit-based variant filters. A-C)** ADSP WES.
682 **D-E)** ADSP WGS. **A & D)** Variants that passed filters showed largely consistent P-values across model-1
683 and model-2 case-control association analyses, with only few variants remaining that reach suggestive
684 significance in model-1 but lose suggestive significance upon center/kit adjustment in model-2 (lower
685 right quadrant). **B & E)** Variants that were removed by filters showed many inconsistent P-values across
686 model-1 and model-2, consistent with center/kit-related variant artifacts that could not fully be
687 accounted for by model-2. **C & F)** Frequency density plots, comparing variants that were
688 filtered/removed to those that were not filtered. Note that variants were consistently filtered across the
689 full frequency range, with increased density at frequencies <1% or >10% in ADSP WES.