

1 **Molecular estimation of neurodegeneration pseudotime in older brains**

2
3 Sumit Mukherjee^{1,2*}, Laura Heath^{1*}, Christoph Preuss³, Suman Jayadev⁴, Gwenn A. Garden⁴, Anna K
4 Greenwood¹, Solveig K Sieberts¹, Phillip L De Jager^{5,6}, Nilufer Ertekin-Taner^{7,8}, Gregory W Carter³, Lara
5 M Mangravite¹, Benjamin A Logsdon^{1,9}

6
7 ¹Sage Bionetworks, Seattle, WA, USA.

8 ²Microsoft, Redmond, WA, USA[†].

9 ³The Jackson Laboratory, Bar Harbor, ME, USA.

10 ⁴Department of Neurology, University of Washington, Seattle, WA, USA.

11 ⁵Center for Translational & Computational Neuroimmunology, Department of Neurology, Columbia
12 University Irving Medical Center, New York City, USA

13 ⁶Taub Institute, Columbia University Irving Medical Center, New York City, USA

14 ⁷Mayo Clinic Florida, Department of Neurology, Jacksonville, FL, USA.

15 ⁸Mayo Clinic Florida, Department of Neuroscience, Jacksonville, FL, USA.

16

17 [†]Current Address.

18 ^{*}These authors contributed equally to this work.

19

20 ⁹Correspondence:

21 Benjamin A Logsdon, PhD

22 Sage Bionetworks

23 2901 3rd Avenue Suite 330

24 Seattle, WA 98121

25 ben.logsdon@sagebionetworks.org

26

27 **Abstract**

28 Therapeutic treatments for late-onset Alzheimer's disease (LOAD) are hindered by an incomplete
29 understanding of the temporal molecular changes that lead to disease onset and progression. Here, we
30 evaluate the ability of manifold learning to develop a molecular model for the unobserved temporal
31 disease progression from RNA-Seq data collected from human postmortem brain samples collected
32 within the ROS/MAP and Mayo Clinic RNA-Seq studies of the AMP-AD consortium. This approach
33 defines a cross-sectional ordering across samples based on their relative similarity in RNA-Seq profiles
34 and uses this information to define an estimate of molecular disease stage – or disease progression
35 pseudotime - for each sample. This transcriptional estimate of disease progression is strongly concordant
36 with burden of tau pathology (Braak score, $P = 1.0 \times 10^{-5}$), amyloid pathology (CERAD score, $P = 1.8 \times 10^{-5}$),
37 and cognitive diagnosis ($P = 3.5 \times 10^{-7}$) of LOAD. Further, the disease progression estimate
38 recapitulates known changes in cell type abundance and impact of genes that harbor known AD risk loci.
39 Samples estimated to reside early in disease progression were enriched for control and early stage AD
40 cases, and demonstrated changes in basic cellular functions. Samples estimated to reside late in disease
41 progression were enriched for late-stage AD cases, and demonstrated changes in known disease processes
42 including neuroinflammation and amyloid pathology. We also identified a set of control samples with
43 late-stage estimated disease progression who also showed compensatory changes in genes involved in
44 affected pathways are protein trafficking, splicing, regulation of apoptosis, and prevention of amyloid
45 cleavage. In summary, we present a disease specific method for ordering patients based on their LOAD
46 disease progression from brain transcriptomic data.

47 **Introduction**

48 Late-onset Alzheimer's disease (LOAD) is a devastating illness with no effective disease modifying
49 therapy, owing to a 99.6% failure of clinical trials¹. There is a growing consensus that the most effective
50 treatments will intervene early in disease progression and halt disease pathophysiological processes prior
51 to conversion to LOAD². In addition, there is increasing recognition that LOAD may in fact be a

52 spectrum of related diseases that have similar clinical and neuropathological manifestations^{3,4}. Devising
53 successful therapeutic strategies will likely require targeting potentially diverse early-stage disease
54 processes that occur prior to a high burden of neuropathology or cognitive impairment.
55
56 Current approaches to identify AD affected individuals include *in vivo* measures of the pathological
57 hallmarks of disease – amyloid, tau, and neurodegeneration – via CSF biomarkers for amyloid and tau⁵,
58 positron emission tomography for amyloid and tau (PET)⁶, and structural and functional MRI of
59 neurodegeneration. Cognitive assessments are used to estimate disease burden⁷, although measurable
60 cognitive impairment generally indicates a sustained burden of neuropathology and advanced
61 neurodegeneration. Based on biomarker studies of AD, by the time cognitive decline becomes detectable,
62 neuropathological changes of AD have already occurred, first in A β and subsequently in tau related
63 measures⁸ and therefore cannot be used to select patients for early disease stage studies. Furthermore,
64 while these measures of disease progression capture the overall increase in burden of pathology and
65 cognitive decline, they do not necessarily identify the dysfunctional molecular mechanisms that lead to
66 neuropathology and cognitive decline. There are likely many independent patient specific molecular
67 pathways present at an early stage in disease that then contribute to later stage disease progression^{9,10}.
68 This motivates the need to identify these early stage molecular mechanisms driving disease progression.
69
70 The Accelerating Medicines Project for Alzheimer’s Disease (AMP-AD) consortia have generated
71 genome-wide transcriptomics of post-mortem brain tissue from patients across a broad range of
72 Alzheimer’s disease neuropathological progression – including individuals with various stages of AD
73 neuropathology and those who lack AD neuropathology, but who may in fact harbor early stage disease
74 molecular processes. We therefore sought to chart the molecular progression of the disease as reflected in
75 the aggregate behavior of the CNS transcriptome across these individuals. While standard approaches
76 such as differential expression or coexpression analyses have proven informative^{11–15}, these analyses do
77 not infer the relative stage of disease progression or identify distinct disease subtypes. Here we propose

78 an approach to analyze population level RNA-seq data from post-mortem brain tissue to learn a tree
79 structured progression (**Figure 1**) that represents distinct sub-types of disease and the relative progression
80 of disease across patients. With this approach, we identify potentially generalizable trajectories of LOAD
81 across heterogeneous patient populations at all stages of disease. Furthermore, we characterize molecular
82 pathways that define disease stages – a potential source of new biomarkers and therapeutic interventions
83 for early-stage disease processes along multiple different disease trajectories.

84

85 To learn the molecular disease staging and neuropathologic progression tree we use a manifold learning
86 method¹⁶. Manifold learning refers to a group of algorithms that aim to recover the low dimensional
87 subspace underlying a high dimensional dataset. Previous authors use manifold learning to estimate
88 disease progression from neuroimaging data¹⁷ and to study lineage commitment of cells during
89 differentiation from single cell RNA-seq (scRNA-seq)¹⁸⁻²¹. To our knowledge, manifold learning has not
90 been used to estimate disease progression and/or disease stages from bulk RNA-seq data derived from
91 post-mortem tissue. Henceforth, we refer to manifold learning, lineage inference interchangeably in
92 reference to the construction of the inference of a neuropathologic progression tree. We demonstrate that
93 these tools can estimate the disease staging and progression tree (**Figure 2**) from bulk RNA-Seq data
94 collected from post-mortem brain tissues in a case/control cohort. Moreover, these trees show clear
95 LOAD staging, enable the study of cell type specific effects of LOAD, and allow the identification of
96 genetic factors driving disease progression.

97 **Results**

98 *Manifold learning distinguishes pathologically defined LOAD from control*

99 We first quantify the bulk RNA-Seq data from the ROS/MAP and Mayo Clinic cohorts into gene counts
100 and remove any batch effects introduced due to sequencing runs using standard count normalization (see
101 **Methods**). The data from the ROS/MAP cohort is sampled from the dorsolateral prefrontal cortex

102 (DLPFC), and the data from the Mayo Clinic cohort is sampled from the temporal cortex (TCX). Patient
103 clinical characteristics are reported in **Table S1** and described in Methods. The full pipeline we use for
104 RNA-Seq data generation and quality control was recently reported²². The entire transcriptome comprises
105 many genes which do not have measurable expression or vary across case/control samples, which we
106 remove in order to reduce the noise in manifold learning¹⁹. To do this, we first perform differential
107 expression analysis between case/control samples separately for each study and retain genes that reach an
108 FDR of 0.10. To test if this biased the disease lineage inference, we also perform manifold learning using
109 only genes with high variance across samples, and we see a strong concordance with disease lineages
110 inferred with differentially expressed genes (**Figure S1**). Changing the significance threshold for the
111 differential expression analysis to $FDR < 0.01$ did not materially change these results (**Figure S2**). We
112 infer the disease lineage for each brain region on this subset of retained genes (**Figure 2A-B**). Adjusting
113 for post-mortem interval (PMI) (**Figure S3A, S4A**), 10 principal components from a principal component
114 analysis of genotype data to account for ancestry effects (**Figure S3B, S4B**), RNA integrity number
115 (RIN) (**Figure S3C, S4C**), or all of these variables (**Figure S3D, S4D**) did not materially change the
116 overarching ordering of patients for either the TCX or DLPFC regions. Furthermore, to assess the general
117 robustness of the results we apply leave one out cross validation to infer disease pseudotime for both
118 DLPFC and TCX brain regions and find strong correlations between lineages inferred with each sample
119 removed, and the lineage for the entire sample set (**Figure S5**).

120
121 We visualize the clinical diagnosis of the samples on the inferred disease staging tree to verify that there
122 is indeed separation of AD patients across the tree. To determine if inferred tree structure is an accurate
123 model of disease progression, we introduce the notion of disease pseudotime which is the geodesic
124 distance along the tree from an inferred initial point to the point of interest as a quantitative linear
125 measure of LOAD stage. We scale this estimated disease pseudotime to lie in the range $[0,1]$ to make the
126 effects comparable between the two studies (and brain regions). We show that for LOAD cases compared
127 to controls there is a significant association ($P = 0.02$ in Mayo and $P = 2.0 \times 10^{-6}$ in ROS/MAP, logistic

128 regression) between the estimated pseudotime and AD case/control status (**Figure 2C**). These effects are
129 not abrogated by adjusting for RIN, PMI, or ancestry in either tissue (**Figure S6A-D, Table S2**). To
130 assess whether the association between inferred disease pseudotime is a phenomena in only the Mayo
131 RNAseq and ROS/MAP RNA-seq data, we also apply the lineage inference approach to expression array
132 data from the Mayo eQTL study²³ (**Methods**). These samples are derived from a completely independent
133 set of donors than the Mayo RNA-seq study²⁴. Similarly, we restrict to only female samples, and test for
134 an association between inferred disease pseudotime and disease status (**Figure S7A-B**). We see a
135 significant association between disease pseudotime and neuropathological AD diagnosis ($P = 2.2 \times 10^{-8}$).
136
137 Furthermore, we observe strong evidence of sex heterogeneity when performing the manifold learning
138 approach, and find that the manifolds inferred for female only samples show stronger association with
139 pseudotime than for male samples. This matches previous observations concerning disease specific sex
140 heterogeneity²². As such, we do not see as statistically significant of an association between pseudotime
141 and disease diagnosis in male samples ($P = 0.040$ in Mayo and $P = 0.11$ in ROS/MAP, logistic regression,
142 **Figure S8A-D**). Similarly, the association between pseudotime and amyloid, tau, and cognitive diagnosis
143 is attenuated in male samples in ROS/MAP (**Figure S9A-D**). For the combined samples we see moderate
144 evidence of disease association with pseudotime ($P = 0.003$ in Mayo and $P = 0.003$ in ROS/MAP, logistic
145 regression, **Figure S10A-D**). The association with neuropathological measures of disease is more robust
146 in the combined sample (**Figure S11A-F**), but not as strong as in females only, hence we restrict to
147 female only analyses for all subsequent reported results.

148
149 We test whether genes in loci that have been implicated in genome wide association studies of LOAD are
150 associated with inferred disease pseudotime. We use the prioritized LOAD GWAS genes²⁵, **Table S3**,
151 and compute the correlation between their expression and inferred pseudotime (**Figure 2D**). When
152 compared to the background of all genes, we see that there is a significant increase in positive correlation
153 with disease pseudotime for implicated LOAD GWAS genes (P -value: 7.3×10^{-5} in Mayo and 5.6×10^{-3} in

154 ROS/MAP). This effect is robust to adjustments for PMI, RIN, or ancestry (**Figure S12A-D**).

155 Furthermore, this does not appear to be driven by a small subset of outlier genes, but by the majority of

156 the distribution of LOAD GWAS genes. The fact that AD GWAS loci genes have expression

157 associations with pseudotime likely implies that the AD risk variants at these are also eQTL as previously

158 shown²⁶⁻²⁹ and/or are members of co-expression networks that are differentially expressed in AD^{13,30}.

159

160 To further explore the relationship between inferred disease stage and LOAD, we test for its association

161 with neuropathological and clinical measures of LOAD severity, namely: i) Braak score, ii) CERAD

162 score, and iii) cognitive diagnosis. The ROS/MAP study has numeric scores for these categories available

163 as covariates for each sample. Braak is a semi-quantitative measure that increases with tau pathology³¹

164 and CERAD is a semi-quantitative measure of density of neuritic plaques³². We overlay these scores on

165 the inferred manifold for the DLPFC brain region (**Figure 3A**). We observe a progressive increase in tau,

166 amyloid, and cognitive burden as we traverse the inferred disease manifold (**Figure 3A**). This is further

167 quantified by characterizing the relationship between branches of the inferred manifold and Braak,

168 CERAD, and cognitive diagnosis (**Figure 3B**). We observe significant associations between pseudotime

169 and Braak score ($P=1.0 \times 10^{-5}$), CERAD score ($P=1.8 \times 10^{-5}$), and cognitive diagnosis ($P=3.5 \times 10^{-7}$). In

170 ROS/MAP, adjustment for Braak score when fitting the DDRTree method attenuates the association

171 between pseudotime and disease states (P-value: 0.214, **Figure S13A-B**), though there is still evidence of

172 association with cognitive diagnosis (P-value: 0.03, **Figure S13C-D**). In the Mayo RNA-seq study we

173 have Braak score and Thal Amyloid scores for only a subset of samples, but observe a similar pattern as

174 in ROS/MAP (**Figure S14A-B**) for the samples that we do have data. There is a significant association

175 between Braak score and pseudotime (P-value: 5×10^{-5}) as well as Thal amyloid (P-value: 1.7×10^{-5}) within

176 this subset with available neuropathology data.

177

178 *Comparison to other unsupervised learning approaches.* We compare the manifold learning approach to

179 other unsupervised learning approaches including principal component analysis (PCA), t-distributed

180 Stochastic Neighbor Embedding (tSNE)³³, and Uniform Manifold Approximation and Projection
181 (UMAP)³⁴. Correlations between the first two dimensions of each of these approaches and the DDRTree
182 learned pseudotimes are shown in **Figures S15A-F, S16A-F** for DLPFC and TCX brain tissues
183 respectively. We see the strongest correlations between PCA1 and UMAP2 and pseudotime in both data-
184 sets, increasing our confidence that the overarching ordering of patients along a disease pseudotime is a
185 robust characteristic of the disease progression as reflected in the gene expression changes as a function
186 of disease, and not dependent on the underlying manifold learning approach. This is further supported by
187 inspecting the association between these approaches and Braak, CERAD, and cognitive diagnosis (cogdx)
188 scores for the DLPFC tissue (**Table S4**). Furthermore, the manifold learning approaches (DDRTree and
189 UMAP) have much stronger associations with Braak, CERAD, and cogdx scores than either PCA or
190 tSNE. In fact, UMAP has been proposed for lineage inference³⁵, and when we apply UMAP with lineage
191 inference using Monocle3, we observe similar results (**Figure S17-S18**) as with DDRTree and Monocle2
192 (**Figure 2-3**), though the inferred pseudotimes from Monocle3 are not quite as significant as the
193 association with UMAP2 or from Monocle 2 with Braak, CERAD, and cogdx (**Table S4**).

194 *Inferred staging recapitulates known biology of AD*

195 To demonstrate that the inferred disease pseudotime recapitulates known biology of LOAD, we test for
196 association between inferred disease stage and both the cellular response to disease and the genetics of the
197 disease. A prominent hypothesis in AD is that the effects of the disease vary across different brain cell
198 types, specifically neurons and glial subtypes. Current understanding of the cell biology of the disease
199 implicates progressive neuronal loss and increase in gliosis³⁶. To test if the inferred pseudotime aligns
200 with existing cell type specific hypotheses regarding AD, we first selected from the genes used in lineage
201 construction the marker genes for four key cell types: neurons, astrocytes, microglia, and
202 oligodendrocytes based on a previously published brain cell atlas³⁷ (**Table S5**). We then calculate the
203 normalized mean expression for the marker genes of each cell type and fitted a linear model to the mean
204 expression with disease pseudotime as the dependent variable. We find that, in both studies, the cell

205 specific marker gene levels show a statistically significant linear dependence on pseudotime (**Table S6**).
206 Fitted effects recapitulate known neuropathologic changes which occur in AD, namely: i) a reduction in
207 the neuronal populations as AD progresses, and ii) an increase in expression associated with activation of
208 microglia, astrocytes, and oligodendrocytes as AD progresses (**Figure 4**).

209

210 Next, we test for association between assigned lineage state in ROS/MAP (DLPFC) and Mayo (TCX) and
211 APOE e4 status (**Figure S19**). For reference, the inferred trees for TCX and DLPFC each resolve into 6
212 branches (**Figure 5A,S20**). Carriers of the APOE e4 allele are significantly enriched on the State 4
213 branch in TCX (P-value = 0.027, unadjusted), and suggestively enriched on the State 5 branch (P-value =
214 0.06, unadjusted), compared to the State 1 branch (logistic regression). Similarly, in the Mayo eGWAS
215 study, when we perform an ordinal logistic regression of APOE e4 dosage and disease pseudotime we see
216 a significant positive association as a function of pseudotime (P-value = 6.9×10^{-4} , **Figure S7C**).

217

218 *Genetic factors associated with inferred disease staging*

219 Lineage inference of LOAD transcriptomes provides a quantitative measure of disease progression for
220 genetic associate testing, and the significantly greater correlation between pseudotime and gene
221 expression for known LOAD risk genes (**Figure 2D**) suggests that the observed differences in disease
222 trajectories are influenced by genetic factors. To test this hypothesis, we perform single variant analysis
223 using whole-genome sequencing data for 305 patients from the ROS/MAP and 131 patients from the
224 Mayo cohort. Despite the limited sample size, resulting in lack of statistical power to discover genome-
225 wide significant associations, multiple variants reach a genome-wide suggestive threshold of $p < 1 \times 10^{-5}$
226 (**Table S7**). We do not see evidence of population stratification in the analysis (**Figure S21-S22**).

227 Notably, the most significant association with pseudotime for the ROS/MAP cohort is observed at the
228 *PTRPD* locus (rs7870388, $p = 1.31 \times 10^{-6}$) (**Figure S23, Table S7**). The *PTRPD* locus is associated with
229 the susceptibility to neurofibrillary tangle independent of amyloid deposition in the ROS/MAP cohort³⁸.

230 For the Mayo Clinic cohort, known LOAD variants in the *APOE* (rs6857, $p = 9.18 \times 10^{-6}$) and *BINI*

231 (rs62158731, $p = 4.68 \times 10^{-5}$) loci overlap with variants associated with inferred disease stage (**Figure**
232 **S23, Table S7**)³⁹. When comparing our association results for inferred disease stage with summary
233 statistics from a large-scale case-control approach, we identify multiple variants which have been
234 previously associated with LOAD in the IGAP cohort (**Table S8**). Furthermore, we identify several
235 potential novel candidate genes associated with inferred disease stage (*ADAMTS14*, *IL7*, *MAN2B1*) linked
236 to immune and lysosomal storage function (**Figure S23, Table S7**). *IL-7* has been proposed as an
237 inflammatory biomarker for LOAD that correlates with disease outcome and severity⁴⁰. *ADAMTS14* is
238 part of a locus that has been previously linked with Alzheimer susceptibility and plays an important role
239 in the regulation of immune function via TGF-beta signaling.

240 *New disease insights identified from inferred disease lineages*

241 Another important direction of study in the field of Alzheimer's is the identification of disease subtypes,
242 which has so far predominantly been done using imaging data⁴¹. The branches of the inferred disease trees
243 provide a new transcriptomic-based approach to identify disease subtypes. In both brain regions and in
244 two separate cohorts, there were two distinct early-lineage branches corresponding to predominantly
245 control samples, which we interpret as different initial paths towards the disease. Similarly, both brain
246 regions feature several distinct branches with predominantly LOAD samples (**Figure 2A-B**).

247

248 *Branch-specific differential expression patterns.* To study the genes and pathways specific to each branch,
249 we perform a branch-specific differential expression analysis with an ANOVA model using the branches
250 with the highest proportion of controls as the reference branch for DLPFC (**Table S9**) and TCX (**Table**
251 **S10**). We see many genes are differentially expressed between the control branch and branches that are
252 enriched in the affected individuals (**Table S11**). We test for overlap between the differentially expressed
253 gene sets between the two studies (**Figure S24**), and find significant overlaps in branches enriched for
254 late stage disease cases, especially between up-regulated genes in State 6 of DLPFC and up regulated
255 genes in State 5 of TCX (P-value: 4.1×10^{-108} , OR: 4.5, Fisher's exact test), as well as genes that are up-

256 regulated in State 6 of DLPFC and in State 4 of TCX (P-value: 1.8×10^{-14} , OR: 1.9, Fisher's exact test),
257 and more modestly for genes that are down regulated in State 6 of DLPFC and down regulated in State 3
258 of TCX (P-value: 1.1×10^{-6} , OR: 1.6, Fisher's exact test). Next, we performed an enrichment analysis on
259 each of these differentially expressed gene sets with the enrichR⁴² package for Gene Ontology⁴³
260 annotations (**Methods**). The results of this enrichment analysis for DLPFC and TCX tissues (**Table S12-**
261 **S13**). Only gene sets with significant enrichment are shown (FDR adjusted p-value < 0.05). Overall, we
262 see a pattern of loss of expression of basic cell biology mechanisms in early-stage branches including
263 RNA splicing, mitochondria function, protein transport, and DNA repair. Late-stage branches were
264 characterized by increased immune response (e.g. TGFb/WNT signaling) and apoptotic activity (**Table**
265 **1**).

266

267 While studying the different branches in the two brain regions, we observe a branch (branch 5) that
268 corresponds to a group of predominantly neuropathological control samples from the Mayo RNA-seq
269 cohort that were in close proximity to a branch with predominantly LOAD samples (branch 4) on the
270 inferred disease lineage (**Figure 5A**). However, most of the samples on branch 5 are neuropathological
271 controls as defined by the Mayo diagnostic criteria. We bi-cluster the mean expression of genes in each
272 branch and the branches themselves (**Methods**). This clustering analysis (**Figure 5B**) shows that the
273 closest branch to this potentially disease resistant branch contains the highest proportion of AD samples.
274 While the stage proximity implies some transcriptomic similarity between these controls and nearby
275 cases, we also see a secondary cluster of genes with increased expression in the resistant state while
276 having reduced expression in all other states. We perform an enrichment analysis on this set of genes and
277 find significant GO terms corresponding to: protein transport (GO:0015031), regulation of mRNA
278 splicing, via spliceosome (GO:0048024), negative regulation of apoptotic process (GO:0043066), and
279 regulation of amyloid-beta clearance (GO:1900221) (Cluster4, **Table S14**). It is possible that these
280 potentially disease resistant individuals have compensatory mechanisms which suppress the hallmarks of
281 disease despite sharing gene expression patterns with pathologically affected individuals.

282

283 To replicate this observation, we perform a differential expression analysis on individuals in the Mayo
284 eGWAS study where we consider individuals that are in the top quintile of pseudotime, but are classified
285 as neuropathological controls as resistant individuals (**Figure S7D, Methods**). To test if these individuals
286 also have a similar resistant molecular endophenotype, we compare the overlap between various
287 differential expressed gene sets derived from these resistant individuals and the gene sets identified in the
288 biclustering of the Mayo RNA-seq data (**Figure S25**). We observe that there is a highly statistically
289 significant overlap between genes that are upregulated in these Mayo eGWAS resistant individuals (P-
290 value: 2.6×10^{-51} , OR: 2.9, Fisher's exact test), and the Cluster 4 genes that are upregulated in the branch 5
291 Mayo RNA-seq samples (**Figure S25**).

292 **Discussion**

293 Here we proposed a novel approach to infer the Alzheimer's disease severity and disease subtypes in an
294 unsupervised manner from post-mortem bulk RNA-seq data that gets directly at the challenge of
295 identifying the temporal progression of disease in the disease resistant tissue. Our strategy utilized a
296 manifold learning approach to infer a disease progression tree from cross-sectionally collected patient
297 samples from two different brain regions. The underlying assumption of our approach is that the inferred
298 disease progression from cross sectional samples serves as a proxy for the unobserved progression of the
299 disease across subtypes of LOAD. We validated this hypothesis through comparisons with
300 neuropathological measures of disease stage severity and against known cell type specific effects caused
301 by the disease. While one could argue that the method is merely classifying patients as either disease
302 cases or controls based on expression signatures of the hallmarks of disease, we see at least three
303 advantages of this approach beyond that interpretation. First the application of this method appears to be
304 produce a more quantitative measure of disease state than strictly neuropathological assessments - as born
305 out through the identification of novel and distinct genetic loci that replicate based on IGAP summary
306 statistics. This suggests that it may be adding information related to other aspects of disease such as the

307 effect of neuroinflammation or neuronal injury. In addition we see evidence of neuropathological controls
308 that are disease resistant given their molecular state in two independent studies – which would not be
309 detectable with standard neuropathological or clinical assessments – and could provide important
310 molecular clues to mechanisms of disease resistance. Finally, there is the potential that specific pathways
311 associated with early stage disease processes can be characterized which is desperately needed for
312 hypothesis generation in the field. Furthermore, this approach provides clues to better understanding the
313 molecular heterogeneity of disease by identifying specific pathways that are dysregulated in subsets of
314 patients at different disease stages. This opens up the possibility of better patient stratification and
315 precision medicine.

316

317 We observed that different biological processes vary as a function of inferred disease stage, and that
318 early-stage disease processes potentially include RNA-splicing, mitochondrial function, and protein
319 transport – implicating multiple basic cell biology mechanisms as potential early stage disease processes
320 for further study in relevant model systems. Additionally, the manifold learning method identified 6
321 potential subtypes of LOAD from RNA-seq (i.e. branches) suggesting the LOAD populations should be
322 stratified by better biomarkers with tailored treatment strategies. To identify and test these stratifications
323 future studies should focus on longitudinal cohorts of patients with rich molecular and imaging data to be
324 able to identify biomarkers that can accurately and precisely stratify patients into the underlying
325 molecular subtypes in terms of the molecular characteristics of their transcriptome and different relative
326 stages of disease. Furthermore, we observe a potential disease resistant sub-type of patients. This disease
327 resistance should be tested in disease model systems, to identify if neuropathological readouts can be
328 modified by altering the function of the pathways identified in our analysis (e.g. APP processing, RNA
329 splicing, apoptosis, protein trafficking). While this preliminary observation needs to be validated in
330 another cohort, it has the potential to be a novel source of hypotheses concerning new therapeutic
331 development. Specifically, for constructing better combination therapy hypotheses that may confer
332 neuroprotection, even in patients that are mildly affected by disease.

333
334 LOAD is a complex and heterogeneous disease encompassing a broad spectrum of clinical symptoms.
335 Disease progression can vary widely between patients leading to different rates of cognitive decline.
336 Several lines of evidence suggest that these differences in progression are modified by multiple genetic
337 factors affecting the transition from one pathological state to another^{44,45}. However, it has remained
338 difficult to assess the role of genetic variants affecting disease trajectories by case-control approaches
339 alone. Here, we showed that our novel expression trait pseudotime might be used as a molecular
340 phenotype to identify known and novel AD loci associated with different disease progression states across
341 AD patients. Despite a limited sample size, we identified previously associated AD candidate loci in the
342 ROSMAP (*PTPRD*) and Mayo (*BINI*, *APOE*) cohorts with suggestive significance ($p < 1 \times 10^{-5}$).
343 Variants in *PTPRD* have been associated with the susceptibility to neurofibrillary tangles, independent
344 of amyloid burden. This is in line with the results from the differential gene expression analysis of
345 pseudotime branches showing an enrichment of molecular pathways implicated in TAU pathology.
346 Furthermore, our analysis revealed several novel loci linked to immune function (*ADAMTS14*, *IL7*) and
347 neurotransmitter signaling (*CHRM2*, *CHRM3*) processes associated with disease pseudotime (**Table S4**).
348 Future studies will be needed to replicate these findings in independent cohorts of LOAD and validate the
349 role of candidate genes in LOAD related disease progression by first identifying peripheral biomarkers
350 that correspond to this molecular definition of disease stage, and then testing for GWAS association with
351 that disease stage. Subsequent results can improve functional interpretation by linking candidate genes
352 with ordered pathological processes.

353

354 **Methods**

355 **ROS/MAP and Mayo RNAseq study population characteristics**

356 Detailed descriptions of cohort and patient characteristics included in this study can be found in
357 previously published work^{46,47}. Patient characteristics included in this study are summarized in **Table S1**,

358 stratified by sex. In brief: for Mayo samples, AD diagnosis was performed according to NINCDS-
359 ADRDA criteria (probable or possible AD); control individuals had Braak NFT stage ≤ 3 , CERAD score
360 < 2.5 , and lacked other pathologic diagnoses; Path.Aging are individuals who lacked any pathologic
361 diagnoses and had Braak NFT ≤ 3 and CERAD score ≥ 2 . Progressive supranuclear palsy (PSP)
362 individuals were diagnosed neuropathologically by a single neuropathologist. (For further details, see
363 Allen et al 2016). For ROSMAP samples, AD diagnosis was according to NIA Reagan criteria, which
364 combines neuropathology and clinical data; Control individuals had no signs of cognitive impairment;
365 and Other individuals had MCI, mixed pathology, or other form of dementia. Age at death was collected
366 for all patients in the ROSMAP study, though age at first AD diagnosis had a high degree of missingness
367 and thus was not used as a variable in follow-up analyses (see **Table S1**). Braak stage indicates the
368 measure of severity of NFT pathology. Stages I and II indicate NFTs confined mainly to the entorhinal
369 region of the brain, stages III and IV indicate involvement of limbic regions, and V and VI indicate
370 moderate to severe neocortical involvement. CERAD score is a semiquantitative measure of neuritic
371 plaques. 1=definite AD; 2=probable AD; 3=possible AD; 4=No AD. Cognitive diagnostic category (as
372 determined by neurologist): 1) NCI; 2) MCI and no other cause of CI; 3) MCI and another cause of CI; 4)
373 AD and no other cause of CI; 5) AD and another cause of CI; 6) Other primary cause of dementia⁴⁸. Thal
374 amyloid stages: phases of amyloid deposition. 0) no amyloid; 1) isocortical phase; 2) limbic phase; 3)
375 basal ganglia phase; 4) basal forebrain and midbrain phase; 5) pons/medulla oblongata and cerebellum
376 phase. The list of differentially expressed genes used to create the Monocle objects were based on results
377 which included the whole data set.

378

379 *RNA sequencing*

380 The details of the sample collections, postmortem sample characteristics, the tissue and RNA
381 preparations, the library preparations and sequencing technology and parameters, and sample quality
382 control filters are provided in previously published work^{46,49}. Furthermore, details of the bioinformatic
383 pipeline used to generate count level data has been previously described²². Briefly, reads were aligned to

384 the GENCODE24 (GRCh38) reference genome with STAR⁵⁰, and gene counts generated using the
385 HTSeq algorithm⁵¹. Genes that had more than one counts per million total reads total reads in at least
386 50% of samples in each tissue and diagnosis category were used for further analysis.

387

388 *Differential Expression analysis on Mayo and ROS/MAP cohorts*

389 For gene filtering we used false discovery rate of 0.05 from the previously published differential
390 expression analysis of Mayo and ROS/MAP RNA seq data²². Briefly, case control status was harmonized
391 across the Mayo and ROS/MAP cohorts, where controls were defined as individuals with a low burden of
392 amyloid and tau based on CERAD and Braak scores, and cases with a high burden. Furthermore in
393 ROS/MAP, clinical diagnosis was also used with controls having to have no cognitive impairment, and
394 cases have probably AD²². Differential expression analysis was run on suitably normalized data – using
395 conditional quantile normalization to account for variation in gene length and GC content, removing
396 sample outliers, covariate identification adjustment, with sampling abundance confidence estimated using
397 a weighted linear model with the voom-limma package^{22,52,53}. A fixed/mixed effect linear model is used
398 to fit the differential expression model on the normalized data²².

399

400 *Manifold learning for LOAD*

401 Manifold learning refers to a group of machine learning algorithms that recover a low dimensional
402 subspace underlying a high dimensional dataset. Manifold learning approaches are typically used in
403 datasets or applications where data samples lie on an underlying low dimensional latent space (e.g. a tree,
404 a line, a curved plane). The low dimensional space is learned via a projection from the high dimensional
405 space of the observed data (e.g. RNA-seq profiles across hundreds of patient samples) down to a low
406 dimensional space with suitable regularization constraints to enforce smoothness and the structural
407 constraints of the low dimensional space. (**Figure 1A**). Due to the necessary assumption of an underlying
408 latent subspace, manifold learning is commonly used in applications where it is known that the observed
409 data is obtained from a progression of some kind; e.g., i) to infer the temporal ordering of a sequence of

410 images, or ii) to infer the approximate lineage of cells in a differentiation trajectory using single cell
411 RNA-Seq data (**Figure 1B-C**).
412
413 Here, we repurpose methods originally developed for learning cell lineage using scRNA-Seq data, to infer
414 the staging of Alzheimer's disease (AD) using bulk RNA-Seq data from post-mortem brain samples with
415 known AD diagnosis status. Since bulk RNA-Seq has many of the same sampling and distributional
416 properties as scRNA-Seq, we observe that scRNA-Seq methods are applicable with no additional
417 modifications. As such, we use the DDRTree manifold learning approach available in the Monocle2 R
418 package⁵⁴. However, we also show that the estimated staging of disease is quite similar across some of the
419 other common methods used for scRNA-Seq lineage estimation (**Figures S26-S28**) including Monocle1⁵⁵
420 and diffusion pseudotime (DPT)⁵⁶.

421
422 The RNA-Seq data used in this study was generated from post-mortem brain homogenate samples, and
423 obtained from two separate studies that are a part of the Accelerating Medicines Partnership in
424 Alzheimer's Disease (AMP-AD) consortium, namely: i) the Religious Orders Study and the Memory and
425 Aging Project (ROSMAP)^{57,58}, and ii) the Mayo RNA-seq study⁴⁹. For this paper, we focused our analysis
426 on the temporal cortex (TCX) and dorsolateral prefrontal cortex (DLPFC) tissue samples. Within the
427 Mayo RNA-seq study the TCX samples are derived from individuals neuropathologically defined as
428 either aged controls, LOAD cases, Progressive Supranuclear Palsy (PSP) cases, or pathological aging
429 (PA) cases⁴⁹. The ROSMAP study is a prospective longitudinal cohort of an aging population, and has
430 samples from participants with clinical and neuropathological diagnoses of LOAD⁴⁶, aged controls, and
431 individuals with mild cognitive impairment. Furthermore, most results presented in the main paper are
432 from female samples only unless indicated otherwise, as we observed significant sex differences in the
433 transcriptomic data consistent with current knowledge of sex differences in LOAD^{59,60}. For replication,
434 we also consider microarray data generated from the Illumina DASL gene expression platform from the
435 Mayo eGWAS study from TCX for N=186 patients, of which 108 were neuropathologically confirmed

436 AD Cases and 78 were controls²³. Probes were mapped to genes using BioMart. Data was adjusted for
 437 plate using ordinary least squares regression prior to manifold learning.

438

439 *Manifold learning using Discriminative Dimensionality Reduction Tree (DDRTree)*

440 DDRTree is a manifold learning algorithm that infers a smooth low dimensional manifold by an approach
 441 called reverse graph embedding. Briefly, the algorithm simultaneously learns a non-linear projection to a
 442 latent space where the points lie on a spanning tree. A reverse embedding is also simultaneously learned
 443 from the latent space to the high dimensional data. Mathematically, the DDRTree algorithm can be posed
 444 as the following optimization problem:

$$\sum_{i=1}^N ||x_i - Wz_i||^2 + \frac{\lambda}{2} \sum_{k,k'} b_{k,k'} ||Wy_k - Wy_{k'}||^2$$

$$+ \gamma \left[\sum_{k=1}^K \sum_{i=1}^N r_{i,k} ||z_i - y_k||^2 + \sigma r_{i,k} \log (r_{i,k}) \right]$$

s. t. B represents a spanning tree,

$$W^T W = I, r_{i,k} \geq 0, \sum_{k=1}^K r_{i,k} = 1$$

445

446 Here, $\{z_i\}_{i=1}^N \in R^{genes}$ represents RNA-Seq data from each patient sample, $\{z_i\}_{i=1}^N \in R^2$ represents the

447 latent representation of each sample as inferred by the algorithm, $\{y_k\}_{k=1}^K$ represents the centers of

448 clusters in the dataset, $W \in R^{2 \times genes}$ represents an inverse mapping from the latent space to the high

449 dimensional space of RNA-Seq data, $B \in R^{K \times K}$ represents a spanning tree on which the centers of the

450 clusters lie and $R \in R^{N \times K}$ captures the soft clustering information of samples in the dataset. The first term

451 of the optimization problem is responsible for learning a low dimensional representation of the data such

452 that an inverse mapping exists to the high dimensional data points, the second term learns the tree

453 structure of the points and the third term learns a soft clustering for the latent dimension points as well as

454 the centers of the clusters. Despite the non-convexity of the problem, each individual optimization
455 variable can be solved for efficiently using alternative minimization as described previously⁶¹. This
456 algorithm was implemented using the Monocle package in R¹⁹. When fitting the Monocle objects we also
457 considered various adjustments to the expression data prior to manifold learning for 10 principal
458 components of genetic ancestry, RNA integrity number, post mortem interval, age, Braak score - among
459 other potential confounds. The code to infer the lineage in MayoRNAseq is available here
460 [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_GenerateMonocleDS_new.R)
461 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_GenerateMonocleDS_new.R](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_GenerateMonocleDS_new.R), and code used
462 to infer the lineage in ROSMAP is available here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_GenerateMonocleDS_new.R)
463 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_GenerateMonocleDS_new.R](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_GenerateMonocleDS_new.R). Code to
464 perform analysis on Mayo eGWAS study is available here: [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_ben_april_2020/mayo_egwas_GenerateMonocleDS_new.R)
465 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_ben_april_2020/mayo_egwas_GenerateMonocleDS](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_ben_april_2020/mayo_egwas_GenerateMonocleDS_new.R)
466 [_new.R](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_ben_april_2020/mayo_egwas_GenerateMonocleDS_new.R).

467

468 *Branch assignment and pseudotime calculation for samples*

469 Branch assignment and pseudotime calculation was also performed using the *Monocle* package using
470 techniques described previously¹⁹. Briefly, pseudotime is calculated by first identifying a root point on
471 one of the two ends of the maximum diameter path in the tree. Then the pseudotime of each point is
472 calculated by projecting it to its closest point on the spanning tree and calculating the geodesic distance to
473 the root point. Assigning samples to branches is done by first identifying the branches of the spanning
474 tree and then assigning samples to the branch on which their projection to the spanning tree lies on.
475 Robustness of pseudotime was assessed with leave-one-out cross validation by dropping one sample at a
476 time, running the DDRTree method with Monocle, and then computing the absolute value of the
477 correlation between the pseudotime estimated with the reduced data-set, and the pseudotime estimated
478 with the full data set. Alternative approaches for performing dimensionality reduction included principal
479 component analysis (PCA), t-stochastic neighborhood embedding (tSNE)³³, and uniform manifold

480 approximation and projection (UMAP)³⁴, which were all run on the same data set as the DDRTree method
481 was run in R (here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_GenerateMonocleDS_new.R)
482 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_GenerateMonocleDS_new.R](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_GenerateMonocleDS_new.R) and
483 [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_GenerateMonocleDS_new.R)
484 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_GenerateMonocleDS_new.R](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_GenerateMonocleDS_new.R)).

485

486 *Association of pseudotime with AD status, hallmarks of Alzheimer's disease, and cognitive diagnosis*
487 We test for association between disease pseudotime and AD case or control status with logistic regression
488 with AD case or control status as the outcome and inferred pseudotime as the dependent variable in both
489 the Mayo and ROS/MAP studies. We test for association between pseudotime and hallmarks of disease
490 in the ROS/MAP studies for both Braak (measure of tau pathology) score and CERAD score (measure of
491 amyloid pathology) with an ordinal logistic regression model, with the neuropath score as the ordered
492 outcome, and pseudotime as the dependent variable. Finally, we test for association between disease
493 pseudotime and cognitive diagnosis for the following ordered clinical diagnoses of no cognitive
494 impairment, mild cognitive impairment, and probable Alzheimer's disease with an ordinal logistic
495 regression model. All code for running these association tests is available [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/paper_figures.Rmd)
496 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/paper_figures.Rmd](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/paper_figures.Rmd).

497

498 *Inferring cell type specific expression patterns given marker gene expression as a function of pseudotime*
499 List of marker genes for different major cell types in the brain was curated from a previously published
500 brain cell expression signature study³⁷. The marker gene list was then pruned to include only genes that
501 were included in lineage construction. Each gene's expression as a function of pseudotime was then
502 obtained by smoothing using a smoothing spline of degree of freedom = 3 and normalized to lie in [0,1].
503 The smoothing was done to remove the effects of technical noise introduced due to RNA-Seq and the
504 normalization was done since the absolute expression levels of genes might be very different from each
505 other. The smoothed and normalized expression of marker genes for each category were then averaged to

506 obtain the average marker gene expression as a function of pseudotime. A linear model was used to test
507 for association between average expression of a given cell type expression signature and pseudotime.

508

509 *Association between GWAS loci and correlation with pseudotime*

510 To test for association between pseudotime and LOAD GWAS genes, we computed the Spearman's
511 correlation between each gene's expression and pseudotime in the Mayo and ROS/MAP studies. Next,
512 we considered the 60 highly prioritized genes (priority score greater than 4) identified within AD GWAS
513 loci by the International Genetics of Alzheimer's Project (IGAP)²⁵. We test for a difference between the
514 correlation with pseudotime of background of all other genes and the IGAP AD genes using a linear
515 model, and see a significant increase in correlation between gene expression and pseudotime in both the
516 Mayo and ROS/MAP study for AD GWAS genes.

517

518 *Branch specific differential expression analysis*

519 We perform a state specific differential expression analysis using a one-way ANOVA model in both the
520 Mayo and ROS/MAP studies. The branch with the highest proportion of AD controls is defined as the
521 reference branch for all analyses. We use Tukey's honest significant difference method to compute P-
522 values for the test for change in expression of a given gene compared to the reference branch. Genes are
523 grouped based on their branch and direction of change in expression for further downstream pathway
524 enrichment analyses. Overlap between differential expressed genes was depicted using UpSet plots⁶²
525 (**Figure S24**). Code to run analyses are available here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_DE_Anova.R)
526 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_DE_Anova.R](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/DLPFC_DE_Anova.R) for ROS/MAP and here
527 https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_DE_Anova.R for
528 Mayo.

529

530 *Branch specific gene expression signatures*

531 Branch specific expression signature was obtained by first calculating the average normalized expression
532 for all genes in each state/branch. This was followed by performing a bi-clustering using the pheatmap
533 package in R (<https://cran.r-project.org/web/packages/pheatmap/index.html>) which uses hierarchical
534 clustering on both samples and genes. We also used the pheatmap R package to visualize the state
535 specific expression signatures.

536

537 *Disease resistant subgroup validation analysis*

538 After identifying potentially disease resistant individuals in the Mayo RNAseq study based on the TCX
539 brain region (individuals from branch 5, **Figure 5A**), we considered the Mayo eGWAS TCX data, and
540 defined neuropathological controls with pseudotimes in the top quintile of all pseudotimes as disease
541 resistant (n=9). We then performed a differential expression analysis using linear regression to identify
542 array probes that were differentially expressed between resistant and non-resistant individuals, of which
543 there were more than 5000 probes that were either up or down regulated at an FDR of 0.05. Overlaps
544 were explored between the branch specific gene clusters from Mayo RNAseq (Figure 5B) and these Mayo
545 eGWAS resistance differential expressed probes using UpSet plots⁶² (**Figure S25**).

546

547 *Gene set enrichment analyses*

548 For each branch specific differential expression gene set (DEGs) in both Mayo RNAseq and ROS/MAP
549 we perform a gene set enrichment analysis against Gene Ontology pathways using the enrichR⁴² R
550 package. Only pathways with FDR < 0.05 are reported. The code we used to run the ROS/MAP DEG
551 enrichments are available here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/lineage.Rmd)
552 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/lineage.Rmd](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/lineage.Rmd), the code we used to run the Mayo
553 DEG enrichments are available here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/lineageTCX.Rmd)
554 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/lineageTCX.Rmd](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/lineageTCX.Rmd), and the code we used to run the
555 branch specific gene expression signature pathway enrichments is available here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/resilience.Rmd)
556 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/resilience.Rmd](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/resilience.Rmd).

557

558 *Whole-genome sequencing*

559 Whole-genome sequencing was performed at the New York Genome Center for all individuals from the
560 ROS/MAP and Mayo cohorts. Detailed information for both data sets can be accessed via synapse
561 (DOI:10.7303/syn2580853). Briefly, 650ng of genomic DNA from whole blood was sheared using a
562 Covaris LE220 sonicator. DNA fragments underwent bead-based size selection and were subsequently
563 end-repaired, adenylated, and ligated to Illumina sequencing adapters. Libraries were sequenced on an
564 Illumina HiSeq X sequencer using 2 x 150bp cycles. Paired-end reads were aligned to the GRCh37
565 (hg19) human reference genome using the Burrows-Wheeler Aligner (BWA-MEM v0.7.8) and processed
566 using the GATK best-practices workflow^{63,64}. This included marking of duplicate reads by the use of
567 Picard tools v1.83, local realignment around indels, and base quality score recalibration (BQSR) via
568 Genome Analysis Toolkit (GATK v3.4.0). Joint variant calling files (vcfs) for whole-genome sequencing
569 data for the Mayo and ROS/MAP cohort were obtained through the AMP-AD knowledge portal
570 (DOI:10.7303/syn10901595).

571

572 *Single variant association with pseudotime in two independent cohorts*

573 Likelihood ratio tests within a linear regression framework were used to model the relationship between
574 the quantitative expression trait pseudotime and genetic variants in 436 AD cases. Genome-wide genetic
575 association analysis was performed for 305 female patients in the ROS/MAP cohort and 131 female
576 patients in the Mayo cohort for which both genotyping and post-mortem RNA-seq data was available. An
577 efficient mixed model approach, implemented in the EMMAX software suite, was used to account for
578 potential biases and cryptic relatedness among individuals⁶⁵. Only variants with MAF > 0.05, genotyping
579 call rates > 95%, minimum sequencing depth of 20 reads and Hardy-Weinberg equilibrium $p > 10^{-4}$ were
580 considered for analysis. Quantile-quantile plots (**Figure S21-S22**) for the test statistics showed no
581 significant deviation between expected and observed p-values, highlighting that there is no consistent
582 differences across cases and controls except for the small number of significantly associated variants.

583 Furthermore, the genomic inflation factor (λ) was determined to be 0.99 for the Mayo and 0.98 for
584 the ROS/MAP single variant association tests. This highlights that potential confounding factors, such as
585 population stratification have been adequately controlled.

586

587

588 **Data Availability**

589 All data analyzed in the study is publicly available and described in detail in previous publications^{22-24,46}.

590 Specifically, we use a version of the RNA-seq data from the ROS/MAP study

591 (DOI:10.7303/syn8456638.22) and RNA-seq data from the Mayo RNAseq (10.7303/syn8466816.19) run

592 through the same bioinformatic processing pipeline²². The array expression data from the Mayo eGWAS

593 study is available at DOI:10.7303/syn3617054.1.

594

595 **Code Availability**

596 All code is publicly available (https://github.com/Sage-Bionetworks/AMPAD_Lineage). References to

597 code to perform specific analyses are described in detail in the **Methods**.

598

599 References

- 600 1. Cummings, J. L., Morstorf, T. & Zhong, K. Alzheimer's disease drug-development pipeline: few
601 candidates, frequent failures. *Alzheimers. Res. Ther.* **6**, 37 (2014).
- 602 2. Cummings, J. L., Doody, R. & Clark, C. Disease-modifying therapies for Alzheimer disease:
603 Challenges to early intervention. *Neurology* **69**, 1622–1634 (2007).
- 604 3. Ferreira, D. *et al.* Distinct subtypes of Alzheimer's disease based on patterns of brain atrophy:
605 longitudinal trajectories and clinical applications. *Sci. Rep.* **7**, 46263 (2017).
- 606 4. Bredesen, D. E. Metabolic profiling distinguishes three subtypes of Alzheimer's disease. *Aging*
607 (*Albany NY*) **7**, 595 (2015).
- 608 5. Brier, M. R. *et al.* Tau and A β imaging, CSF measures, and cognition in Alzheimer's disease.
609 *Sci. Transl. Med.* **8**, 338ra66--338ra66 (2016).
- 610 6. Gordon, B. A. *et al.* The relationship between cerebrospinal fluid markers of Alzheimer pathology
611 and positron emission tomography tau imaging. *Brain* **139**, 2249–2260 (2016).
- 612 7. Dichgans, M. *et al.* METACOHORTS for the study of vascular disease and its contribution to
613 cognitive decline and neurodegeneration: An initiative of the Joint Programme for
614 Neurodegenerative Disease Research. *Alzheimer's Dement.* **12**, 1235–1249 (2016).
- 615 8. Jack, C. R. *et al.* Tracking pathophysiological processes in Alzheimer's disease: an updated
616 hypothetical model of dynamic biomarkers. *Lancet. Neurol.* **12**, 207–16 (2013).
- 617 9. Au, R., Piers, R. J. & Lancashire, L. Back to the future: Alzheimer's disease heterogeneity
618 revisited. *Alzheimer's Dement. (Amsterdam, Netherlands)* **1**, 368–370 (2015).
- 619 10. Carrasquillo, M. M. *et al.* Late-onset Alzheimer's risk variants in memory decline, incident mild
620 cognitive impairment, and Alzheimer's disease. *Neurobiol. Aging* **36**, 60–7 (2015).
- 621 11. Zhang, B. *et al.* Integrated systems approach identifies genetic nodes and networks in late-onset
622 Alzheimer's disease. *Cell* **153**, 707–20 (2013).
- 623 12. Mostafavi, S. *et al.* A molecular network of the aging human brain provides insights into the
624 pathology and cognitive decline of Alzheimer's disease. *Nat. Neurosci.* **21**, 811–819 (2018).

- 625 13. Conway, O. J. *et al.* ABI3 and PLCG2 missense variants as risk factors for neurodegenerative
626 diseases in Caucasians and African Americans. *Mol. Neurodegener.* **13**, 53 (2018).
- 627 14. Allen, M. *et al.* Conserved brain myelination networks are altered in Alzheimer’s and other
628 neurodegenerative diseases. *Alzheimers. Dement.* **14**, 352–366 (2018).
- 629 15. Allen, M. *et al.* Divergent brain gene expression patterns associate with distinct cell-specific tau
630 neuropathology traits in progressive supranuclear palsy. *Acta Neuropathol.* **136**, 709–727 (2018).
- 631 16. Bengio, Y., Courville, A. & Vincent, P. Representation learning: A review and new perspectives.
632 *IEEE Trans. Pattern Anal. Mach. Intell.* **35**, 1798–1828 (2013).
- 633 17. Wolz, R., Aljabar, P., Hajnal, J. V. & Rueckert, D. Manifold Learning for Biomarker Discovery in
634 MR Imaging. in 116–123 (Springer, Berlin, Heidelberg, 2010). doi:10.1007/978-3-642-15948-
635 0_15
- 636 18. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
637 pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, (2014).
- 638 19. Qiu, X. *et al.* Reversed graph embedding resolves complex single-cell trajectories. *Nat. Methods*
639 **14**, 979–982 (2017).
- 640 20. Rosenberg, A. B. *et al.* Single-cell profiling of the developing mouse brain and spinal cord with
641 split-pool barcoding. *Science (80-.).* **360**, 176–182 (2018).
- 642 21. Mukherjee, S., Zhang, Y., Fan, J., Seelig, G. & Kannan, S. Scalable preprocessing for sparse
643 scRNA-seq data exploiting prior knowledge. *Bioinformatics* **34**, i124–i132 (2018).
- 644 22. Wan, Y.-W. *et al.* Meta-Analysis of the Alzheimer’s Disease Human Brain Transcriptome and
645 Functional Dissection in Mouse Models. *Cell Rep.* **32**, 107908 (2020).
- 646 23. Zou, F. *et al.* Brain Expression Genome-Wide Association Study (eGWAS) Identifies Human
647 Disease-Associated Variants. *PLoS Genet.* **8**, e1002707 (2012).
- 648 24. Allen, M. *et al.* Human whole genome genotype and transcriptome data for Alzheimer’s and other
649 neurodegenerative diseases. *Sci. Data* **3**, (2016).
- 650 25. Kunkle, B. W. *et al.* Genetic meta-analysis of diagnosed Alzheimer’s disease identifies new risk

- 651 loci and implicates A β , tau, immunity and lipid processing. *Nat. Genet.* **51**, 414–430 (2019).
- 652 26. Allen, M. *et al.* Novel late-onset Alzheimer disease loci variants associate with brain gene
653 expression. *Neurology* **79**, 221–8 (2012).
- 654 27. Allen, M. *et al.* Late-onset Alzheimer disease risk variants mark brain regulatory loci. *Neurol.*
655 *Genet.* **1**, e15 (2015).
- 656 28. Allen, M. *et al.* Association of MAPT haplotypes with Alzheimer’s disease risk and MAPT brain
657 gene expression levels. *Alzheimers. Res. Ther.* **6**, 39 (2014).
- 658 29. Carrasquillo, M. M. *et al.* A candidate regulatory variant at the TREM gene cluster associates with
659 decreased Alzheimer’s disease risk and increased TREML1 and TREM2 brain gene expression.
660 *Alzheimers. Dement.* **13**, 663–673 (2017).
- 661 30. Sims, R. *et al.* Rare coding variants in PLCG2, ABI3, and TREM2 implicate microglial-mediated
662 innate immunity in Alzheimer’s disease. *Nat. Genet.* **49**, 1373–1384 (2017).
- 663 31. Braak, H., Alafuzoff, I., Arzberger, T., Kretschmar, H. & Del Tredici, K. Staging of Alzheimer
664 disease-associated neurofibrillary pathology using paraffin sections and immunocytochemistry.
665 *Acta Neuropathol.* **112**, 389–404 (2006).
- 666 32. Mirra, S. S. *et al.* The Consortium to Establish a Registry for Alzheimer’s Disease (CERAD). Part
667 II. Standardization of the neuropathologic assessment of Alzheimer’s disease. *Neurology* **41**, 479–
668 86 (1991).
- 669 33. Maaten, L. van der & Hinton, G. Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9**, 2579–
670 2605 (2008).
- 671 34. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for
672 Dimension Reduction. (2018).
- 673 35. Cao, J. *et al.* The single-cell transcriptional landscape of mammalian organogenesis. *Nature* **566**,
674 496–502 (2019).
- 675 36. De Strooper, B. & Karran, E. The Cellular Phase of Alzheimer’s Disease. *Cell* **164**, 603–615
676 (2016).

- 677 37. Zhang, Y. *et al.* An RNA-sequencing transcriptome and splicing database of glia, neurons, and
678 vascular cells of the cerebral cortex. *J. Neurosci.* **34**, 11929–47 (2014).
- 679 38. Chibnik, L. B. *et al.* Susceptibility to neurofibrillary tangles: role of the PTPRD locus and limited
680 pleiotropy with other neuropathologies. *Mol. Psychiatry* **23**, 1521 (2018).
- 681 39. Lambert, J.-C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for
682 Alzheimer’s disease. *Nat. Genet.* **45**, 1452 (2013).
- 683 40. Janelidze, S. *et al.* CSF biomarkers of neuroinflammation and cerebrovascular dysfunction in early
684 Alzheimer disease. *Neurology* **91**, e867–e877 (2018).
- 685 41. Whitwell, J. L. *et al.* Neuroimaging correlates of pathologically defined subtypes of Alzheimer’s
686 disease: a case-control study. *Lancet Neurol.* **11**, 868–877 (2012).
- 687 42. Chen, E. Y. *et al.* Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool.
688 *BMC Bioinformatics* **14**, 128 (2013).
- 689 43. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology
690 Consortium. *Nat. Genet.* **25**, 25–9 (2000).
- 691 44. Dong, A. *et al.* Heterogeneity of neuroanatomical patterns in prodromal Alzheimer’s disease: links
692 to cognition, progression and biomarkers. *Brain* **140**, aww319 (2016).
- 693 45. Wang, X. *et al.* Genetic determinants of disease progression in Alzheimer’s disease. *J. Alzheimers.*
694 *Dis.* **43**, 649–55 (2015).
- 695 46. De Jager, P. L. *et al.* A multi-omic atlas of the human frontal cortex for aging and Alzheimer’s
696 disease research. *Sci. Data* **5**, 180142 (2018).
- 697 47. Allen, M. *et al.* Human whole genome genotype and transcriptome data for Alzheimer’s and other
698 neurodegenerative diseases. *Sci. data* **3**, 160089 (2016).
- 699 48. Schneider, J. A., Arvanitakis, Z., Bang, W. & Bennett, D. A. Mixed brain pathologies account for
700 most dementia cases in community-dwelling older persons. *Neurology* **69**, 2197–2204 (2007).
- 701 49. Allen, M. *et al.* Human whole genome genotype and transcriptome data for Alzheimer’s and other
702 neurodegenerative diseases. *Sci. Data* **3**, 160089 (2016).

- 703 50. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 704 51. Anders, S., Pyl, P. T. & Huber, W. HTSeq—a Python framework to work with high-throughput
705 sequencing data. *Bioinformatics* **31**, 166–169 (2015).
- 706 52. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: Precision weights unlock linear model
707 analysis tools for RNA-seq read counts. *Genome Biol.* **15**, (2014).
- 708 53. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and
709 microarray studies. *Nucleic Acids Res.* **43**, e47–e47 (2015).
- 710 54. Qiu, X. *et al.* Reversed graph embedding resolves complex single-cell trajectories. *Nat. Methods*
711 **14**, 979–982 (2017).
- 712 55. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
713 pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
- 714 56. Haghverdi, L., Büttner, M., Wolf, F. A., Buettner, F. & Theis, F. J. Diffusion pseudotime robustly
715 reconstructs lineage branching. *Nat. Methods* **13**, 845–848 (2016).
- 716 57. Bennett, D. *et al.* Overview and findings from the rush Memory and Aging Project. *Curr.*
717 *Alzheimer Res.* **9**, 646–663 (2012).
- 718 58. Bennett, D. A. *et al.* Religious orders study and rush memory and aging project. *J. Alzheimer's*
719 *Dis.* 1–28 (2018).
- 720 59. Ferretti, M. T. *et al.* Sex differences in Alzheimer disease—the gateway to precision medicine.
721 *Nat. Rev. Neurol.* 1 (2018).
- 722 60. Deming, Y. *et al.* Sex-specific genetic predictors of Alzheimer's disease biomarkers. *Acta*
723 *Neuropathol.* **136**, 857–872 (2018).
- 724 61. Mao, Q., Wang, L., Goodison, S. & Sun, Y. Dimensionality Reduction Via Graph Structure
725 Learning. in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge*
726 *Discovery and Data Mining - KDD '15* 765–774 (ACM Press, 2015).
727 doi:10.1145/2783258.2783309
- 728 62. Conway, J. R., Lex, A. & Gehlenborg, N. UpSetR: an R package for the visualization of

- 729 intersecting sets and their properties. *Bioinformatics* **33**, 2938–2940 (2017).
- 730 63. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
731 generation DNA sequencing data. *Genome Res.* **20**, 1297–303 (2010).
- 732 64. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
733 *Bioinformatics* **25**, 1754–60 (2009).
- 734 65. Kang, H. M. *et al.* Variance component model to account for sample structure in genome-wide
735 association studies. *Nat. Genet.* **42**, 348–354 (2010).
- 736

737 **Acknowledgements**

738 This work was supported by NIA grants U54AG054345 and RFIAG057443. The ROSMAP Study data
739 were provided by the Rush Alzheimer's Disease Center, Rush University Medical Center, Chicago. Data
740 collection was supported through funding by NIA grants P30AG10161, R01AG15819, R01AG17917,
741 R01AG30146, R01AG36836, U01AG32984, U01AG46152, the Illinois Department of Public Health,
742 and the Translational Genomics Research Institute. Mayo RNAseq Study data were provided by the
743 following sources: The Mayo Clinic Alzheimer's Disease Genetic Studies, led by Dr. Nilufer Ertekin-
744 Taner and Dr. Steven G. Younkin, Mayo Clinic, Jacksonville, FL using samples from the Mayo Clinic
745 Study of Aging, the Mayo Clinic Alzheimer's Disease Research Center, and the Mayo Clinic Brain Bank.
746 Data collection was supported through funding by NIA grants P50 AG016574, R01 AG032990, U01
747 AG046139, R01 AG018023, U01 AG006576, U01 AG006786, R01 AG025711, R01 AG017216, R01
748 AG003949, NINDS grant R01 NS080820, CurePSP Foundation, and support from Mayo Foundation.
749 Study data includes samples collected through the Sun Health Research Institute Brain and Body
750 Donation Program of Sun City, Arizona. The Brain and Body Donation Program is supported by the
751 National Institute of Neurological Disorders and Stroke (U24 NS072026 National Brain and Tissue
752 Resource for Parkinson's Disease and Related Disorders), the National Institute on Aging (P30 AG19610
753 Arizona Alzheimer's Disease Core Center), the Arizona Department of Health Services (contract 211002,
754 Arizona Alzheimer's Research Center), the Arizona Biomedical Research Commission (contracts 4001,
755 0011, 05-901 and 1001 to the Arizona Parkinson's Disease Consortium) and the Michael J. Fox
756 Foundation for Parkinson's Research. MSBB data were generated from postmortem brain tissue collected
757 through the Mount Sinai VA Medical Center Brain Bank and were provided by Dr. Eric Schadt from
758 Mount Sinai School of Medicine.

759

760 **Author Contributions**

761 S.M. and B.A.L. designed the study. S.M., L.H., and B.A.L. performed the analyses. S.M., L.H., C.P.,

762 S.J., G.G., A.K.G., S.K.S., P.L.D.J., N.E.T., G.W.C., L.M.M., and B.A.L. wrote the manuscript.

763

764 **Tables**

765 **Table 1** - Representative significant Gene Ontology pathway enrichments (FDR<0.05) of differentially
 766 expressed genes for each branch (FDR < 0.05). Differential expressed genes are identified with an
 767 ANOVA analysis, with Branch 1 as the reference.

Brain Region	Direction	Branch	Representative Enriched Gene Ontology Terms
TCX	Down	2	prespliceosome (GO:0071010), mitochondrial electron transport, cytochrome c to oxygen (GO: 0006123)
		3	negative regulation of microtubule polymerization or depolymerization (GO:0031111)
		4	mitochondrial electron transport, NADH to ubiquinone (GO: 0006120), spliceosomal tri-snRNP complex (GO:0097526), negative regulation of microtubule depolymerization (GO:0007026)
		5	axon (GO:0030424), protein kinase C activity (GO:0004697),
		6	gamma-tubulin large complex (GO:0000931), U1 snRNP (GO:0005685), mitochondrial respiratory chain complex IV (GO:0005751), response to cadmium ion (GO:0046686)
		Up	2
	3	fatty acid elongase activity (GO:0009922), ubiquitin protein ligase activity (GO:0061630)	
	4	transforming growth factor beta-activated receptor activity (GO:0005024), hippo signaling (GO:0035329), regulation of extrinsic apoptotic signaling pathway via death domain receptors (GO: 1902041), regulation of DNA repair (GO: 0006282)	

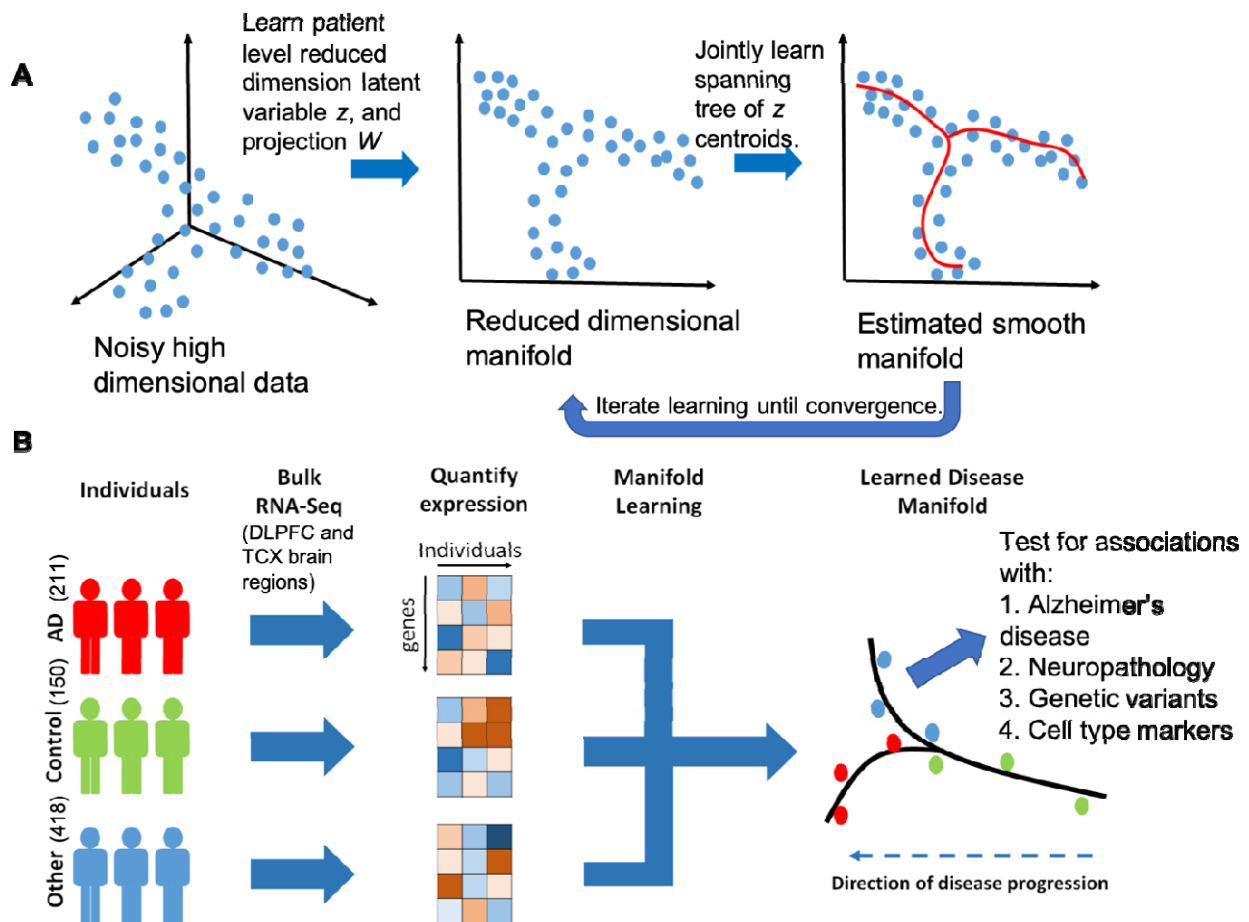
		5	regulation of apoptotic process (GO:0042981), leptin mediated signaling pathway (GO:0033210), negative regulation of hippo signaling (GO:0035331), small GTPase binding (GO:0031267)
		6	extracellular ligand-gated ion channel activity (GO:0005230), integral component of mitochondrial inner membrane (GO:0031305)
DLPFC	Down	2	DNA repair (GO:0006281), intracellular protein transport (GO:0006886)
		3	mismatch repair complex binding (GO:0032404)
		4	
		5	mitochondrial respiratory chain complex assembly (GO: 0033108)
		6	
	Up	2	racemase and epimerase activity (GO: 0016857)
		3	racemase and epimerase activity (GO: 0016857)
		4	vesicle mediated transport (GO: 0016192)
		5	NuRD complex (GO: 0016581)
		6	microtubule motor activity (GO:0003777), AP-2 adaptor complex binding (GO:0035612)

768

769

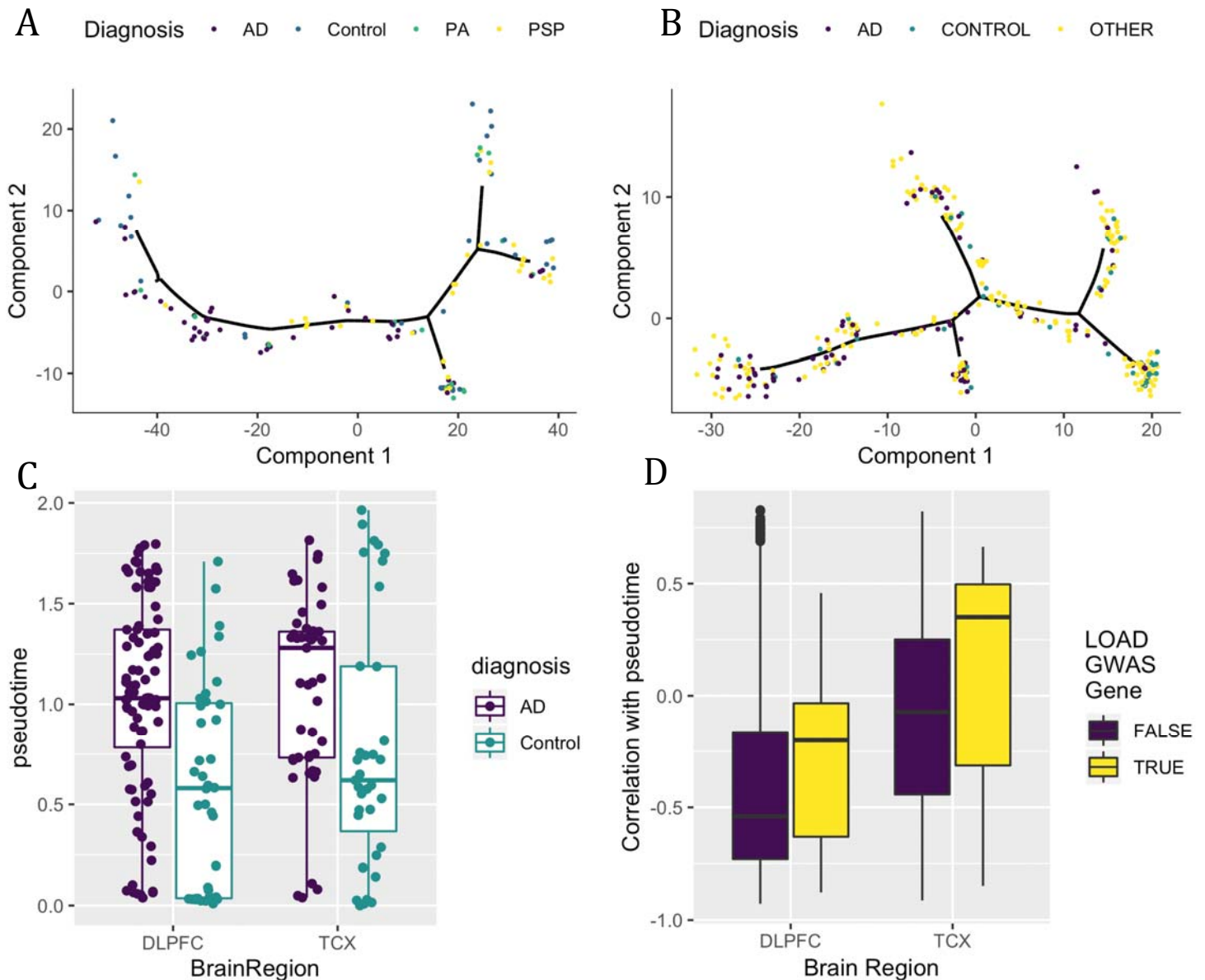
770 **Figures**

771 **Figure 1** - Overview of manifold learning for unraveling staging in Alzheimer's disease. A) Illustration of
 772 steps in manifold learning using reverse graph embedding DDRTree method. B) Illustration of lineage
 773 inference process for LOAD. RNA-seq samples with different disease diagnoses were pooled, batch
 774 normalized, and a smooth manifold was learned for each brain region across individuals (each point is an
 775 individual). Total sample numbers are indicated across Mayo RNaseq TCX and ROS/MAP DLPFC for
 776 the different diagnoses in parentheses.

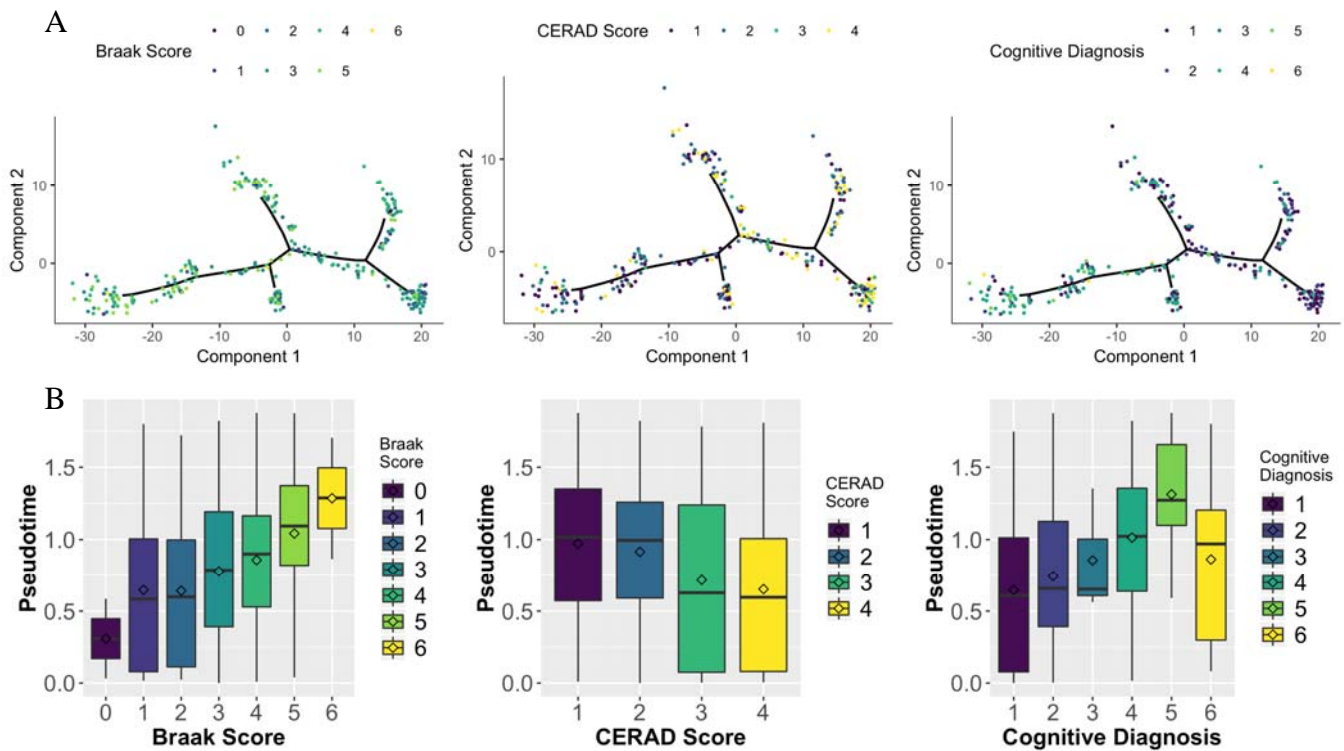


777
778

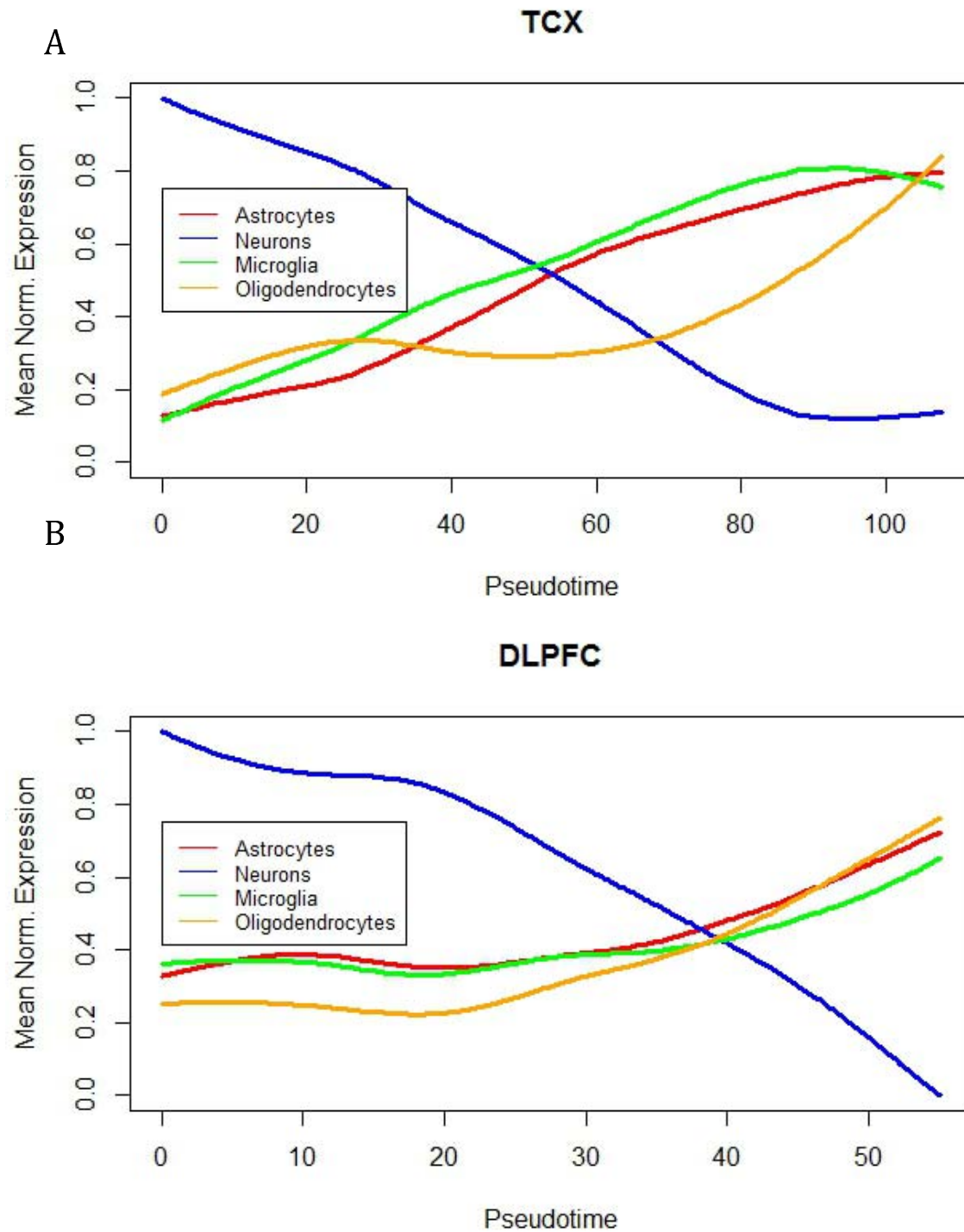
779 **Figure 2** - Manifold learning accurately infers disease states and stages from RNA-seq samples. A)
780 Estimated disease progression trees from temporal cortex (TCX) and B) dorsolateral prefrontal cortex
781 (DLPFC) brain regions showing localization of identified LOAD samples on particular branches. C)
782 Distribution of pseudotime for AD cases and controls for both DLPFC and TCX. D) Distribution of
783 expression correlation with pseudotime for both LOAD GWAS genes and non-LOAD GWAS genes.



784 **Figure 3** - Manifold learning replicates existing measures of staging in LOAD in DLPFC samples. A)
785 Samples colored by 3 different external measures of LOAD staging namely: Braak Score (tau pathology),
786 CERAD Score (amyloid pathology) and Cognitive Diagnosis (Clinical measure of disease severity).
787 Black lines denote inferred lineages. B) Distribution of samples by inferred stage for different distinct
788 stages in each of the three methods of measuring LOAD severity. Inferred disease stages generally
789 corresponded with all methods, and Cognitive diagnosis demonstrated the strongest alignment.



790 **Figure 4** – Cell type gene expression signatures as a function of disease pseudotime. A) Mean expression
791 of cell markers for astrocytes, neurons, microglia, and oligodendrocytes as a function of pseudotime for
792 TCX brain region, B) mean expression of cell markers for astrocytes, neurons, microglia,
793 oligodendrocytes as a function of pseudotime for DLPFC (B) brain region.
794



795 **Figure 5** – Disease resistant state. A) The inferred manifold from the TCX region with samples colored
796 by their inferred disease subtype/state. State 5 (dots, circled) lies at the late end of the disease trajectory,
797 indicating a strong disease-like transcriptomic phenotype, yet most samples in the group did not have
798 pathologically diagnosed AD (Figure 2A). We hypothesize this group represents a disease resistant state
799 to the disease. B) Biclustering results of average expression from each disease state, with increased
800 expression of a gene cluster (Cluster 4) unique to State 5.

