

1 Transcriptomic Stratification of Late-Onset Alzheimer's Cases Reveals Novel
2 Genetic Modifiers of Disease Pathology

3
4 Nikhil Milind^{1,2*}, Christoph Preuss^{1*}, Annat Haber¹, Guru Ananda¹, Shubhabrata
5 Mukherjee³, Cai John¹, Sarah Shapley^{1,4}, Anna L. Tyler¹, Benjamin A. Logsdon⁵, Paul K.
6 Crane³, Gregory W. Carter^{1†}

7
8 *These authors contributed equally to this work.

9
10 ¹The Jackson Laboratory, Bar Harbor, ME, USA 04609

11 ²Program in Genetics, Department of Biological Science, North Carolina State University,
12 Raleigh, NC, USA 27695

13 ³Department of Medicine, School of Medicine, University of Washington, Seattle, WA, USA
14 98195

15 ⁴Program in Neuroscience, Department of Biology & Geology, Baldwin Wallace University,
16 Berea, OH, USA 44017

17 ⁵Sage Bionetworks, Seattle, WA, USA 98121

18

19

20

21

22

23

24 Keywords: Alzheimer' disease, gene expression, co-expression modules, reproducibility

25

26 †Correspondence:

27 Gregory Carter, Ph.D.

28 The Jackson Laboratory

29 600 Main Street

30 04609 Bar Harbor, Maine

31 Phone: (207)288-6025

32 Email: Gregory.Carter@jax.org

33

34 **ABSTRACT**

35

36 Late-Onset Alzheimer's disease (LOAD) is a common, complex genetic disorder well-
37 known for its heterogeneous pathology. The genetic heterogeneity underlying common
38 complex diseases poses a major challenge for targeted therapies and the identification of
39 novel disease-associated variants. Case-control approaches are often limited to examining
40 a specific outcome in a group of heterogenous patients with different clinical characteristics.
41 Here, we developed a novel approach to define relevant transcriptomic endophenotypes
42 and stratify decedents based on molecular profiles in three independent human LOAD
43 cohorts. By integrating post-mortem brain gene co-expression data from 2114 human
44 samples with LOAD, we developed a novel quantitative, composite phenotype that can
45 better account for the heterogeneity in genetic architecture underlying the disease. We
46 used iterative weighted gene co-expression network analysis (WGCNA) analysis to reduce
47 data dimensionality and to isolate gene sets that are highly co-expressed within disease
48 subtypes and represent specific molecular pathways. We then performed single variant
49 association testing using whole genome-sequencing data for the novel composite
50 phenotype in order to identify genetic loci that contribute to disease heterogeneity. Distinct
51 LOAD subtypes were identified for all three study cohorts (two in ROSMAP, three in Mayo
52 Clinic, two in Mount Sinai Brain Bank). Single variant association analysis identified a
53 genome-wide significant variant in *TMEM106B* (p-value $< 5 \times 10^{-8}$, rs1990620^G) in the
54 ROSMAP cohort that confers protection from the inflammatory LOAD subtype. Taken
55 together, our novel approach can be used to stratify LOAD into distinct molecular subtypes
56 based on affected disease pathways.

57

58 INTRODUCTION

59

60 Late-onset Alzheimer's disease (LOAD) is the most common form of dementia in the
61 elderly. The clinical features associated with LOAD are an amnesic type of memory
62 impairment, deterioration of language, and visuospatial deficits. In the later stages of the
63 disease, symptoms may include motor and sensory abnormalities, gait disturbances, and
64 seizures. Without advances in therapy, the number of symptomatic cases in the United
65 States is predicted to rise to 13.2 million by 2050¹.

66

67 Many common, complex diseases such as LOAD present with heterogeneous phenotypes
68 due to interactions between genetic and environmental factors affecting a range of
69 pathways and processes. LOAD has no simple form of inheritance and is governed by a
70 common set of risk alleles across multiple genes that, in combination, have a substantial
71 effect on disease predisposition and age of onset². Genome-Wide Association Studies
72 (GWAS) have become an important tool for identifying variants in complex diseases^{3,4}.
73 GWAS for LOAD have identified variants in over 500 genes as potential risk factors with the
74 $\epsilon 4$ variant in *APOE* as the strongest contributor to overall disease risk^{2,5}. LOAD has a
75 strong polygenic component and an estimated heritability of up to 80%⁶. It has been
76 challenging to transition from the identification of associated genetic variants to the
77 molecular mechanisms that lead to the accumulation of amyloid plaques and helical tau
78 filaments⁷. Furthermore, there is mounting evidence that the observed heterogeneity in
79 LOAD is associated with multiple distinct subtypes^{8,9}.

80

81 Gene co-expression modules tend to consist of genes that belong to the same cellular
82 pathways or programs and help explain the global properties of the transcriptome as it

83 relates to disease risk¹⁰. Networks-based co-expression module approaches have been
84 used to identify causal variants in Late-Onset Alzheimer's disease^{7,11}. However, such
85 studies have failed to account for the heterogeneity of mechanisms that lead to complex
86 diseases. Here, we analyze whole genome sequencing (WGS) and whole transcriptome
87 data from three independent human cohorts from the Accelerating Medicines Partnership -
88 Alzheimer's Disease (AMP-AD) Consortium. We use gene co-expression modules to
89 develop quantitative phenotypes that account for the complex genetic architecture and
90 heterogeneity of LOAD to more effectively map associated variants using genome-wide
91 association. Furthermore, the method presented in this paper can be used to identify
92 variants in other complex diseases.

93

94 **METHODS**

95

96 **Whole genome sequencing and RNA sequencing data**

97 We obtained whole-genome sequencing and RNA sequencing (RNA-Seq) data from
98 Synapse (<https://www.synapse.org/>) for three cohorts from the AMP-AD consortium, from
99 the Mayo Clinic, Mount Sinai Brain Bank, and Rush University. The Mayo Clinic (Mayo)
100 cohort consists of 276 temporal cortex (TCX) samples from 312 North American Caucasian
101 subjects consisting of cases characterized with LOAD, pathological aging (PA), progressive
102 supranuclear palsy (PSP), or elderly controls¹² (Synapse:syn5550404). The Mount Sinai
103 Brain Bank (MSBB) cohort consists of 214 frontopolar prefrontal cortex (FP), 187 inferior
104 temporal gyrus (IFG), 160 parahippocampal gyrus (PHG), and 187 superior temporal gyrus
105 (STG) samples characterized with LOAD, elderly control, or mild cognitive impairment
106 (MCI) (Synapse: syn3159438). The Rush University's Religious Orders Study and Memory
107 and Aging Project (ROSMAP) cohort consists of 623 dorsolateral prefrontal cortex (DLPFC)

108 samples of individuals from 40 groups of religious orders from across the United States
109 (ROS) and older adults in retirement communities in the Chicago area (MAP),
110 characterized with LOAD, elderly control, or MCI^{7,13} (Synapse:syn3219045). A summary of
111 samples from each of the cohorts is provided in Table S1 and Table S2. Sex, age of death,
112 and batch were used as covariates for normalization in the ROSMAP and Mayo data. Sex,
113 age of death, race, and batch were used as covariates for normalization in the MSBB data.
114 Details on post-mortem brain sample collection, tissue and RNA preparation, sequencing,
115 and sample quality control can be found in published work related to each cohort^{12,14,15}.

116

117 **Co-expression modules and iterativeWGCNA**

118 Data on human AMP-AD co-expression modules were obtained from Synapse (Synapse:
119 syn11932957.1). The modules derive from the three independent LOAD cohorts used in
120 this study. A detailed description on how co-expression modules were identified can be
121 found in a recent study that identified the human co-expression modules as part of a
122 transcriptome wide LOAD meta-analysis¹⁶. In brief, a modified procedure using five
123 different co-expression analysis protocols followed by merging by graph clustering methods
124 was performed to obtain 30 modules across all three cohorts (Synapse: syn2580853), 26 of
125 which corresponded to the six tissue regions used in this study. A summary of these
126 modules is provided in Table S3. We focused on tissues from the frontal cortex, temporal
127 cortex, and hippocampus due to their relevance to LOAD neuropathology¹⁷. These modules
128 are generally large, containing thousands of genes that represent multiple functions¹⁶. In
129 order to construct more functionally-specific submodules from these AMP-AD co-
130 expression modules, we subjected them to a repeated pruning process called
131 iterativeWGCNA¹⁸. Briefly, iterativeWGCNA performed WGCNA on each AMP-AD co-
132 expression module independently. The gene sets produced by this process were then

133 pruned to ensure that only highly correlated genes remained by evaluating the connectivity
134 of the genes to the gene set eigengene. The resulting gene sets, containing highly
135 correlated genes, were combined and the process was repeated until the gene sets
136 converged. The algorithm then attempted to reclassify genes from the residual gene set.
137 We specified a soft-threshold power of six, a minimum eigengene connectivity of 0.6, and a
138 required module size of 100 to promote the generation of submodules that capture
139 pathway-level signals. The final set of 68 submodules consisted of highly correlated and
140 cell-type specific genes. The submodules were mutually exclusive for a given cohort but
141 overlapped with submodules from other cohorts. A summary of these submodules is
142 provided in Table S4. An eigengene for a given submodule is defined as the first principle
143 component of gene expression data within each submodule.

144

145 **Stratification of LOAD cases based on clustering of human co-expression** 146 **submodules**

147 Eigengene expression data for TCX, PHG, FP, and DLPFC regions was used to stratify
148 LOAD cases in separate analyses. Clustering was performed on submodule eigengenes to
149 determine subtypes of LOAD cases in each brain region. The NbClust R package
150 determined the optimal number of clusters for different clustering methods by polling with
151 the majority rule across 30 indices¹⁹. We tested agglomerative hierarchical approaches
152 (Ward, UPGMA, WPGMA) and a reallocation approach (K-means) on the eigengene
153 expression data and evaluated the within-cluster similarity of cases using silhouettes. The
154 silhouette for a given object is a measure that simultaneously assesses how similar the
155 object is to its cluster and how different the object is from all the other clusters²⁰. Prior
156 analysis of simulated genome-wide methylation data suggests that no one clustering
157 method outperforms the other consistently and that mean silhouette widths can be used to

158 pick the ideal clustering method²¹. The silhouette plots revealed that different methods were
159 required for the different regions to generate clusters with the largest average silhouette
160 widths. We determined that K-means was an optimal approach for DLPFC, Ward was
161 optimal for PHG and TCX, and UPGMA was optimal for FP after analyzing silhouette plots
162 of clusters generated by each method for each region. An example of silhouettes used to
163 determine the ideal clustering method for the DLPFC region is shown in Figure S1. A
164 summary of the clusters for each brain region, considered case subtypes, is provided in
165 Table S5. In the subtypes generated for the DLPFC region from the ROSMAP cohort, we
166 assessed each subtype for enrichment of cognitive and pathological measures. We used
167 Braak stages as a measure of neurofibrillary tangle burden and CERAD scores as a
168 measure of neuritic plaque burden^{22,23}. We also assessed the rate of decline in memory,
169 executive function, visuospatial function, and language across the subtypes. Definitions,
170 collection, and standardization of these decline measures can be found in previously
171 published work²⁴.

172

173 **Differential expression analysis of case subtypes**

174 For differential expression analysis, control decedents were defined as cognitively-normal
175 and MCI decedents for PHG, FP, and DLPFC. In the case of TCX, control decedents were
176 defined as cognitively normal, PSP, and PA decedents. For each of the regions used to
177 stratify LOAD cases (TCX, PHG, FP, and DLPFC), we performed differential expression
178 analysis to compare gene expression in case subtypes with control decedents as described
179 above. We used the limma R package to perform the differential expression analysis
180 between subtype and control decedents²⁵. We used the clusterProfiler R package to
181 perform KEGG and Reactome pathway analysis on differentially expressed genes to
182 determine the signal captured by clustering on eigengene expression data²⁶.

183

184 **Single-variant association of eigengene expression and subtype specificity**

185 We used EMMAX, a variance component linear mixed model, to perform single-variant
186 association of our newly derived quantitative traits²⁷. Each submodule eigengene was used
187 as a quantitative trait in single-variant association for its respective brain region. For each
188 region, we also developed a subtype specificity metric by calculating the Euclidean
189 distance between the eigengene expression profile of each decedent and the centroid of
190 each subtype cluster. This resulted in a vector of scores for each subtype that was mapped
191 separately. All quantitative trait mapping results had a genomic inflation factor near one,
192 indicating that there was no significant population substructure effect on the mapping. QQ
193 plot analysis on the p-values showed no evidence of population substructure or
194 confounding effects (Figure S2).

195

196 **Replication of suggestive and significant SNPs in other cohorts**

197 The ROSMAP cohort represented the most adequately powered cohort in the study and
198 was used as a baseline for assessing replication of suggestive and significant SNPs in the
199 other cohorts. SNPs were considered suggestive if quantitative trait mapping with either the
200 submodule eigengenes or the subtype specificity metric resulted in a p-value smaller than
201 1×10^{-5} and genome-wide significant if they resulted in a p-value smaller than 5×10^{-8} , which
202 are standard cutoffs for GWAS. Suggestive and significant SNPs from the DLPFC region in
203 ROSMAP were considered replicated in the TCX, FP, and PHG regions if the SNPs were
204 associated with the submodule eigengenes or subtype specificity metric of the given region
205 at a p-value of 0.05. Summary statistics of prior association studies were obtained from the
206 NHGRI-EBI catalog²⁸. Loci were considered replicated if suggestive and significant SNPs
207 from the ROSMAP cohort were reported in these studies at a p-value smaller than 5×10^{-8} .

208 A summary of the entire analysis is provided in Figure S3.

209

210 **RESULTS**

211

212 **Refinement of 26 human co-expression modules identifies disease-associated** 213 **transcriptomic signals**

214 We performed an iterative gene list pruning process using the iterativeWGCNA approach to
215 refine the 26 human co-expression modules from the AMP-AD consortium. This resulted in
216 subsets, or submodules, of highly correlated genes that were exclusive to each module.

217 Genes that were not highly correlated to any submodule were removed since they are less
218 likely to contribute to the overall signal of the submodule and more likely to introduce noise.

219 We compared the submodules and detected specific LOAD-associated molecular pathways
220 and processes that are shared across the three post-mortem brain cohorts and six brain
221 regions (Figure S4). Furthermore, incorporating information from previously defined cell-
222 type specific markers derived from bulk RNA-Seq and single cell RNA-Seq²⁹ showed that
223 pruning the 26 co-expression modules into 68 submodules resulted in multiple novel cell-
224 type specific submodules (Figure 2, Figure S5). Taken together, these novel 68
225 submodules reflect 15 specific functional consensus clusters that are associated with
226 distinct pathways and processes related to LOAD (Figure S4).

227

228 **Submodule gene sets capture biological signals specific to LOAD pathology**

229 We annotated submodules using GO term enrichment, KEGG pathway enrichment, and
230 Reactome pathway enrichment to highlight the biological specificity of co-expression
231 signals captured by the different submodules (Table S6, Table S7, Table S8). While the 26
232 harmonized co-expression modules were associated with five distinct consensus clusters

233 that captured a broader signal, the submodule associations were more specific in terms of
234 functional enrichment (Figure S4). The 15 functional consensus clusters associated with
235 the 68 submodules revealed cell-type specific signatures and elucidated gene sets for
236 specific biological pathways, including tau-protein kinase activity, neuroinflammation,
237 myelination, and cytoskeletal reorganization (Figure S4).

238

239 **Single-variant association mapping of submodule eigengenes**

240 To map the genetic drivers of biological disease-associated signals resolved by
241 submodules, we performed single-variant association mapping of submodule eigengenes.
242 Eigengenes were defined as the first principle component of the gene expression data
243 associated with each submodule. They capture the variation of gene co-expression and
244 reduce noise associated with the transcriptomic data. Genome-wide suggestive and
245 significant loci were detected for submodules in all four brain regions (Table S9, Table S10,
246 Table S11, Table S12). We identified multiple loci that were replicated across the cohorts at
247 a genome-wide significant level. For instance, rs1990620 is a known variant in *TMEM106B*
248 that was identified as genome-wide significant in the DLPFC region from the ROSMAP
249 cohort was replicated ($p < 5 \times 10^{-2}$) in all other brain regions from the Mayo and MSSM
250 cohorts.

251

252 **Stratification of LOAD cases based on 68 AMP-AD co-expression submodules**

253 Clustering LOAD cases in subtypes based on eigengenes provided a method of assessing
254 genetic drivers of heterogeneity in the transcriptome of LOAD cases. The NbClust package
255 chose between two and three clusters for each region and the number of cases in each
256 cluster was balanced (Table S5). The subtypes were not enriched for common LOAD-
257 associated covariates, such as sex, *APOEε4* genotype, or years of education (Figure 4).

258 Eigengene expression profiles for each subtype were used to assess the association of
259 each subtype with molecular and biological pathways associated with submodules. An
260 example for the ROSMAP cohort is shown in Figure 4. We observed no significant
261 enrichment of cognitive or neuropathological measures between the subtypes for the
262 DLPFC region (Figure S6).

263

264 **ROSMAP subtypes differ in inflammatory response**

265 In order to better understand the underlying molecular differences across the novel LOAD
266 associated subtypes in the ROSMAP cohort and to identify potential subtype specific
267 candidate markers, differential expression analysis was performed for each of the
268 previously defined subtypes against a set of controls (Figure 5a). Each of the two subtypes
269 was compared to a set of 471 decedents from the ROSMAP cohort that were either
270 cognitively normal or had mild cognitive impairment. The Venn diagram in Figure 5b depicts
271 the comparison across the different subtypes. Interestingly, cases associated with Subtype
272 A showed a stronger transcriptional response with 127 differentially expressed genes
273 (adjusted p-values < 0.05, absolute log fold change > 0.5) when compared with controls. Of
274 these genes, 86 were up-regulated and 41 were down-regulated. Among the most
275 significantly down-regulated genes associated with Subtype A cases was the stress-
276 response mediator corticotropin-releasing hormone (*CRH*). Overacting *CRH* signaling has
277 been implicated in inflammatory disorders and LOAD where it has been proposed as a
278 therapeutic target to reduce the negative effects of chronic stress related to memory
279 function and amyloid beta ($A\beta$) production³⁰. Cases associated with Subtype B had 40
280 differentially expressed genes (adjusted p-values < 0.05, absolute log fold change > 0.5),
281 39 of which were down-regulated when compared to controls. Notably, two key pro-
282 inflammatory mediators of amyloid deposition (*S100A8*, *S100A9*) were among the most

283 significantly down-regulated genes in Subtype B decedents when compared to controls
284 (Figure 5a). Both genes, which are established inflammatory biomarkers, are part of a
285 complex that serves as a critical link between the amyloid cascade and inflammatory
286 events in LOAD³¹. Furthermore, multiple pathways linked to *S100A8/9* activation, including
287 IL-10 signaling and complement activation were enriched across down-regulated genes in
288 Subtype B but not in Subtype A decedents as highlighted in Figure 5c. In addition,
289 molecular pathways linked to microglia activation (Figure S8), the immune response, and
290 the stress response were found among the most significant pathways and gene sets (Table
291 S13, Table S14) that differ across subtypes. Gene set enrichment analysis revealed a
292 subset of genes linked to the KEGG osteoclast differentiation pathway (Figure S8),
293 including known AD risk markers such as *TREM2*, *TYROBP*, and *CCL2* among others
294 which were highly up-regulated in Subtype A cases compared to Subtype B cases. This
295 highlights that both molecularly defined LOAD subtypes differ in their immune response
296 and that known LOAD biomarkers, including *S100A8/A9*³², *TREM2*, and *CCL2* might be
297 used to stratify patients based upon their inflammatory response to the observed disease
298 state. These results were consistent with the functional annotations of the previously
299 defined submodules that define both subtypes (Figure 4C).

300

301 **Single variant association mapping for ROSMAP decedents**

302 Genome wide association mapping revealed a differential enrichment of significant variants
303 across subtypes (Figure 6, Table S9, Table S10, Table S11, Table S12). Loci were
304 associated with one or more submodule eigengenes, as shown in Figure 6. One genome-
305 wide suggestive allele in *TMEM106B* was identified for Subtype B (p-value < 4×10^{-6} ,
306 rs1990620^G). This association was replicated at a genome-wide suggestive level in
307 association with the DLPFCbrown_2 eigengene and at a genome-wide significant level with

308 the DLPFCbrown_1 and DLPFCyellow_2 eigengenes (Figure 3). DLPFCbrown_1 contains
309 genes related to myelination and lysosomal activity (KEGG pathways hsa00600 and
310 hsa04142), while DLPFCyellow_2 contains genes related to endocytosis and potassium
311 channel activity (KEGG pathway hsa04144 and Reactome pathway R-HSA-1296071).
312 *TMEM106B* is a known modifier of neurodegenerative disease and cognitive aging, which
313 has been previously linked with cognitive performance³³. Loss of *TMEM106B* function has
314 been shown to rescue lysosomal phenotypes related to frontotemporal dementia³⁴. The
315 identified protective allele rs1990620^G is a known CCCTC-binding factor (CTCF) site, which
316 has been shown to modify the inflammatory response in the course of aging³⁵. Besides the
317 association with *TMEM106B* in Subtype B, protective variants near *MTUS2* were identified
318 which are in close vicinity to *HMGB1*, a locus that has been previously implicated in brain
319 atrophy³⁶. A differential expression analysis of haplotype carriers of the protective
320 rs1990620^G variant in *TMEM106B* showed an up-regulation of neuroactive ligand receptor
321 interactions, while decedents carrying the risk variant showed significant up-regulation for
322 pathways related to Osteoclast differentiation (KEGG pathway hsa04380) and
323 neuroinflammation (data not shown).

324

325 **Suggestive SNPs in ROSMAP are replicated in other cohorts**

326 A total of 1326 unique SNPs representing 163 loci were genome-wide suggestive or
327 significant (p-value < 1×10^{-5}) in the DLPFC region when pooled from all 11 DLPFC
328 eigengenes and two subtype-specific variant mapping analyses. Of these, 645 SNPs were
329 replicated in the PHG analyses, 762 SNPs were replicated in the FP analysis, and 482
330 SNPs were replicated in the TCX analyses (p-value < 1×10^{-2}). The *TMEM106B* variant
331 associated with dementia, rs1990620, was replicated in all cohorts. Of the 163 loci, 29 loci

332 across 27 studies had been previously reported in the NHGRI-EBI catalog such that the
333 most significant SNP from the prior study was a suggestive SNP in the DLPFC region
334 (Table S15, Table S16, Table S17, Table S18).

335

336 **DISCUSSION**

337

338 Common complex diseases such as LOAD are characterized by phenotypic heterogeneity
339 and the presence of multiple common variants affecting disease risk. In this study, we
340 present an analysis that uses transcriptomic co-expression data and whole-genome
341 sequencing from multiple cohorts to dissect phenotypic heterogeneity and identify potential
342 genetic drivers of complex trait pathology in LOAD.

343

344 Here, we used an iterative pruning approach based on 26 human post-mortem co-
345 expression modules to generate 68 novel submodules that contained genes associated
346 with LOAD specific biological pathways and molecular processes. Indeed, we observed
347 that genes in the novel submodules are enriched for functional terms that were specific to
348 pathways associated with LOAD, such as lipid modification, the TREM2/TYROBP pathway,
349 and tau-protein kinase activity. Furthermore, submodules from all six brain regions
350 clustered independently of the co-expression module of origin and brain region, suggesting
351 that the genes captured in each submodule represented signals that were associated with
352 LOAD pathology rather than cohort- or tissue-specific factors. Notably, submodules were
353 much more specific for markers of different brain cell types, suggesting that the processes
354 associated with submodules represent the pathological signals from these specific cell
355 types. This is in line with recent studies showing that different cell types in the brain play
356 specific roles at different stages in the pathogenesis of LOAD³⁷. Taken together, our results

357 demonstrate that the novel human co-expression submodules identified in this study
358 capture cell-type specific pathways associated with LOAD pathogenesis in the brain.

359

360 Mapping the eigengene expression for individual submodules represents a pathway- or
361 process-level alternative to expression quantitative trait locus (eQTL) mapping for each
362 individual transcript. Since the human co-expression submodules represented pathological,
363 cell-type specific pathways in LOAD brain tissue, mapping eigengene expression for
364 decedents was expected to identify genetic drivers of LOAD pathology. RNA-Seq data from
365 post-mortem brain tissue in human cohorts contains a strong immune signal, as evidenced
366 by repeated identification of genetic loci related to microglial response in meta-analyses
367 with increasingly large cohorts^{5,38}. Using submodule eigengenes as quantitative traits for
368 single-variant association provided an opportunity to identify genetic drivers of biological
369 processes that are known to be drivers of early LOAD pathogenesis, such as astrogliosis,
370 neuronal plasticity, myelination, and vascular blood brain barrier interactions³⁷. Suggestive
371 variants identified were unique to subsets of submodules. For instance, the *TMEM106B*
372 locus was associated at a genome-wide significant level with the DLPFCbrown_1 and
373 DLPFCyellow_2 eigengenes (Figure 3), representing processes related to oligodendrocytic
374 myelination, lysosomal activity, endocytosis, and potassium channel activity. The
375 *TMEM106B* locus has been implicated in cognitive aging, with functional consequences in
376 frontotemporal dementia related to lysosomal activity³³⁻³⁵. A submodule of particular
377 interest is the microglia-associated submodule DLPFCblue_3, which contains genes
378 related to the TREM2/TYROBP cascade. The *FAM110A* locus is close to rs1014897 and
379 the *CNTNAP5* locus is close to rs76854344, both variants have been previously associated
380 with posterior cortical atrophy and LOAD³⁹. The *NTM* locus is close to rs1040103, a variant
381 that has been associated with white blood cell count⁴⁰. Thus, quantitative trait mapping of

382 single variants using eigengene expression for submodules presented in this study can
383 elucidate genetic factors specific to associated pathological pathways.

384

385 Furthermore, eigengenes represent a dimensional reduction of transcriptomic data onto
386 axes of pathological relevance. Thus, we expected that clustering on the eigengene
387 expression of LOAD cases would generate pathway-level profiles of putative molecular
388 LOAD subtypes based on case heterogeneity. As anticipated, we observed that average
389 eigengene expression was enriched by subtype for multiple submodules in all four brain
390 regions tested. Strikingly, these enrichments were diametric in the subtypes generated for
391 LOAD cases, an example of which is presented for the DLPFC region in Figure 4. Similar
392 enrichment patterns were identified in the other three brain regions. These results suggest
393 that the biological programs identified by submodules in this study align themselves along
394 the heterogeneity of transcriptomic data present in LOAD cases across multiple cohorts
395 rather than differentiating solely based on cases and controls. Furthermore, the
396 stratification of patients based on submodule expression profiles demonstrated that there is
397 significant variation in immune response in post-mortem brain tissue, a process that is
398 considered a hallmark of LOAD pathogenesis (Figure 5, Figure S8). Variants associated
399 with the subtype specificity metric overlapped with the variants associated with individual
400 submodule eigengenes (Figure 6). This suggests that the genetic factors that influenced
401 subtypes can be dissected into loci driving specific submodules. Furthermore, the
402 deconstruction of genetic loci can provide the basis for more targeted treatment of
403 dysfunctional pathways that contribute to different subtypes of LOAD.

404

405 Our subtypes in the DLPFC brain region of the ROSMAP cohort represent differences in
406 transcriptomic profiles of LOAD cases derived from post-mortem RNA-Seq data. A lack of

407 temporal data makes it challenging to decisively interpret these profiles. The subtypes may
408 represent distinct LOAD endpoints, differences in disease severity, environmental effects,
409 or phases of molecular pathology. Neither subtype was associated with cognitive or
410 neuropathological outcome (Figure S6). Furthermore, covariates such as sex, *APOE*
411 genotype, and years of education were not significantly enriched in any given subtype
412 (Figure 4). This suggests that the transcriptomic profiles do not represent transitions in
413 disease severity and that there are overall risk factors not reflected in transcriptomic
414 subtypes. Furthermore, both subtypes are associated with unique loci that belong to the
415 same community of loci detected by submodule mapping (Figure 6), indicating that the
416 subtypes capture various combinations of genetic elements that lead to LOAD pathology.
417 While suggestive, these transcriptomic LOAD subtypes will require further validation in
418 cohorts that adequately control for disease progression.

419

420 The methodology presented in this study is not limited to RNA-Seq data and can be
421 performed on other omics, such as proteomics or metabolomics. As such data become
422 available for the decedents in these cohorts, this analysis can be expanded across these
423 additional informative dimensions.

424

425 **ACKNOWLEDGEMENTS**

426

427 This study was supported by the National Institutes of Health grant U54 AG 054354.
428 Additional funding was provided by the Barbara H. Sanford Endowed Scholarship Fund, the
429 Robert E. Garrity Cooperative Education Fund as well as the Park and Goldwater
430 Scholarship funds. We thank A. Saykin and K. Nho for helpful conversations.

431

432 The results published here are in whole or in part based on data obtained from the AMP-
433 AD Knowledge Portal ([doi:10.7303/syn2580853](https://doi.org/10.7303/syn2580853)). Study data were provided by the Rush
434 Alzheimer's Disease Center, Rush University Medical Center, Chicago. Data collection was
435 supported through funding by NIA grants P30AG10161, R01AG15819, R01AG17917,
436 R01AG30146, R01AG36836, U01AG32984, U01AG46152, the Illinois Department of
437 Public Health, and the Translational Genomics Research Institute.

438

439 The results published here are in whole or in part based on data obtained from the AMP-
440 AD Knowledge Portal ([doi:10.7303/syn2580853](https://doi.org/10.7303/syn2580853)). Study data were provided by the following
441 sources: The Mayo Clinic Alzheimer's Disease Genetic Studies, led by Dr. Nilufer Taner
442 and Dr. Steven G. Younkin, Mayo Clinic, Jacksonville, FL using samples from the Mayo
443 Clinic Study of Aging, the Mayo Clinic Alzheimer's Disease Research Center, and the Mayo
444 Clinic Brain Bank. Data collection was supported through funding by NIA grants P50
445 AG016574, R01 AG032990, U01 AG046139, R01 AG018023, U01 AG006576, U01
446 AG006786, R01 AG025711, R01 AG017216, R01 AG003949, NINDS grant R01
447 NS080820, CurePSP Foundation, and support from Mayo Foundation. Study data includes
448 samples collected through the Sun Health Research Institute Brain and Body Donation
449 Program of Sun City, Arizona. The Brain and Body Donation Program is supported by the
450 National Institute of Neurological Disorders and Stroke (U24 NS072026 National Brain and
451 Tissue Resource for Parkinson's Disease and Related Disorders), the National Institute on
452 Aging (P30 AG19610 Arizona Alzheimer's Disease Core Center), the Arizona Department
453 of Health Services (contract 211002, Arizona Alzheimer's Research Center), the Arizona
454 Biomedical Research Commission (contracts 4001, 0011, 05-901 and 1001 to the Arizona
455 Parkinson's Disease Consortium) and the Michael J. Fox Foundation for Parkinson's
456 Research. The results published here are in whole or in part based on data obtained from

457 the AMP-AD Knowledge Portal ([doi:10.7303/syn2580853](https://doi.org/10.7303/syn2580853)). These data were generated
458 from postmortem brain tissue collected through the Mount Sinai VA Medical Center Brain
459 Bank and were provided by Dr. Eric Schadt from Mount Sinai School of Medicine.

460

461 **AUTHOR'S CONTRIBUTIONS**

462

463 NM, CP, AH, GA and CJ performed the genetic and transcriptomic analysis. SS annotated
464 the functional variants. SM and PC provided cognitive and phenotype data for the analysis.
465 BL and AT provided additional transcriptomic data for the analysis. GWC supervised and
466 designed the project. NM, CP and GWC wrote the manuscript. All authors read and
467 approved the final manuscript.

468

469 **CONSENT FOR PUBLICATION**

470

471 All authors have approved of the manuscript and agree with its submission.

472

473 **COMPETING INTERESTS**

474

475 *Not applicable.*

476

477 **REFERENCES**

478

- 479 1. Cummings, J. L. Alzheimer's Disease. *N. Engl. J. Med.* **351**, 56–67 (2004).
- 480 2. Bertram, L. & Tanzi, R. E. Thirty years of Alzheimer's disease genetics: the
481 implications of systematic meta-analyses. *Nat. Rev. Neurosci.* **9**, 768–778 (2008).
- 482 3. Kilpinen, H. & Barrett, J. C. How next-generation sequencing is transforming complex
483 disease genetics. *Trends Genet.* **29**, 23–30 (2013).
- 484 4. Risch, N. & Merikangas, K. The future of genetic studies of complex human diseases.
485 *Science* (1996). doi:10.1126/science.273.5281.1516
- 486 5. Jansen, I. E. *et al.* Genome-wide meta-analysis identifies new loci and functional
487 pathways influencing Alzheimer's disease risk. *Nat. Genet.* **51**, 404–413 (2019).
- 488 6. Verheijen, J. & Sleegers, K. Understanding Alzheimer Disease at the Interface
489 between Genetics and Transcriptomics. *Trends Genet.* **34**, 434–447 (2018).
- 490 7. Mostafavi, S. *et al.* A molecular network of the aging human brain provides insights
491 into the pathology and cognitive decline of Alzheimer's disease. *Nat. Neurosci.* **21**,
492 811–819 (2018).
- 493 8. Mukherjee, S. *et al.* Genetic data and cognitively-defined late-onset Alzheimer's
494 disease subgroups. *Mol. Psychiatry* 1–10 (2018). doi:10.1101/367615
- 495 9. Ferreira, D. *et al.* Distinct subtypes of Alzheimer's disease based on patterns of brain
496 atrophy: longitudinal trajectories and clinical applications. *Sci. Rep.* **7**, 1–13 (2017).
- 497 10. Langfelder, P. & Horvath, S. Eigengene networks for studying the relationships
498 between co-expression modules. *BMC Syst. Biol.* **1**, 54 (2007).
- 499 11. Zhang, B. *et al.* Integrated systems approach identifies genetic nodes and networks
500 in late-onset Alzheimer's disease. *Cell* **153**, 707–20 (2013).
- 501 12. Allen, M. *et al.* Human whole genome genotype and transcriptome data for

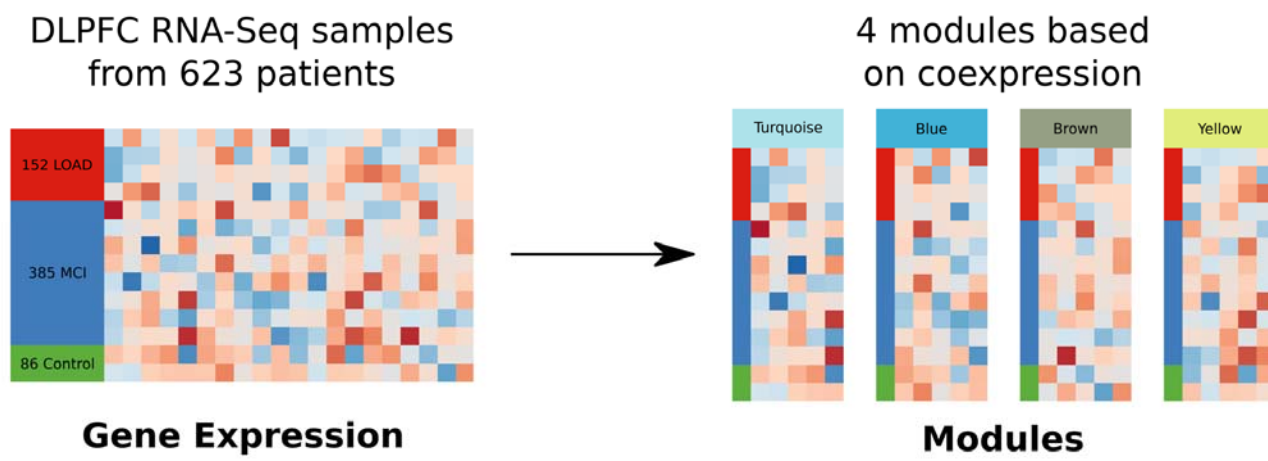
- 502 Alzheimer's and other neurodegenerative diseases. *Sci. Data* **3**, (2016).
- 503 13. Chibnik, L. B. *et al.* Susceptibility to neurofibrillary tangles: role of the PTPRD locus
504 and limited pleiotropy with other neuropathologies. *Mol. Psychiatry* **23**, 1521–1529
505 (2018).
- 506 14. Wang, M. *et al.* The Mount Sinai cohort of large-scale genomic, transcriptomic and
507 proteomic data in Alzheimer's disease. *Sci. Data* **5**, 1–16 (2018).
- 508 15. De Jager, P. L. *et al.* A multi-omic atlas of the human frontal cortex for aging and
509 Alzheimer's disease research. *Sci. Data* **5**, 1–13 (2018).
- 510 16. Logsdon, B. *et al.* Meta-analysis of the human brain transcriptome identifies
511 heterogeneity across human AD coexpression modules robust to sample collection
512 and methodological approach. *bioRxiv* 510420 (2019). doi:10.1101/510420
- 513 17. DeTure, M. A. & Dickson, D. W. The neuropathological diagnosis of Alzheimer
514 disease. *Mol. Neurodegener.* **14**, 1–18 (2019).
- 515 18. Greenfest-Allen, E., Cartailier, J.-P., Magnuson, M. A. & Stoeckert, C. J.
516 iterativeWGCNA: iterative refinement to improve module detection from WGCNA co-
517 expression networks. *bioRxiv* 234062 (2017). doi:10.1101/234062
- 518 19. Charrad, M., Ghazzali, N., Boiteau, V. & Niknafs, A. **NbClust**: An R Package for
519 Determining the Relevant Number of Clusters in a Data Set. *J. Stat. Softw.* **61**, 1–36
520 (2014).
- 521 20. Rousseeuw, P. J. Silhouettes: A graphical aid to the interpretation and validation of
522 cluster analysis. *J. Comput. Appl. Math.* **20**, 53–65 (1987).
- 523 21. Clifford, H., Wessely, F., Pendurthi, S. & Emes, R. D. Comparison of Clustering
524 Methods for Investigation of Genome-Wide Methylation Array Data. *Front. Genet.* **2**,
525 88 (2011).
- 526 22. Braak, H., Thal, D. R., Ghebremedhin, E. & Del Tredici, K. Stages of the Pathologic

- 527 Process in Alzheimer Disease. *J. Neuropathol. Exp. Neurol.* **70**, 960–969 (2011).
- 528 23. Wilson, R. S., Arnold, S. E., Schneider, J. A., Li, Y. & Bennett, D. A. Chronic Distress,
529 Age-Related Neuropathology, and Late-Life Dementia. *Psychosom. Med.* **69**, 47–53
530 (2007).
- 531 24. Mukherjee, S. *et al.* Genetic data and cognitively defined late-onset Alzheimer’s
532 disease subgroups. *Mol. Psychiatry* (2018). doi:10.1038/s41380-018-0298-8
- 533 25. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-
534 sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47–e47 (2015).
- 535 26. Yu, G., Wang, L.-G., Han, Y. & He, Q.-Y. clusterProfiler: an R Package for
536 Comparing Biological Themes Among Gene Clusters. *Omi. A J. Integr. Biol.* **16**, 284–
537 287 (2012).
- 538 27. Kang, H. M. *et al.* Variance component model to account for sample structure in
539 genome-wide association studies. *Nat. Genet.* **42**, 348–354 (2010).
- 540 28. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide
541 association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.*
542 **47**, D1005–D1012 (2019).
- 543 29. McKenzie, A. T. *et al.* Brain Cell Type Specific Gene Expression and Co-expression
544 Network Architectures. *Sci. Rep.* **8**, 8868 (2018).
- 545 30. Futch, H. S., Croft, C. L., Truong, V. Q., Krause, E. G. & Golde, T. E. Targeting
546 psychologic stress signaling pathways in Alzheimer’s disease. *Mol. Neurodegener.*
547 **12**, 49 (2017).
- 548 31. Vogl, T., Gharibyan, A. L. & Morozova-Roche, L. A. Pro-Inflammatory S100A8 and
549 S100A9 Proteins: Self-Assembly into Multifunctional Native and Amyloid Complexes.
550 *Int. J. Mol. Sci.* **13**, 2893 (2012).
- 551 32. Horvath, I. *et al.* Pro-inflammatory S100A9 Protein as a Robust Biomarker

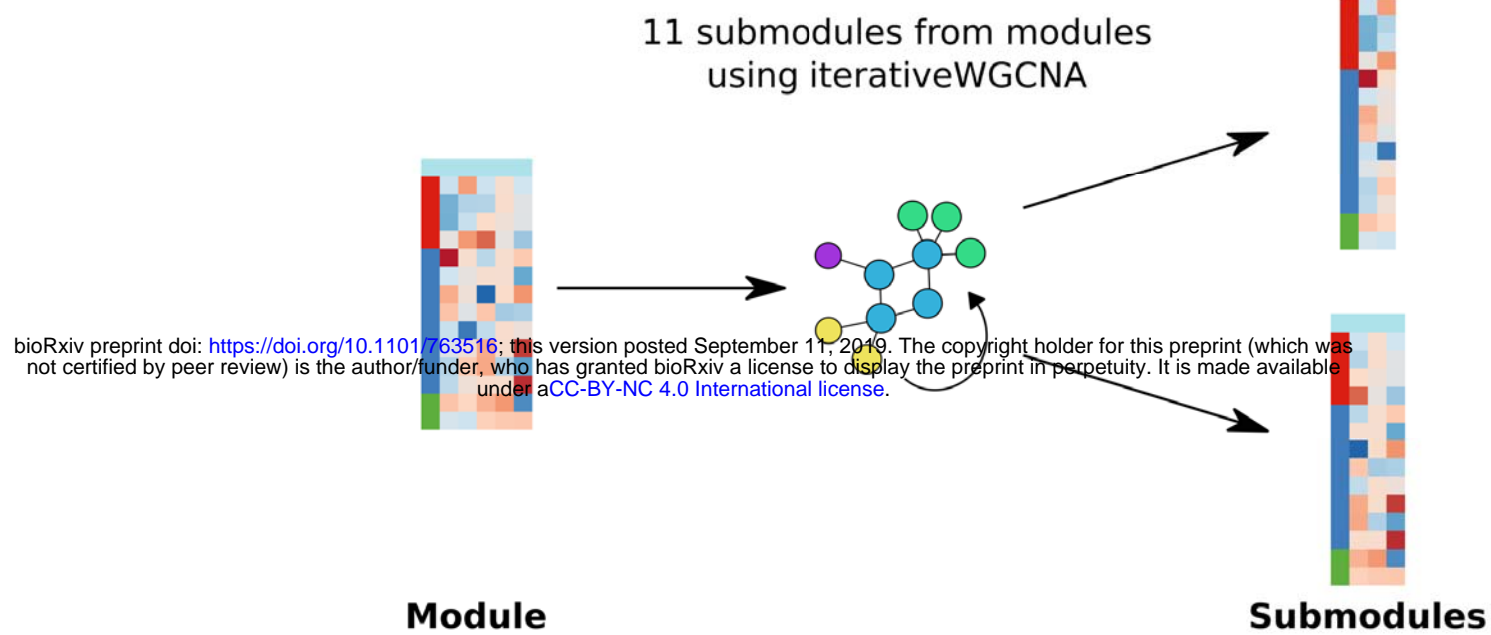
- 552 Differentiating Early Stages of Cognitive Impairment in Alzheimer's Disease. *ACS*
553 *Chem. Neurosci.* **7**, 34–39 (2016).
- 554 33. White, C. C. *et al.* Identification of genes associated with dissociation of cognitive
555 performance and neuropathological burden: Multistep analysis of genetic, epigenetic,
556 and transcriptional data. *PLOS Med.* **14**, e1002287 (2017).
- 557 34. Klein, Z. A. *et al.* Loss of TMEM106B Ameliorates Lysosomal and Frontotemporal
558 Dementia-Related Phenotypes in Progranulin-Deficient Mice. *Neuron* **95**, 281–296.e6
559 (2017).
- 560 35. Gallagher, M. D. *et al.* A Dementia-Associated Risk Variant near TMEM106B Alters
561 Chromatin Architecture and Gene Expression. *Am. J. Hum. Genet.* **101**, 643–663
562 (2017).
- 563 36. Furney, S. J. *et al.* Genome-wide association with MRI atrophy measures as a
564 quantitative trait locus for Alzheimer's disease. *Mol. Psychiatry* **16**, 1130–8 (2011).
- 565 37. De Strooper, B. & Karran, E. The Cellular Phase of Alzheimer's Disease. *Cell* **164**,
566 603–615 (2016).
- 567 38. Lambert, J. C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new
568 susceptibility loci for Alzheimer's disease. *Nat. Genet.* **45**, 1452–1458 (2013).
- 569 39. Schott, J. M. *et al.* Genetic risk factors for the posterior cortical atrophy variant of
570 Alzheimer's disease. *Alzheimer's Dement.* **12**, 862–871 (2016).
- 571 40. Kichaev, G. *et al.* Leveraging Polygenic Functional Enrichment to Improve GWAS
572 Power. *Am. J. Hum. Genet.* **104**, 65–75 (2019).
- 573
- 574

FIGURES

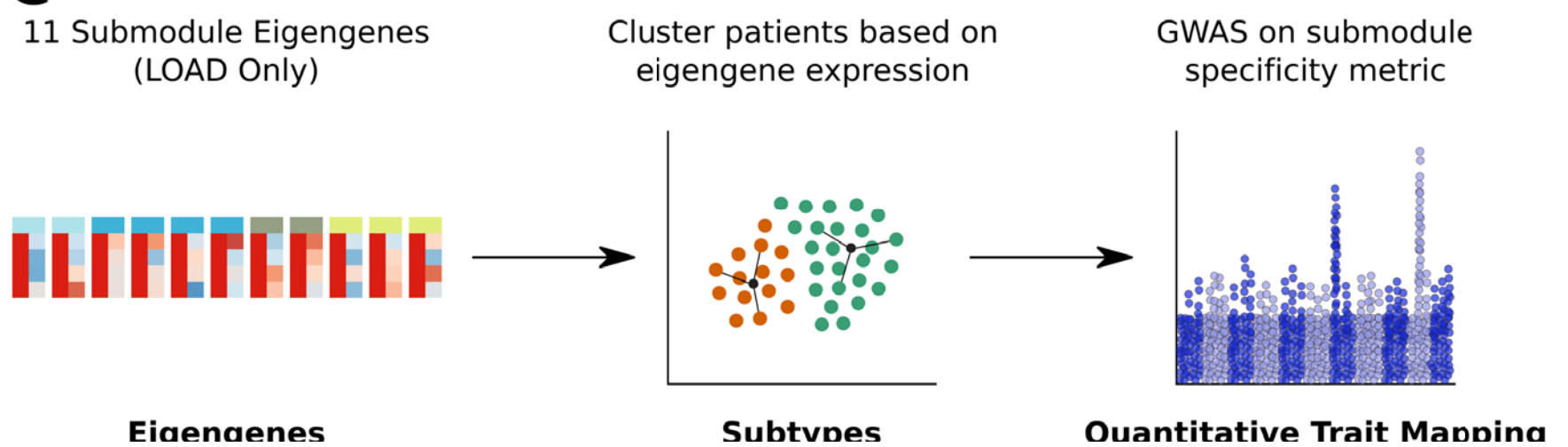
A



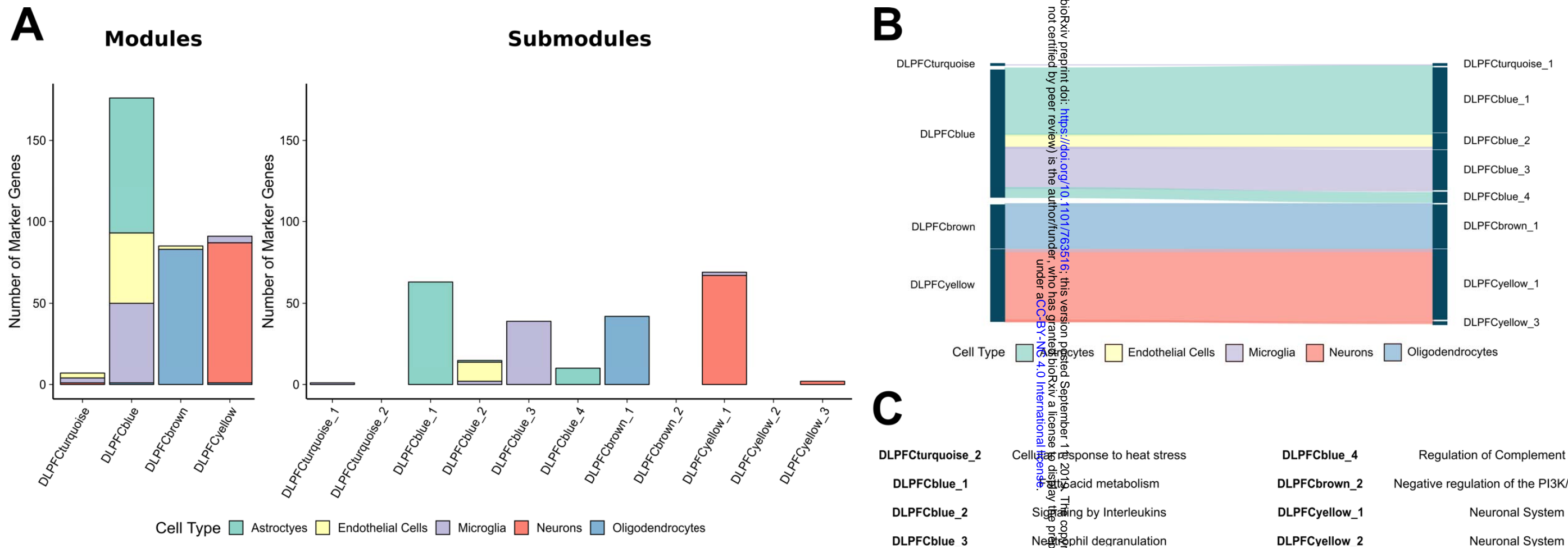
B



C

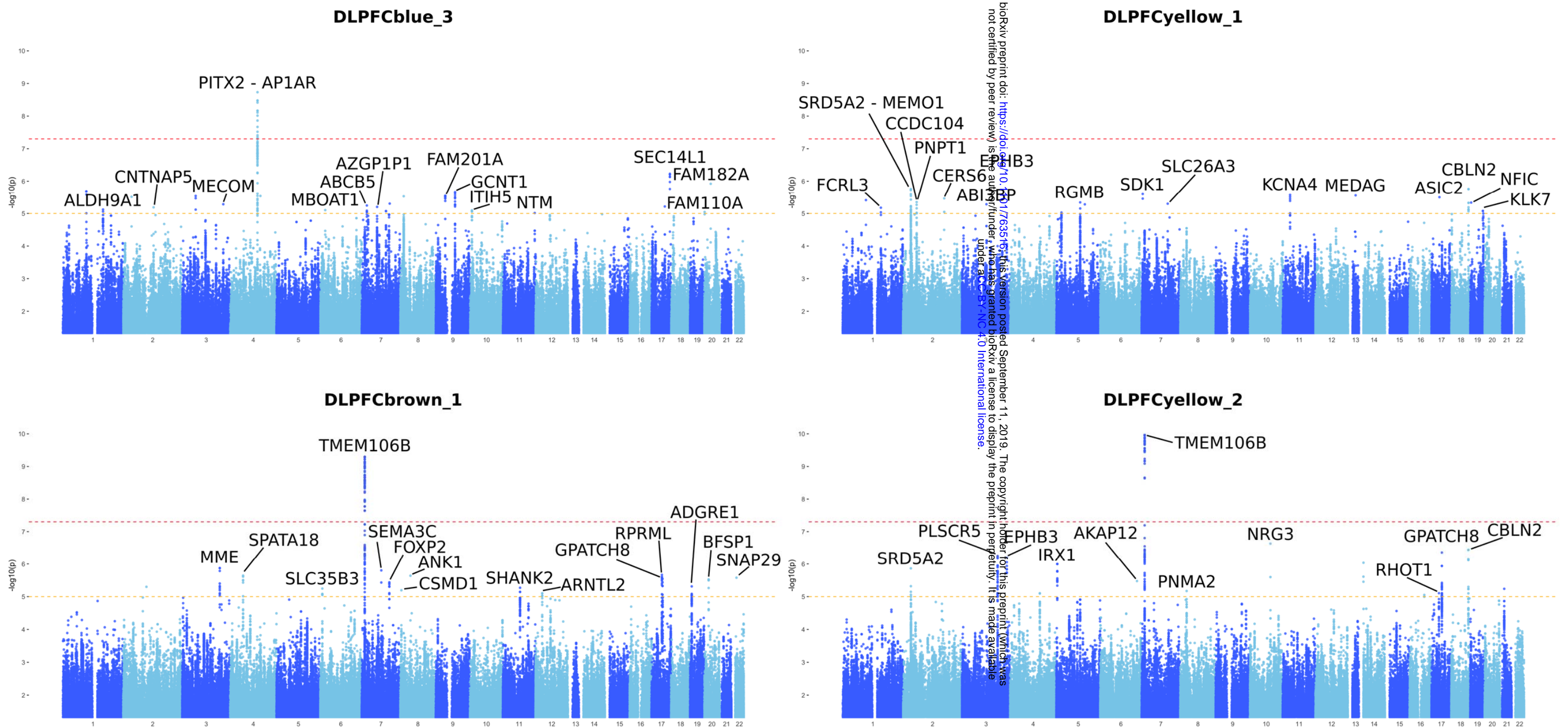


579 **Figure 1: Methodology used in this study to map genetic drivers of LOAD pathology in the ROSMAP cohort.** (A) RNA-
 580 Seq was performed on dorsolateral prefrontal cortex (DLPFC) tissue from 86 control decedents, 385 decedents with Mild
 581 Cognitive Impairment (MCI), and 152 decedents with Late-Onset Alzheimer's Disease (LOAD). A modified procedure using
 582 seven different WGCNA protocols, followed by merging by clustering methods, was performed to obtain 4 modules based on
 583 gene co-expression. (B) Each of the four modules was subjected to iterativeWGCNA, a procedure that repeatedly performs
 584 WGCNA on expression data to generate highly correlated gene sets and exclude weakly correlated genes. 11 submodules
 585 were generated from the 4 modules. (C) The eigengene, or first principal component of each submodule, was calculated for
 586 all 11 submodules and used as a quantitative trait for single-variant association mapping. Furthermore, the eigengene
 587 expression for LOAD cases was used to perform cluster analysis and generate subtypes of LOAD cases. A Euclidean
 588 distance quantitative trait was developed to identify genomic loci for each subtype using single-variant association mapping.
 589

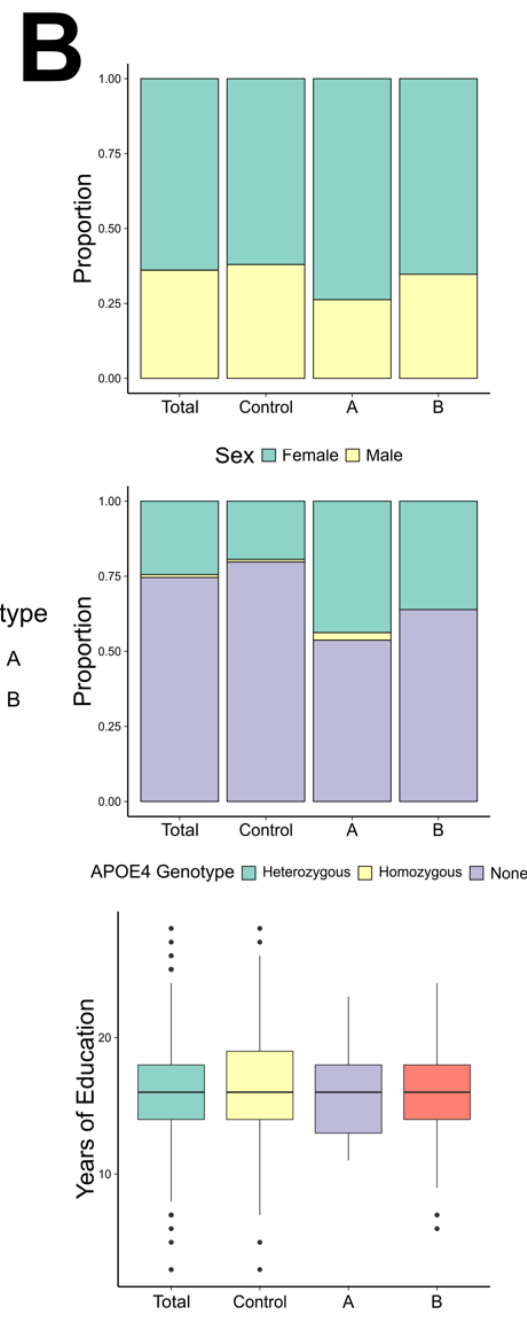
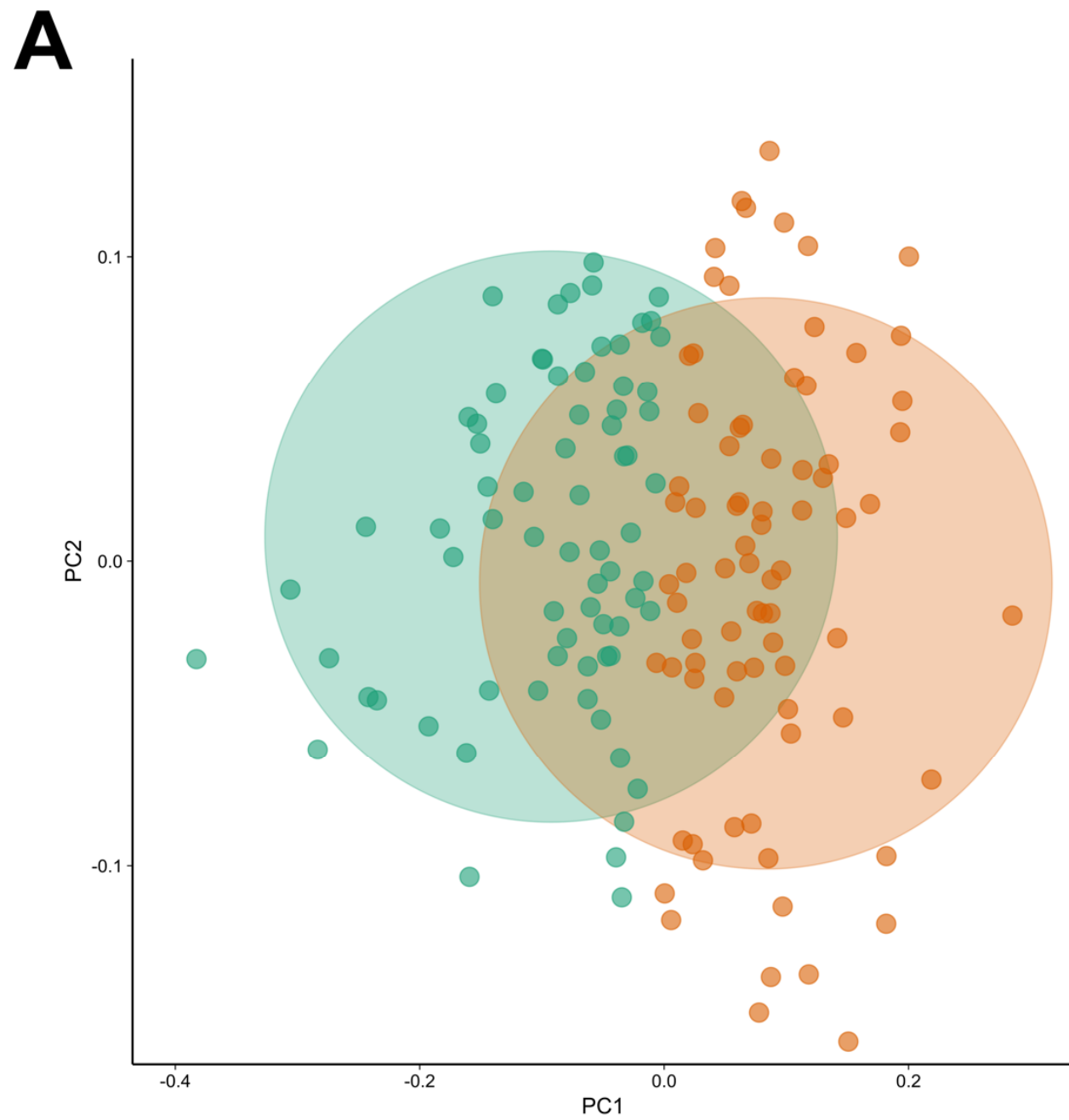


590
591
592
593
594
595
596

Figure 2: Cell-type specificity of modules is refined in submodules. (A) Cell type specific marker genes reported by McKenzie *et al.* were used to annotate modules and submodules for astrocytes, endothelial cells, microglia, neurons, and oligodendrocytes. The top 100 marker genes for each cell type were used. The iterativeWGCNA procedure generated submodules that were more cell-type specific than their modules of origin. (B) A Sankey diagram demonstrates which cell-type specific markers from modules were found in submodules generated using iterativeWGCNA for the ROSMAP cohort. (C) Gene set enrichment analysis for Reactome pathways was performed for each submodule gene list. The top enriched Reactome pathways for submodules are reported.



597
 598
 599 **Figure 3: Manhattan plots of single-variant association of select submodule eigengenes in ROSMAP.** Eigengene expression for each submodule was used as a quantitative trait when performing single-variant
 600 mapping. These Manhattan plots were generated for select DLPFC region submodule eigengenes. Multiple submodule eigengenes were associated with SNPs at a genome-wide significance level of $p = 5e-08$ (red
 601 dotted line). Loci of interest are annotated with the gene closest to the region. Some SNPs were also detected at a genome-wide suggestive level of $p = 1e-05$ (yellow dotted line). DLPFCblue_3 contains genes
 602 related to the TREM2/TYROBP pathway, an important network of genes related to microglial activation during neuroinflammation of the brain. Submodules were associated with both unique and overlapping loci. For
 603 example, DLPFCbrown_1 and DLPFCyellow_2 are derived from separate co-expression modules but were both associated with the *TMEM106B* locus. Similarly, DLPFCyellow_1 and DLPFCyellow_2 were derived
 604 from the same co-expression module but were associated with a mix of overlapping and unique loci.
 605



bioRxiv preprint doi: <https://doi.org/10.1101/763516>; this version posted September 11, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.

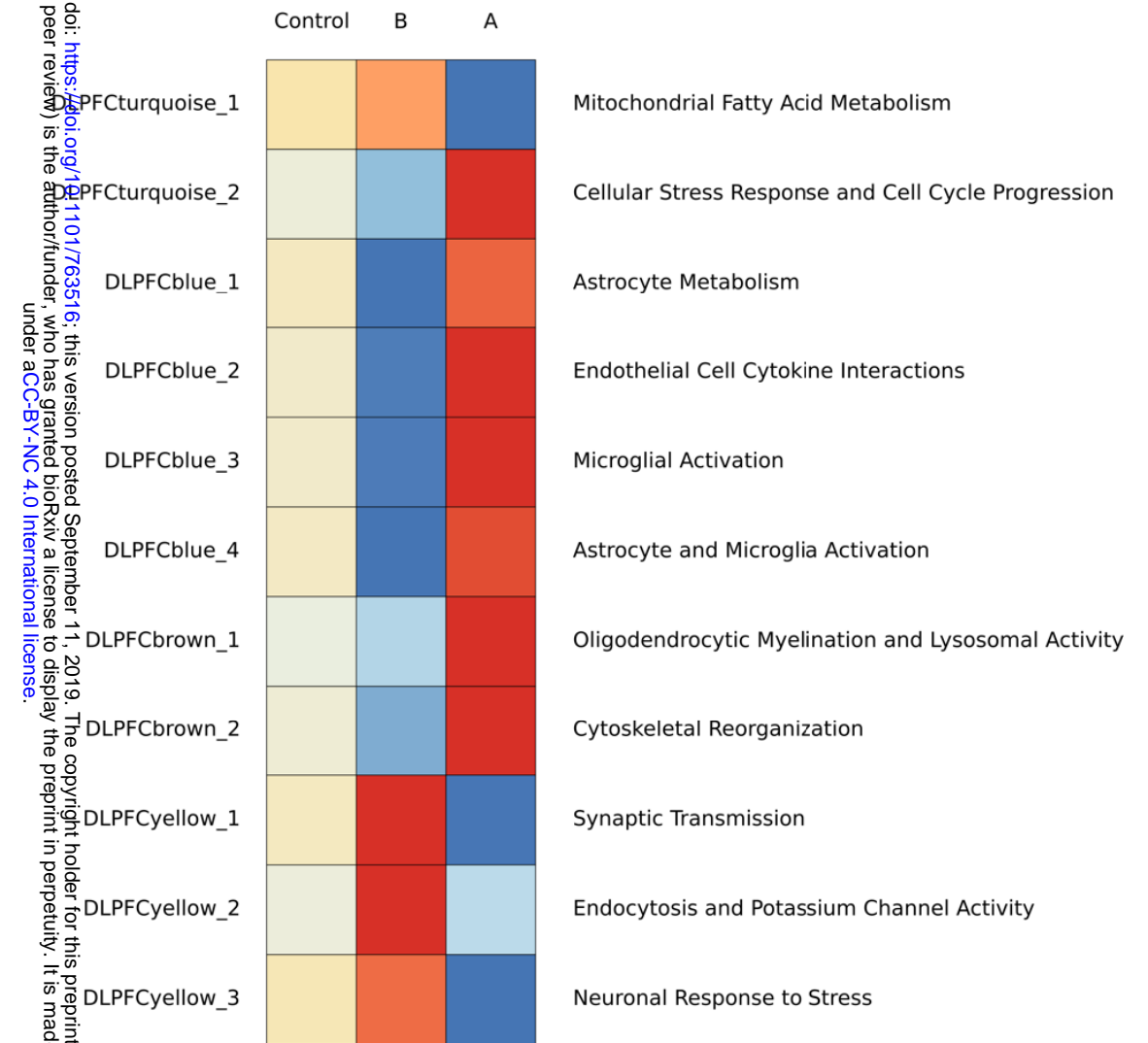
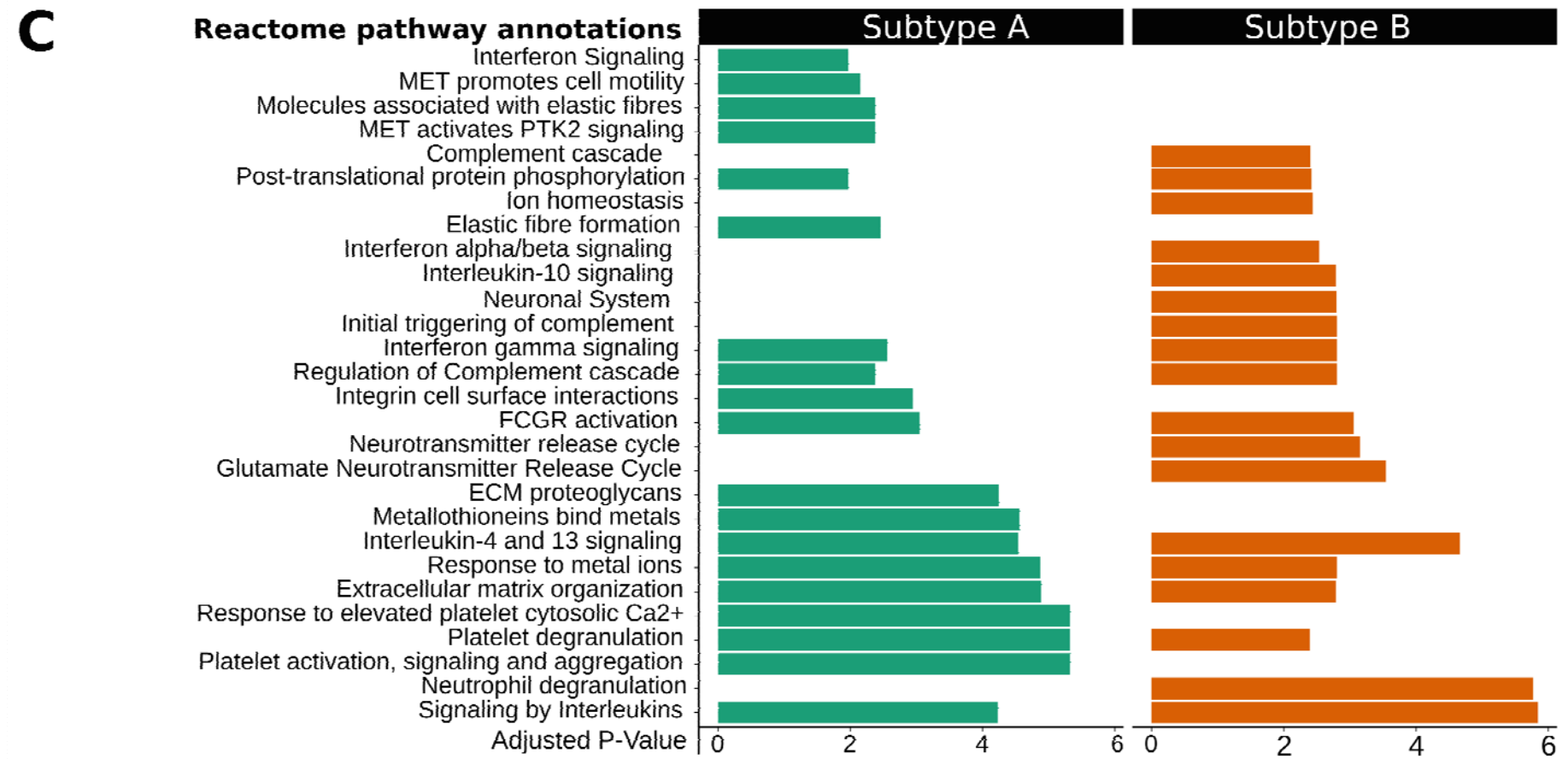
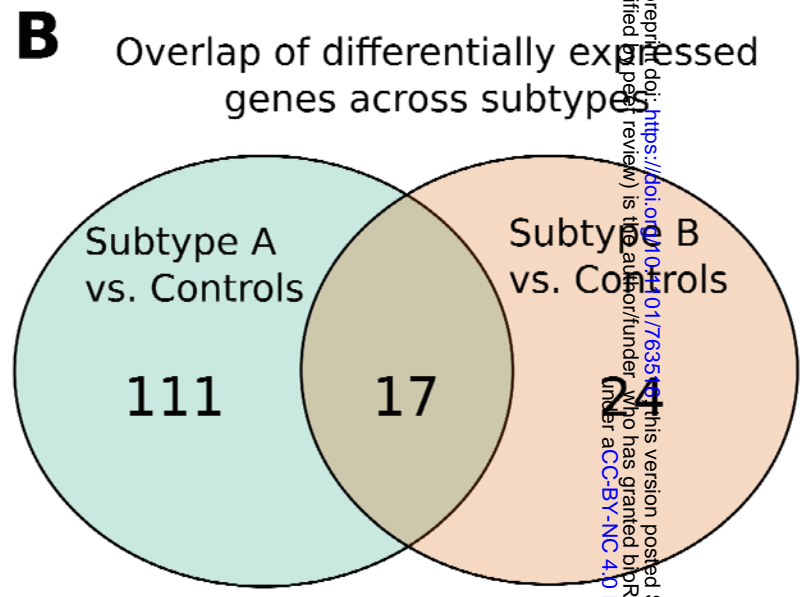
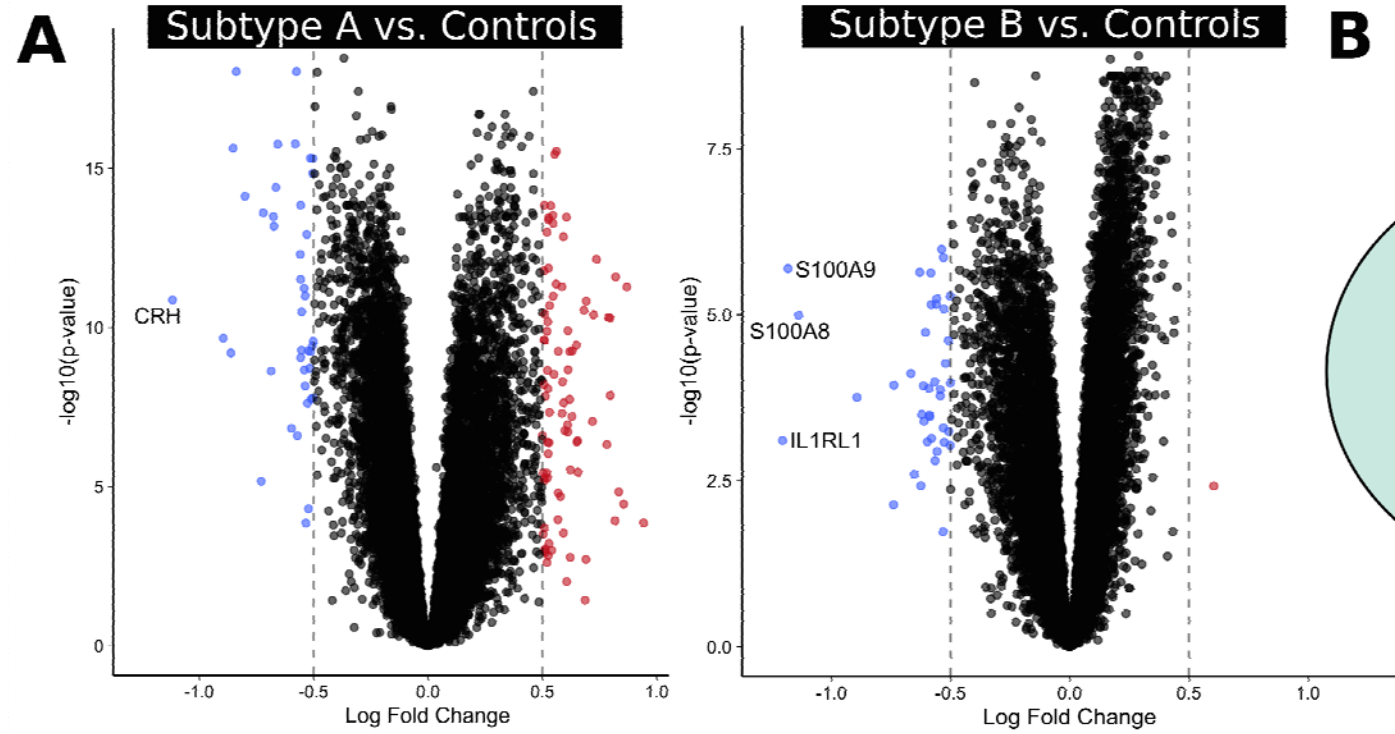


Figure 4: Clustering on eigengene expression in ROSMAP data generates 2 subtypes. (A) Eigengene expression was used to cluster cases into subtypes using K-Means clustering for the DLPFC region. The number of clusters were determined by democratizing results across 30 mathematical indices using the NbClust R package. Two clusters with relatively equal number of cases were generated. (B) There were no significant differences in proportion of sex, *APOEε4* genotype, and years of education between subtypes. (C) The scaled eigengene expression profile of the subtypes demonstrated a strong immune and neuronal signal when compared to control and MCI decedents.

613
614



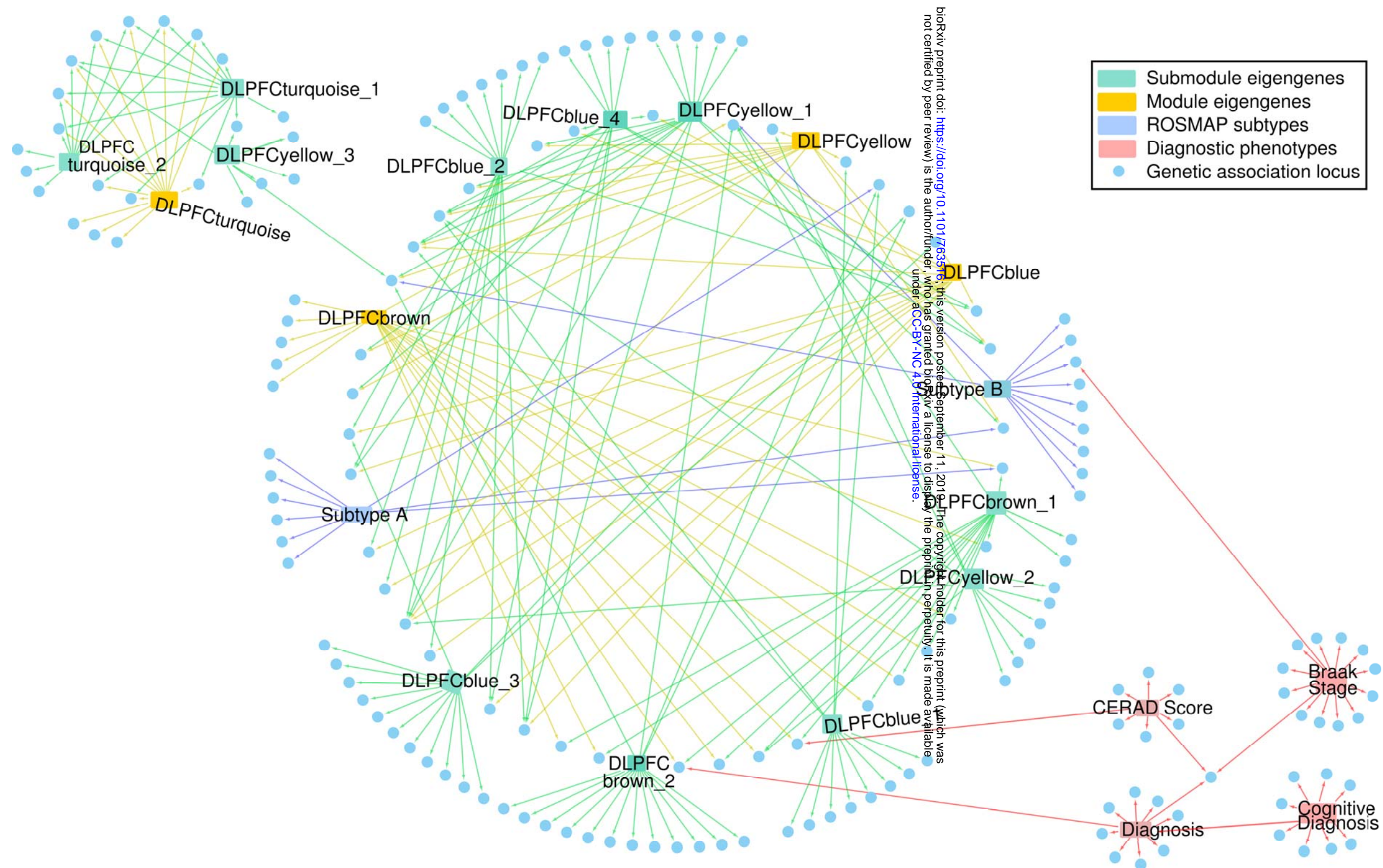
bioRxiv preprint doi: <https://doi.org/10.1101/176351>; this version posted September 11, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.

615
616

617
618
619
620
621
622

Figure 5: Differential expression analysis of ROSMAP subtypes reveals heterogeneity in inflammatory response in LOAD cases. (A) Differential expression analysis comparing each subtype to control decedents for the DLPFC region was performed using the limma R package. We show up-regulated (red, $p < 0.05$, log fold change > 0.5) and down-regulated (blue, $p < 0.05$, log fold change < -0.5) genes in the volcano plot and label genes that have an absolute log fold change of greater than 1 (dotted lines). (B) Differentially expressed genes ($p < 0.05$, absolute log fold change > 0.5) from the analysis show a partial overlap between subtypes. (C) Top Reactome pathways for differentially expressed genes for both subtypes are reported. Subtype A demonstrates an enrichment of immune and stress-response related pathways across up-regulated genes, while Subtype B demonstrates a down-regulation of a set of specific immune-related pathways linked to S100A8/A9 activation.

bioRxiv preprint doi: <https://doi.org/10.1101/763516>; this version posted September 11, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.



623 **Figure 6: Network of phenotypes and associated loci.** We created a directed network describing the loci detected from the multiple analyses in this study. Blue nodes represent loci associated with a phenotype.
 624 Red nodes represent phenotypes. An edge from a phenotype to a genetic locus signifies that the locus is associated with the specified phenotype. Diagnostic phenotypes (red edges) were associated with some of
 625 the loci detected in this study. The module eigengenes (yellow edges), submodule eigengenes (green edges), and subtypes (blue edges) were associated with overlapping and unique loci (center and left). A
 626 community of loci was associated with multiple submodules associated with microglia, endothelial cells, astrocytes, and oligodendrocytes (center). A small community of loci was associated with submodules related
 627 to proteostasis (left). Diagnostic phenotypes included CERAD scores, Braak stages, cognitive diagnosis, and case-control association.
 628