

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

2

3 Sumit Mukherjee^{1,2}, 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
5 Carter³, Lara M Mangravite¹, Benjamin A Logsdon^{1,*}

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,
12 Columbia 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

19 *Correspondence:

20 Benjamin A Logsdon, PhD

21 Sage Bionetworks

22 2901 3rd Avenue Suite 330

23 Seattle, WA 98121

24 ben.logsdon@sagebionetworks.org

25

26 **Abstract**

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

46 **Introduction**

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

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

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

83

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

96 **Results**

97 *Unsupervised manifold learning distinguishes pathologically defined LOAD from control*

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

101 (DLPFC), and the data from the Mayo Clinic cohort is sampled from the temporal cortex (TCX). The full
102 pipeline we use for RNA-Seq data generation and quality control was recently reported²². The entire
103 transcriptome comprises many genes which do not have measurable expression or vary across
104 case/control samples, which we remove in order to reduce the noise in manifold learning¹⁹. To do this, we
105 first perform differential expression analysis between case/control samples separately for each study and
106 retain genes that reach an FDR of 0.10. To test if this biased the disease lineage inference, we also
107 perform manifold learning using only genes with high variance across samples, and we see a strong
108 concordance with disease lineages inferred with differentially expressed genes (**Figure S4**). We infer the
109 disease lineage for each brain region on this subset of retained genes (**Figure 2A-B**). Furthermore, we
110 observe strong evidence of sex heterogeneity when performing the manifold learning approach, and find
111 that the manifolds inferred for female only samples are much more robust than for male samples. This
112 matches previous observations concerning disease specific sex heterogeneity²². We only show results for
113 manifolds inferred on female samples.

114

115 We first visualize the clinical diagnosis of the samples on the inferred disease staging tree to verify that
116 there is indeed separation of AD patients across the tree. To determine if inferred tree structure is an
117 accurate model of disease progression, we introduce the notion of disease pseudotime which is the
118 geodesic distance along the tree from an inferred initial point to the point of interest as a quantitative
119 linear measure of LOAD stage. We scale this estimated disease pseudotime to lie in the range [0,1] to
120 make the effects comparable between the two studies (and brain regions). We show that for LOAD cases
121 compared to controls there is a significant association ($P = 0.02$ in Mayo and $P = 2.0 \times 10^{-6}$ in ROS/MAP,
122 logistic regression) between the estimated pseudotime and AD case/control status (**Figure 2C**).

123

124 We test whether genes in loci that have been implicated in genome wide association studies of LOAD are
125 associated with inferred disease pseudotime. We use the highest ranked LOAD GWAS genes (60 genes
126 in total)²³, **Table S1**, and compute the correlation between their expression and inferred pseudotime

127 **(Figure 2D)**. When compared to the background of all genes, we see that there is a significant increase in
128 positive correlation with disease pseudotime for implicated LOAD GWAS genes (P-value: 7.3×10^{-5} in
129 Mayo and 5.6×10^{-3} in ROS/MAP). Furthermore, this does not appear to be driven by a small subset of
130 outlier genes, but by the majority of the distribution of LOAD GWAS genes. The fact that AD GWAS
131 loci genes have expression associations with pseudotime likely implies that the AD risk variants at these
132 are also eQTL as previously shown²⁴⁻²⁷ and/or are members of co-expression networks that are
133 differentially expressed in AD^{13,28}.

134
135 To further explore the relationship between inferred disease stage and LOAD, we test for its association
136 with neuropathological and clinical measures of LOAD severity, namely: i) Braak score, ii) CERAD
137 score, and iii) cognitive diagnosis. The ROSMAP study has numeric scores for these categories available
138 as covariates for each sample. Braak is a semi-quantitative measure that increases with tau pathology²⁹
139 and CERAD is a semi-quantitative measure of density of neuritic plaques³⁰. We overlay these scores on
140 the inferred manifold for the DLPFC brain region (**Figure 3A**). We observe a progressive increase in tau,
141 amyloid, and cognitive burden as we traverse the inferred disease manifold (**Figure 3A**). This is further
142 quantified by characterizing the relationship between branches of the inferred manifold and Braak,
143 CERAD, and cognitive diagnosis (**Figure 3B**). We observe significant associations between pseudotime
144 and Braak score ($P=1.0 \times 10^{-5}$), CERAD score ($P=1.8 \times 10^{-5}$), and cognitive diagnosis ($P=3.5 \times 10^{-7}$).

145 *Inferred staging recapitulates known biology of AD*

146 To demonstrate that the inferred disease pseudotime recapitulates known biology of LOAD, we test for
147 association between inferred disease stage and both the cellular response to disease and the genetics of the
148 disease. A prominent hypothesis in AD is that the effects of the disease vary across different brain cell
149 types, specifically neurons and glial subtypes. Current understanding of the cell biology of the disease
150 implicates progressive neuronal loss and increase in gliosis³¹. To test if the inferred pseudotime aligns
151 with existing cell type specific hypotheses regarding AD, we first selected from the genes used in lineage

152 construction the marker genes for four key cell types: neurons, astrocytes, microglia, and
153 oligodendrocytes based on a previously published brain cell atlas³² (**Table S2**). We then calculate the
154 normalized mean expression for the marker genes of each cell type and fitted a linear model to the mean
155 expression with disease pseudotime as the dependent variable. We find that, in both studies, the cell
156 specific marker gene levels show a statistically significant linear dependence on pseudotime (**Table S3**).
157 Fitted effects recapitulate known neuropathologic changes which occur in AD, namely: i) a reduction in
158 the neuronal populations as AD progresses, and ii) an increase in expression associated with activation of
159 microglia, astrocytes, and oligodendrocytes as AD progresses (**Figure S5**).

160
161 Next, we test for association between assigned lineage state in ROS/MAP (DLPFC) and Mayo (TCX) and
162 APOE e4 status (**Figure S6**). For reference, the inferred trees for TCX and DLPFC each resolve into 6
163 branches (**Figure 4A,S7**). Carriers of the APOE e4 allele are significantly enriched on the State 4 branch
164 in TCX (P-value = 0.027, unadjusted), and suggestively enriched on the State 5 branch (P-value = 0.06,
165 unadjusted), compared to the State 1 branch (logistic regression).

166
167 *Genetic factors associated with inferred disease staging*

168 Lineage inference of LOAD transcriptomes provides a quantitative measure of disease progression for
169 genetic associate testing, and the significantly greater correlation between pseudotime and gene
170 expression for known LOAD risk genes (**Figure 2D**) suggests that the observed differences in disease
171 trajectories are influenced by genetic factors. To test this hypothesis, we perform single variant analysis
172 using whole-genome sequencing data for 305 patients from the ROS/MAP and 131 patients from the
173 Mayo cohort. Despite the limited sample size, resulting in lack of statistical power to discover genome-
174 wide significant associations, multiple variants reach a genome-wide suggestive threshold of $p < 1 \times 10^{-5}$
175 (**Table S4**). We do not see evidence of population stratification in the analysis (**Figure S8-S9**). Notably,
176 the most significant association with pseudotime for the ROS/MAP cohort is observed at the *PTRPD*
177 locus (rs7870388, $p = 1.31 \times 10^{-6}$) (**Figure S10, Table S4**). The *PTRPD* locus is associated with the

178 susceptibility to neurofibrillary tangle independent of amyloid deposition in the ROS/MAP cohort³³. For
179 the Mayo Clinic cohort, known LOAD variants in the *APOE* (rs6857, $p = 9.18 \times 10^{-6}$) and *BINI*
180 (rs62158731, $p = 4.68 \times 10^{-5}$) loci overlap with variants associated with inferred disease stage (**Figure**
181 **S10, Table S4**)³⁴. When comparing our association results for inferred disease stage with summary
182 statistics from a large-scale case-control approach, we identify multiple variants which have been
183 previously associated with LOAD in the IGAP cohort (**Table S5**). Furthermore, we identify several
184 potential novel candidate genes associated with inferred disease stage (*ADAMTS14*, *IL7*, *MAN2B1*) linked
185 to immune and lysosomal storage function (**Figure S10, Table S4**). *IL-7* has been proposed as an
186 inflammatory biomarker for LOAD that correlates with disease outcome and severity³⁵. *ADAMTS14* is
187 part of a locus that has been previously linked with Alzheimer susceptibility and plays an important role
188 in the regulation of immune function via TGF-beta signaling.

189 *New disease insights identified from inferred disease lineages*

190 Another important direction of study in the field of Alzheimer's is the identification of disease subtypes,
191 which has so far predominantly been done using imaging data³⁶. The branches of the inferred disease trees
192 provide a new transcriptomic-based approach to identify disease subtypes. In both brain regions and in
193 two separate cohorts, there were two distinct early-lineage branches corresponding to predominantly
194 control samples, which we interpret as different initial paths towards the disease. Similarly, both brain
195 regions feature several distinct branches with predominantly LOAD samples (**Figure 2A-B**).

196

197 *Branch-specific differential expression patterns.* To study the genes and pathways specific to each branch,
198 we perform a branch-specific differential expression analysis with an ANOVA model using the branches
199 with the highest proportion of controls as the reference branch for DLPFC (**Table S6**) and TCX (**Table**
200 **S7**). We see many genes are differentially expressed between the control branch and branches that are
201 enriched in the affected individuals (**Table S8**). Next, we performed an enrichment analysis on each of
202 these differentially expressed gene sets with the enrichR³⁷ package for Gene Ontology³⁸ annotations

203 **(Methods)**. The results of this enrichment analysis for DLPFC and TCX tissues (**Table S9-S10**). Only
204 gene sets with significant enrichment are shown (FDR adjusted p-value < 0.05). Overall, we see a pattern
205 of loss of expression of basic cell biology mechanisms in early-stage branches including RNA splicing,
206 mitochondria function, protein transport, and DNA repair. Late-stage branches were characterized by
207 increased immune response (e.g. TGFb/WNT signaling) and apoptotic activity (**Table 1**).

208
209 While studying the different branches in the two brain regions, we observe a branch (branch 5) that
210 corresponds to a group of predominantly neuropathological control samples from the Mayo RNA-seq
211 cohort that were in close proximity to a branch with predominantly LOAD samples (branch 4) on the
212 inferred disease lineage (**Figure 4A**). However, most of the samples on branch 5 are neuropathological
213 controls as defined by the Mayo diagnostic criteria. We bi-cluster the mean expression of genes in each
214 branch and the branches themselves (**Methods**). This clustering analysis (**Figure 4B**) shows that the
215 closest branch to this potentially disease resistant branch contains the highest proportion of AD samples.
216 While the stage proximity implies some transcriptomic similarity between these controls and nearby
217 cases, we also see a secondary cluster of genes with increased expression in the resistant state while
218 having reduced expression in all other states. We perform an enrichment analysis on this set of genes and
219 find significant GO terms corresponding to: protein transport (GO:0015031), regulation of mRNA
220 splicing, via spliceosome (GO:0048024), negative regulation of apoptotic process (GO:0043066), and
221 regulation of amyloid-beta clearance (GO:1900221) (Cluster4, **Table S11**). It is possible that these
222 potentially disease resistant individuals have compensatory mechanisms which suppress the hallmarks of
223 disease despite sharing gene expression patterns with pathologically affected individuals. These
224 observations are preliminary, and would need to be replicated in a second cohort to verify the existence of
225 a disease resistant gene expression signature.

226 **Discussion**

227 Here we proposed a novel approach to infer the Alzheimer's disease severity and disease subtypes in an
228 unsupervised manner from post-mortem bulk RNA-seq data that gets directly at the challenge of
229 identifying the temporal progression of disease in the disease resident tissue. Our strategy utilized a
230 manifold learning approach to infer a disease progression tree from cross-sectionally collected patient
231 samples from two different brain regions. The underlying assumption of our approach is that the inferred
232 disease progression from cross sectional samples serves as a proxy for the unobserved progression of the
233 disease across subtypes of LOAD. We validated this hypothesis through comparisons with
234 neuropathological measures of disease stage severity and against known cell type specific effects caused
235 by the disease. Furthermore, this approach provides clues to better understanding the molecular
236 heterogeneity of disease by identifying specific pathways that are dysregulated in subsets of patients at
237 different disease stages. This opens up the possibility of better patient stratification and precision
238 medicine.

239
240 We observed that different biological processes vary as a function of inferred disease stage, and that
241 early-stage disease processes potentially include RNA-splicing, mitochondrial function, and protein
242 transport – implicating multiple basic cell biology mechanisms as potential early stage disease processes
243 for further study in relevant model systems. Additionally, the manifold learning method identified 6
244 potential subtypes of LOAD from RNA-seq (i.e. branches) suggesting the LOAD populations should be
245 stratified by better biomarkers with tailored treatment strategies. To identify and test these stratifications
246 future studies should focus on longitudinal cohorts of patients with rich molecular and imaging data to be
247 able to identify biomarkers that can accurately and precisely stratify patients into the underlying
248 molecular subtypes in terms of the molecular characteristics of their transcriptome and different relative
249 stages of disease. Furthermore, we observe a potential disease resistant sub-type of patients. This disease
250 resistance should be tested in disease model systems, to identify if neuropathological readouts can be

251 modified by altering the function of the pathways identified in our analysis (e.g. APP processing, RNA
252 splicing, apoptosis, protein trafficking). While this preliminary observation needs to be validated in
253 another cohort, it has the potential to be a novel source of hypotheses concerning new therapeutic
254 development. Specifically, for constructing better combination therapy hypotheses that may confer
255 neuroprotection, even in patients that are mildly affected by disease.
256

257 LOAD is a complex and heterogeneous disease encompassing a broad spectrum of clinical symptoms.
258 Disease progression can vary widely between patients leading to different rates of cognitive decline.
259 Several lines of evidence suggest that these differences in progression are modified by multiple genetic
260 factors affecting the transition from one pathological state to another^{39,40}. However, it has remained
261 difficult to assess the role of genetic variants affecting disease trajectories by case-control approaches
262 alone. Here, we showed that our novel expression trait pseudotime might be used as a molecular
263 phenotype to identify known and novel AD loci associated with different disease progression states across
264 AD patients. Despite a limited sample size, we identified previously associated AD candidate loci in the
265 ROSMAP (*PTPRD*) and Mayo (*BINI*, *APOE*) cohorts with suggestive significance ($p < 1 \times 10^{-5}$).
266 Variants in *PTPRD* have been associated with the susceptibility to neurofibrillary tangles, independent
267 of amyloid burden. This is in line with the results from the differential gene expression analysis of
268 pseudotime branches showing an enrichment of molecular pathways implicated in TAU pathology.
269 Furthermore, our analysis revealed several novel loci linked to immune function (*ADAMTS14*, *IL7*) and
270 neurotransmitter signaling (*CHRM2*, *CHRM3*) processes associated with disease pseudotime (**Table S4**).
271 Future studies will be needed to replicate these findings in independent cohorts of LOAD and validate the
272 role of candidate genes in LOAD related disease progression by first identifying peripheral biomarkers
273 that correspond to this molecular definition of disease stage, and then testing for GWAS association with
274 that disease stage. Subsequent results can improve functional interpretation by linking candidate genes
275 with ordered pathological processes.
276

277 **Methods**

278 *RNA sequencing*

279 The details of the sample collections, postmortem sample characteristics, the tissue and RNA
280 preparations, the library preparations and sequencing technology and parameters, and sample quality
281 control filters are provided in previously published work^{41,42}. Furthermore, details of the bioinformatic
282 pipeline used to generate count level data has been previously described²². Briefly, reads were aligned to
283 the GENCODE24 (GRCh38) reference genome with STAR⁴³, and gene counts generated using the
284 HTSeq algorithm⁴⁴. Genes that had more than one counts per million total reads total reads in at least
285 50% of samples in each tissue and diagnosis category were used for further analysis.

286

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

288 For gene filtering we used false discovery rate of 0.05 from the previously published differential
289 expression analysis of Mayo and ROS/MAP RNA seq data²². Briefly, case control status was harmonized
290 across the Mayo and ROS/MAP cohorts, where controls were defined as individuals with a low burden of
291 amyloid and tau based on CERAD and Braak scores, and cases with a high burden. Furthermore in
292 ROS/MAP, clinical diagnosis was also used with controls having to have no cognitive impairment, and
293 cases have probably AD²². Differential expression analysis was run on suitably normalized data – using
294 conditional quantile normalization to account for variation in gene length and GC content, removing
295 sample outliers, covariate identification adjustment, with sampling abundance confidence estimated using
296 a weighted linear model with the voom-limma package^{22,45,46}. A fixed/mixed effect linear model is used
297 to fit the differential expression model on the normalized data²².

298

299 *Manifold learning for LOAD*

300 Manifold learning refers to a group of machine learning algorithms that recover a low dimensional
301 subspace underlying a high dimensional dataset. Manifold learning approaches are typically used in

302 datasets or applications where data samples lie on an underlying low dimensional latent space (e.g. a tree,
303 a line, a curved plane). The low dimensional space is learned via a projection from the high dimensional
304 space of the observed data (e.g. RNA-seq profiles across hundreds of patient samples) down to a low
305 dimensional space with suitable regularization constraints to enforce smoothness and the structural
306 constraints of the low dimensional space. (**Figure 1A**). Due to the necessary assumption of an underlying
307 latent subspace, manifold learning is commonly used in applications where it is known that the observed
308 data is obtained from a progression of some kind; e.g., i) to infer the temporal ordering of a sequence of
309 images, or ii) to infer the approximate lineage of cells in a differentiation trajectory using single cell
310 RNA-Seq data (**Figure 1B-C**).

311
312 Here, we repurpose methods originally developed for learning cell lineage using scRNA-Seq data, to infer
313 the staging of Alzheimer's disease (AD) using bulk RNA-Seq data from post-mortem brain samples with
314 known AD diagnosis status. Since bulk RNA-Seq has many of the same sampling and distributional
315 properties as scRNA-Seq, we observe that scRNA-Seq methods are applicable with no additional
316 modifications. As such, we use the DDRTree manifold learning approach available in the Monocle 2 R
317 package¹⁹. However, we also show that the estimated staging of disease is quite similar across some of
318 the other common methods used for scRNA-Seq lineage estimation (**Figures S1-S3**).

319
320 The RNA-Seq data used in this study was generated from post-mortem brain homogenate samples, and
321 obtained from two separate studies that are a part of the Accelerating Medicines Partnership in
322 Alzheimer's Disease (AMP-AD) consortium, namely: i) the Religious Orders Study and the Memory and
323 Aging Project (ROSMAP)^{47,48}, and ii) the Mayo RNA-seq study⁴⁹. For this paper, we focused our
324 analysis on the temporal cortex (TCX) and dorsolateral prefrontal cortex (DLPFC) tissue samples. Within
325 the Mayo RNA-seq study the TCX samples are derived from individuals neuropathologically defined as
326 either aged controls, LOAD cases, Progressive Supranuclear Palsy (PSP) cases, or pathological aging
327 (PA) cases⁴⁹. The ROSMAP study is a prospective longitudinal cohort of an aging population, and has

328 samples from participants with clinical and neuropathological diagnoses of LOAD⁴², aged controls, and
 329 individuals with mild cognitive impairment. Furthermore, the results presented in the main paper are
 330 from female samples only, as we observed significant sex differences in the transcriptomic data consistent
 331 with current knowledge of sex differences in LOAD^{50,51}, making a common analysis of both sexes
 332 untenable.

333

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

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

$$\min_{W,Z,B,Y,R} \sum_{i=1}^N \|x_i - Wz_i\|^2 + \frac{\lambda}{2} \sum_{k,k'} b_{k,k'} \|Wy_k - Wy_{k'}\|^2 + \gamma \left[\sum_{k=1}^K \sum_{i=1}^N r_{i,k} \|z_i - y_k\|^2 + \sigma r_{i,k} \log(r_{i,k}) \right]$$

s. t. B represents a spanning tree,

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

340

341 Here, $\{z_i\}_{i=1}^N \in \mathbb{R}^{genes}$ represents RNA-Seq data from each patient sample, $\{z_i\}_{i=1}^N \in \mathbb{R}^2$ represents the
 342 latent representation of each sample as inferred by the algorithm, $\{y_k\}_{k=1}^K$ represents the centers of
 343 clusters in the dataset, $W \in \mathbb{R}^{2 \times genes}$ represents an inverse mapping from the latent space to the high
 344 dimensional space of RNA-Seq data, $B \in \mathbb{R}^{K \times K}$ represents a spanning tree on which the centers of the
 345 clusters lie and $R \in \mathbb{R}^{N \times K}$ captures the soft clustering information of samples in the dataset. The first
 346 term of the optimization problem is responsible for learning a low dimensional representation of the data
 347 such that an inverse mapping exists to the high dimensional data points, the second term learns the tree

348 structure of the points and the third term learns a soft clustering for the latent dimension points as well as
349 the centers of the clusters. Despite the non-convexity of the problem, each individual optimization
350 variable can be solved for efficiently using alternative minimization as described previously⁵². This
351 algorithm was implemented using the Monocle package in R¹⁹. The code to infer the lineage in Mayo is
352 available here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_GenerateMonocleDS_new.R)
353 [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
354 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)
355 [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).

356

357 *Branch assignment and pseudotime calculation for samples*

358 Branch assignment and pseudotime calculation was also performed using the *Monocle* package using
359 techniques described previously¹⁹. Briefly, pseudotime is calculated by first identifying a root point on
360 one of the two ends of the maximum diameter path in the tree. Then the pseudotime of each point is
361 calculated by projecting it to its closest point on the spanning tree and calculating the geodesic distance to
362 the root point. Assigning samples to branches is done by first identifying the branches of the spanning
363 tree and then assigning samples to the branch on which their projection to the spanning tree lies on.

364

365 *Association of pseudotime with AD status, hallmarks of Alzheimer's disease, and cognitive diagnosis*

366 We test for association between disease pseudotime and AD case or control status with logistic regression
367 with AD case or control status as the outcome and inferred pseudotime as the dependent variable in both
368 the Mayo and ROS/MAP studies. We test for association between pseudotime and hallmarks of disease
369 in the ROS/MAP studies for both Braak (measure of tau pathology) score and CERAD score (measure of
370 amyloid pathology) with an ordinal logistic regression model, with the neuropath score as the ordered
371 outcome, and pseudotime as the dependent variable. Finally, we test for association between disease
372 pseudotime and cognitive diagnosis for the following ordered clinical diagnoses of no cognitive
373 impairment, mild cognitive impairment, and probable Alzheimer's disease with an ordinal logistic

374 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)
375 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/paper_figures.Rmd](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/paper_figures.Rmd).

376

377 *Inferring cell type specific expression patterns given marker gene expression as a function of pseudotime*

378 List of marker genes for different major cell types in the brain was curated from a previously published
379 brain cell expression signature study³². The marker gene list was then pruned to include only genes that
380 were included in lineage construction. Each gene's expression as a function of pseudotime was then
381 obtained by smoothing using a smoothing spline of degree of freedom = 3 and normalized to lie in [0,1].

382 The smoothing was done to remove the effects of technical noise introduced due to RNA-Seq and the
383 normalization was done since the absolute expression levels of genes might be very different from each
384 other. The smoothed and normalized expression of marker genes for each category were then averaged to
385 obtain the average marker gene expression as a function of pseudotime. A linear model was used to test
386 for association between average expression of a given cell type expression signature and pseudotime.

387

388 *Association between GWAS loci and correlation with pseudotime*

389 To test for association between pseudotime and LOAD GWAS genes, we computed the Spearman's
390 correlation between each gene's expression and pseudotime in the Mayo and ROS/MAP studies. Next,
391 we identified a set of genes implicated in AD GWAS loci in the International Genetics of Alzheimer's
392 Project (IGAP)²³. We treated the set of genes described in Tables 1-3 of that study as high quality
393 candidate AD GWAS genes²³. We test for a difference between the correlation with pseudotime of
394 background of all other genes and the IGAP AD genes using a linear model, and see a significant increase
395 in correlation between gene expression and pseudotime in both the Mayo and ROS/MAP study for AD
396 GWAS genes.

397

398 *Branch specific differential expression analysis*

399 We perform a state specific differential expression analysis using a one-way ANOVA model in both the
400 Mayo and ROS/MAP studies. The branch with the highest proportion of AD controls is defined as the
401 reference branch for all analyses. We use Tukey's honest significant difference method to compute P-
402 values for the test for change in expression of a given gene compared to the reference branch. Genes are
403 grouped based on their branch and direction of change in expression for further downstream pathway
404 enrichment analyses. 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)
405 [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
406 https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/TCX_DE_Anova.R for
407 Mayo.

408

409 *Estimating branch specific gene expression signatures*

410 Branch specific expression signature was obtained by first calculating the average normalized expression
411 for all genes in each state/branch. This was followed by performing a bi-clustering using the pheatmap
412 package in R which uses hierarchical clustering on both samples and genes. We also used the pheatmap
413 package to visualize the state specific expression signatures.

414

415 *Gene set enrichment analyses*

416 For each branch specific differential expression gene set (DEGs) in both Mayo and ROS/MAP we
417 perform a gene set enrichment analysis against Gene Ontology pathways using the enrichR R package.
418 Only pathways with FDR < 0.05 are reported. The code we used to run the ROS/MAP DEG enrichments
419 are available here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/lineage.Rmd)
420 [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
421 DEG enrichments are available here [https://github.com/Sage-](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/lineageTCX.Rmd)
422 [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
423 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)
424 [Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/resilience.Rmd](https://github.com/Sage-Bionetworks/AMPAD_Lineage/blob/paper_rewrites_1/resilience.Rmd).

425

426 *Whole-genome sequencing*

427 Whole-genome sequencing was performed at the New York Genome Center for all individuals from the
428 ROS/MAP and Mayo cohorts. Detailed information for both data sets can be accessed via synapse
429 (DOI:10.7303/syn2580853). Briefly, 650ng of genomic DNA from whole blood was sheared using a
430 Covaris LE220 sonicator. DNA fragments underwent bead-based size selection and were subsequently
431 end-repaired, adenylated, and ligated to Illumina sequencing adapters. Libraries were sequenced on an
432 Illumina HiSeq X sequencer using 2 x 150bp cycles. Paired-end reads were aligned to the GRCh37
433 (hg19) human reference genome using the Burrows-Wheeler Aligner (BWA-MEM v0.7.8) and processed
434 using the GATK best-practices workflow^{53,54}. This included marking of duplicate reads by the use of
435 Picard tools v1.83, local realignment around indels, and base quality score recalibration (BQSR) via
436 Genome Analysis Toolkit (GATK v3.4.0). Joint variant calling files (vcfs) for whole-genome sequencing
437 data for the Mayo and ROS/MAP cohort were obtained through the AMP-AD knowledge portal
438 (www.synapse.org/#!/Synapse:syn10901595).

439

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

441 Likelihood ratio tests within a linear regression framework were used to model the relationship between
442 the quantitative expression trait pseudotime and genetic variants in 436 AD cases. Genome-wide genetic
443 association analysis was performed for 305 female patients in the ROS/MAP cohort and 131 female
444 patients in the Mayo cohort for which both genotyping and post-mortem RNA-seq data was available. An
445 efficient mixed model approach, implemented in the EMMAX software suite, was used to account for
446 potential biases and cryptic relatedness among individuals⁵⁵. Only variants with MAF > 0.05, genotyping
447 call rates > 95%, minimum sequencing depth of 20 reads and Hardy-Weinberg equilibrium $p > 10^{-4}$ were
448 considered for analysis. Quantile-quantile plots (**Figure S8-S9**) for the test statistics showed no
449 significant deviation between expected and observed p-values, highlighting that there is no consistent
450 differences across cases and controls except for the small number of significantly associated variants.

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

454

455 **References**

- 456 1. Cummings, J. L., Morstorf, T. & Zhong, K. Alzheimer's disease drug-development pipeline: few
457 candidates, frequent failures. *Alzheimers. Res. Ther.* **6**, 37 (2014).
- 458 2. Cummings, J. L., Doody, R. & Clark, C. Disease-modifying therapies for Alzheimer disease:
459 Challenges to early intervention. *Neurology* **69**, 1622–1634 (2007).
- 460 3. Ferreira, D. *et al.* Distinct subtypes of Alzheimer's disease based on patterns of brain atrophy:
461 longitudinal trajectories and clinical applications. *Sci. Rep.* **7**, 46263 (2017).
- 462 4. Bredesen, D. E. Metabolic profiling distinguishes three subtypes of Alzheimer's disease. *Aging*
463 (*Albany NY*) **7**, 595 (2015).
- 464 5. Brier, M. R. *et al.* Tau and A β imaging, CSF measures, and cognition in Alzheimer's disease.
465 *Sci. Transl. Med.* **8**, 338ra66--338ra66 (2016).
- 466 6. Gordon, B. A. *et al.* The relationship between cerebrospinal fluid markers of Alzheimer pathology
467 and positron emission tomography tau imaging. *Brain* **139**, 2249–2260 (2016).
- 468 7. Dichgans, M. *et al.* METACOHORTS for the study of vascular disease and its contribution to
469 cognitive decline and neurodegeneration: An initiative of the Joint Programme for
470 Neurodegenerative Disease Research. *Alzheimer's Dement.* **12**, 1235–1249 (2016).
- 471 8. Jack, C. R. *et al.* Tracking pathophysiological processes in Alzheimer's disease: an updated
472 hypothetical model of dynamic biomarkers. *Lancet Neurol.* **12**, 207–216 (2013).
- 473 9. Au, R., Piers, R. J. & Lancashire, L. Back to the future: Alzheimer's disease heterogeneity
474 revisited. *Alzheimer's Dement. (Amsterdam, Netherlands)* **1**, 368–370 (2015).
- 475 10. Carrasquillo, M. M. *et al.* Late-onset Alzheimer's risk variants in memory decline, incident mild
476 cognitive impairment, and Alzheimer's disease. *Neurobiol. Aging* **36**, 60–7 (2015).
- 477 11. Zhang, B. *et al.* Integrated systems approach identifies genetic nodes and networks in late-onset
478 Alzheimer's disease. *Cell* **153**, 707–20 (2013).
- 479 12. Mostafavi, S. *et al.* A molecular network of the aging human brain provides insights into the

- 480 pathology and cognitive decline of Alzheimer's disease. *Nat. Neurosci.* **21**, 811–819 (2018).
- 481 13. Conway, O. J. *et al.* ABI3 and PLCG2 missense variants as risk factors for neurodegenerative
482 diseases in Caucasians and African Americans. *Mol. Neurodegener.* **13**, 53 (2018).
- 483 14. Allen, M. *et al.* Conserved brain myelination networks are altered in Alzheimer's and other
484 neurodegenerative diseases. *Alzheimers. Dement.* **14**, 352–366 (2018).
- 485 15. Allen, M. *et al.* Divergent brain gene expression patterns associate with distinct cell-specific tau
486 neuropathology traits in progressive supranuclear palsy. *Acta Neuropathol.* **136**, 709–727 (2018).
- 487 16. Bengio, Y., Courville, A. & Vincent, P. Representation learning: A review and new perspectives.
488 *IEEE Trans. Pattern Anal. Mach. Intell.* **35**, 1798–1828 (2013).
- 489 17. Wolz, R., Aljabar, P., Hajnal, J. V. & Rueckert, D. Manifold Learning for Biomarker Discovery in
490 MR Imaging. in 116–123 (Springer, Berlin, Heidelberg, 2010). doi:10.1007/978-3-642-15948-
491 0_15
- 492 18. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
493 pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, (2014).
- 494 19. Qiu, X. *et al.* Reversed graph embedding resolves complex single-cell trajectories. *Nat. Methods*
495 **14**, 979–982 (2017).
- 496 20. Rosenberg, A. B. *et al.* Single-cell profiling of the developing mouse brain and spinal cord with
497 split-pool barcoding. *Science (80-.).* **360**, 176–182 (2018).
- 498 21. Mukherjee, S., Zhang, Y., Fan, J., Seelig, G. & Kannan, S. Scalable preprocessing for sparse
499 scRNA-seq data exploiting prior knowledge. *Bioinformatics* **34**, i124--i132 (2018).
- 500 22. Logsdon, B. *et al.* Meta-analysis of the human brain transcriptome identifies heterogeneity across
501 human AD coexpression modules robust to sample collection and methodological approach.
502 *bioRxiv* 510420 (2019).
- 503 23. Kunkle, B. W. *et al.* Genetic meta-analysis of diagnosed Alzheimer's disease identifies new risk
504 loci and implicates A β , tau, immunity and lipid processing. *Nat. Genet.* **51**, 414–430 (2019).
- 505 24. Allen, M. *et al.* Novel late-onset Alzheimer disease loci variants associate with brain gene

- 506 expression. *Neurology* **79**, 221–8 (2012).
- 507 25. Allen, M. *et al.* Late-onset Alzheimer disease risk variants mark brain regulatory loci. *Neurol.*
508 *Genet.* **1**, e15 (2015).
- 509 26. Allen, M. *et al.* Association of MAPT haplotypes with Alzheimer’s disease risk and MAPT brain
510 gene expression levels. *Alzheimers. Res. Ther.* **6**, 39 (2014).
- 511 27. Carrasquillo, M. M. *et al.* A candidate regulatory variant at the TREM gene cluster associates with
512 decreased Alzheimer’s disease risk and increased TREML1 and TREM2 brain gene expression.
513 *Alzheimers. Dement.* **13**, 663–673 (2017).
- 514 28. Sims, R. *et al.* Rare coding variants in PLCG2, ABI3, and TREM2 implicate microglial-mediated
515 innate immunity in Alzheimer’s disease. *Nat. Genet.* **49**, 1373–1384 (2017).
- 516 29. Braak, H., Alafuzoff, I., Arzberger, T., Kretschmar, H. & Del Tredici, K. Staging of Alzheimer
517 disease-associated neurofibrillary pathology using paraffin sections and immunocytochemistry.
518 *Acta Neuropathol.* **112**, 389–404 (2006).
- 519 30. Mirra, S. S. *et al.* The Consortium to Establish a Registry for Alzheimer’s Disease (CERAD). Part
520 II. Standardization of the neuropathologic assessment of Alzheimer’s disease. *Neurology* **41**, 479–
521 86 (1991).
- 522 31. De Strooper, B. & Karran, E. The Cellular Phase of Alzheimer’s Disease. *Cell* **164**, 603–615
523 (2016).
- 524 32. Zhang, Y. *et al.* An RNA-sequencing transcriptome and splicing database of glia, neurons, and
525 vascular cells of the cerebral cortex. *J. Neurosci.* **34**, 11929–47 (2014).
- 526 33. Chibnik, L. B. *et al.* Susceptibility to neurofibrillary tangles: role of the PTPRD locus and limited
527 pleiotropy with other neuropathologies. *Mol. Psychiatry* **23**, 1521 (2018).
- 528 34. Lambert, J.-C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for
529 Alzheimer’s disease. *Nat. Genet.* **45**, 1452 (2013).
- 530 35. Janelidze, S. *et al.* CSF biomarkers of neuroinflammation and cerebrovascular dysfunction in early
531 Alzheimer disease. *Neurology* **91**, e867–e877 (2018).

- 532 36. Whitwell, J. L. *et al.* Neuroimaging correlates of pathologically defined subtypes of Alzheimer's
533 disease: a case-control study. *Lancet Neurol.* **11**, 868–877 (2012).
- 534 37. Chen, E. Y. *et al.* Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool.
535 *BMC Bioinformatics* **14**, 128 (2013).
- 536 38. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology
537 Consortium. *Nat. Genet.* **25**, 25–9 (2000).
- 538 39. Dong, A. *et al.* Heterogeneity of neuroanatomical patterns in prodromal Alzheimer's disease: links
539 to cognition, progression and biomarkers. *Brain* **140**, aww319 (2016).
- 540 40. Wang, X. *et al.* Genetic determinants of disease progression in Alzheimer's disease. *J. Alzheimers.*
541 *Dis.* **43**, 649–55 (2015).
- 542 41. Allen, M. *et al.* Human whole genome genotype and transcriptome data for Alzheimer's and other
543 neurodegenerative diseases. *Sci. Data* **3**, (2016).
- 544 42. De Jager, P. L. *et al.* A multi-omic atlas of the human frontal cortex for aging and Alzheimer's
545 disease research. *Sci. Data* **5**, 180142 (2018).
- 546 43. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 547 44. Anders, S., Pyl, P. T. & Huber, W. HTSeq—a Python framework to work with high-throughput
548 sequencing data. *Bioinformatics* **31**, 166–169 (2015).
- 549 45. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: Precision weights unlock linear model
550 analysis tools for RNA-seq read counts. *Genome Biol.* **15**, (2014).
- 551 46. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and
552 microarray studies. *Nucleic Acids Res.* **43**, e47–e47 (2015).
- 553 47. Bennett, D. *et al.* Overview and findings from the rush Memory and Aging Project. *Curr.*
554 *Alzheimer Res.* **9**, 646–663 (2012).
- 555 48. Bennett, D. A. *et al.* Religious orders study and rush memory and aging project. *J. Alzheimers's*
556 *Dis.* 1–28 (2018).
- 557 49. Allen, M. *et al.* Human whole genome genotype and transcriptome data for Alzheimer's and other

- 558 neurodegenerative diseases. *Sci. data* **3**, 160089 (2016).
- 559 50. Ferretti, M. T. *et al.* Sex differences in Alzheimer disease—the gateway to precision medicine.
560 *Nat. Rev. Neurol.* **1** (2018).
- 561 51. Deming, Y. *et al.* Sex-specific genetic predictors of Alzheimer’s disease biomarkers. *Acta*
562 *Neuropathol.* **136**, 857–872 (2018).
- 563 52. Mao, Q., Wang, L., Goodison, S. & Sun, Y. Dimensionality Reduction Via Graph Structure
564 Learning. in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge*
565 *Discovery and Data Mining - KDD '15* 765–774 (ACM Press, 2015).
566 doi:10.1145/2783258.2783309
- 567 53. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
568 generation DNA sequencing data. *Genome Res.* **20**, 1297–303 (2010).
- 569 54. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
570 *Bioinformatics* **25**, 1754–60 (2009).
- 571 55. Kang, H. M. *et al.* Variance component model to account for sample structure in genome-wide
572 association studies. *Nat. Genet.* **42**, 348–354 (2010).
- 573

574 **Acknowledgements**

575 This work was supported by NIA grants U54AG054345 and RFIAG057443. The ROSMAP Study data
576 were provided by the Rush Alzheimer's Disease Center, Rush University Medical Center, Chicago. Data
577 collection was supported through funding by NIA grants P30AG10161, R01AG15819, R01AG17917,
578 R01AG30146, R01AG36836, U01AG32984, U01AG46152, the Illinois Department of Public Health,
579 and the Translational Genomics Research Institute. Mayo RNAseq Study data were provided by the
580 following sources: The Mayo Clinic Alzheimer's Disease Genetic Studies, led by Dr. Nilufer Ertekin-
581 Taner and Dr. Steven G. Younkin, Mayo Clinic, Jacksonville, FL using samples from the Mayo Clinic
582 Study of Aging, the Mayo Clinic Alzheimer's Disease Research Center, and the Mayo Clinic Brain Bank.
583 Data collection was supported through funding by NIA grants P50 AG016574, R01 AG032990, U01
584 AG046139, R01 AG018023, U01 AG006576, U01 AG006786, R01 AG025711, R01 AG017216, R01
585 AG003949, NINDS grant R01 NS080820, CurePSP Foundation, and support from Mayo Foundation.
586 Study data includes samples collected through the Sun Health Research Institute Brain and Body
587 Donation Program of Sun City, Arizona. The Brain and Body Donation Program is supported by the
588 National Institute of Neurological Disorders and Stroke (U24 NS072026 National Brain and Tissue
589 Resource for Parkinson's Disease and Related Disorders), the National Institute on Aging (P30 AG19610
590 Arizona Alzheimer's Disease Core Center), the Arizona Department of Health Services (contract 211002,
591 Arizona Alzheimer's Research Center), the Arizona Biomedical Research Commission (contracts 4001,
592 0011, 05-901 and 1001 to the Arizona Parkinson's Disease Consortium) and the Michael J. Fox
593 Foundation for Parkinson's Research. MSBB data were generated from postmortem brain tissue collected
594 through the Mount Sinai VA Medical Center Brain Bank and were provided by Dr. Eric Schadt from
595 Mount Sinai School of Medicine.

596

597 **Author Contributions**

598 S.M. and B.A.L. designed the study. S.M. and B.A.L. performed the analyses. S.M., C.P., S.J., G.G.,

599 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.

600

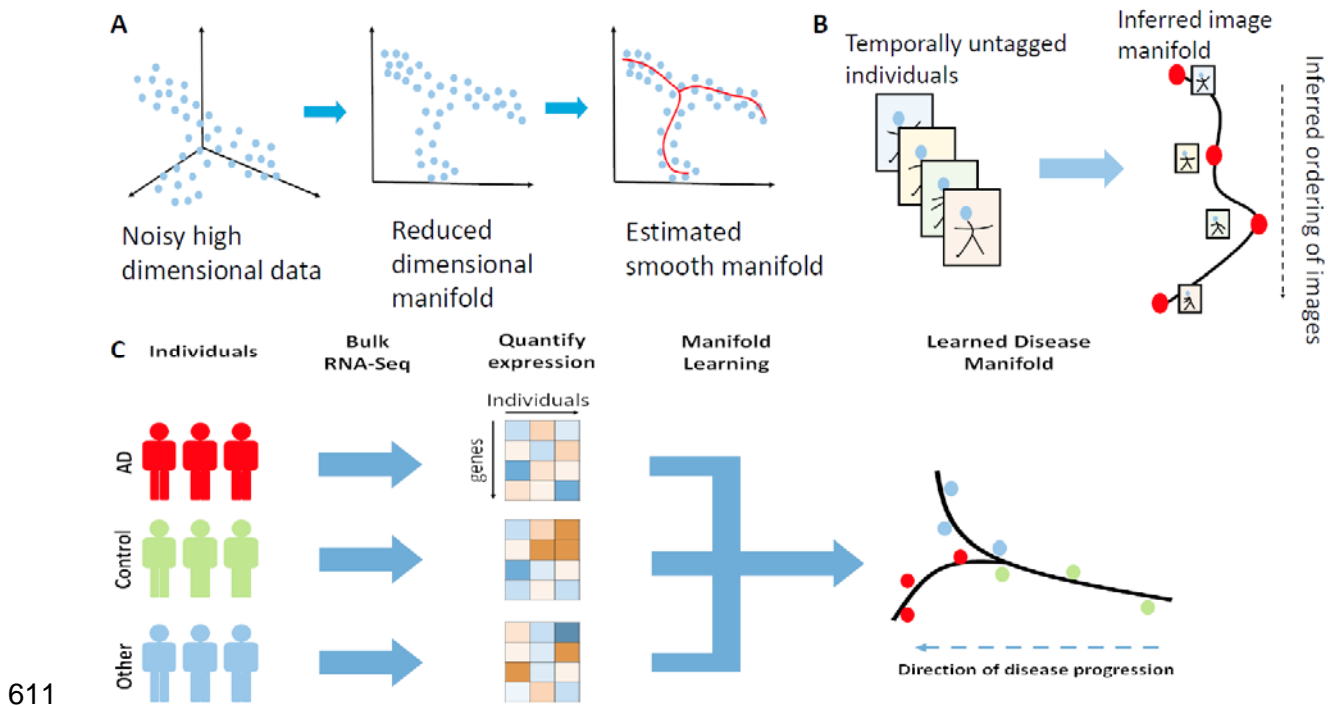
601 **Tables**

602 **Table 1** - Representative significant Gene Ontology pathway enrichments (FDR<0.05) of differentially
 603 expressed genes for each branch (FDR < 0.05). Differentially expressed genes are identified with an
 604 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)		
		2			
	Up	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)			

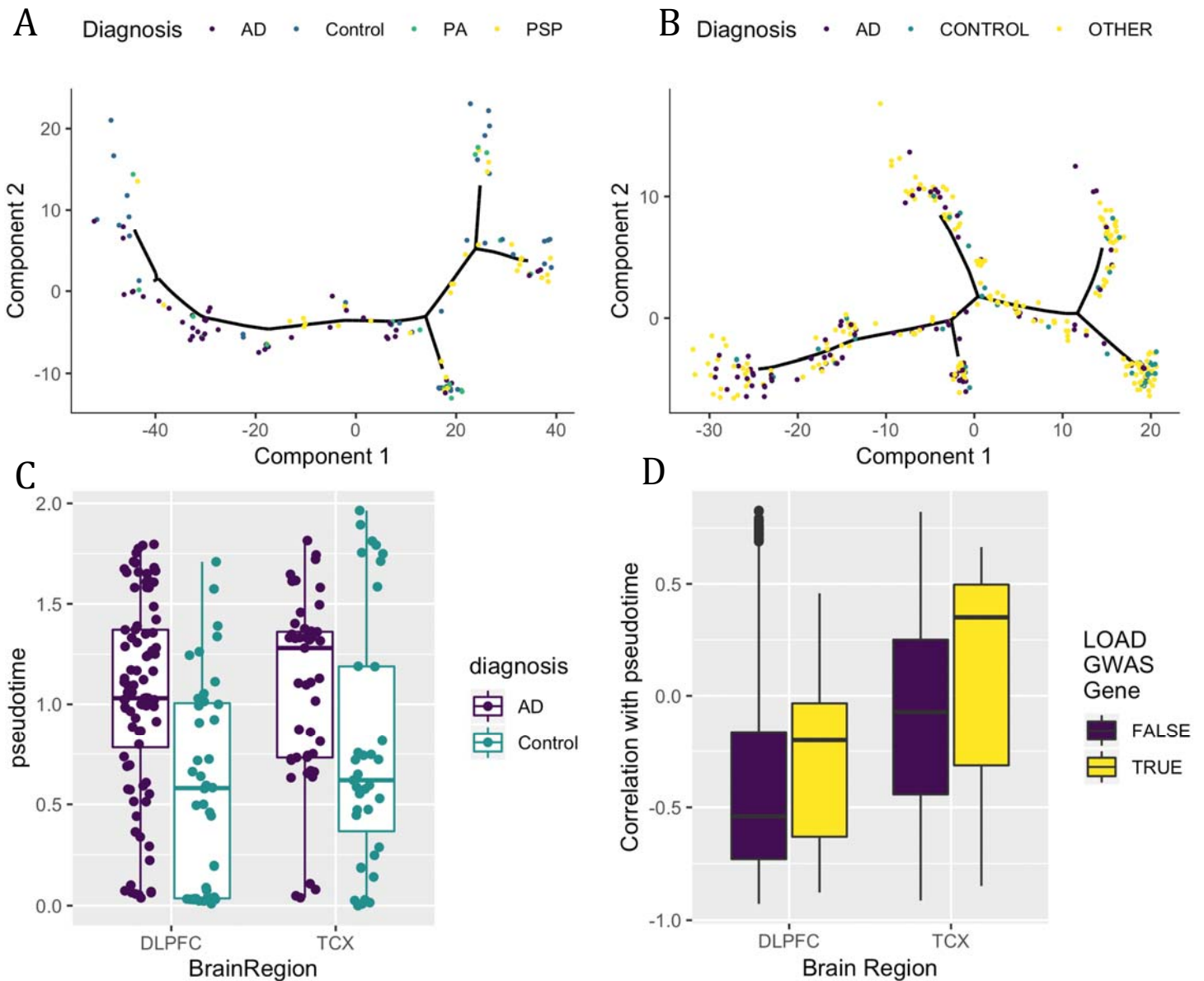
605 **Figures**

606 **Figure 1** - Overview of manifold learning for unraveling staging in Alzheimer's disease. A) Illustration of
607 steps in manifold learning. B) A common application of manifold learning used in computer vision to
608 order temporally untagged images into sequences. C) Illustration of lineage inference process for LOAD.
609 RNA-seq samples with different disease diagnoses were pooled, batch normalized, and a smooth
610 manifold was learned for each brain region across individuals (each point is an individual).

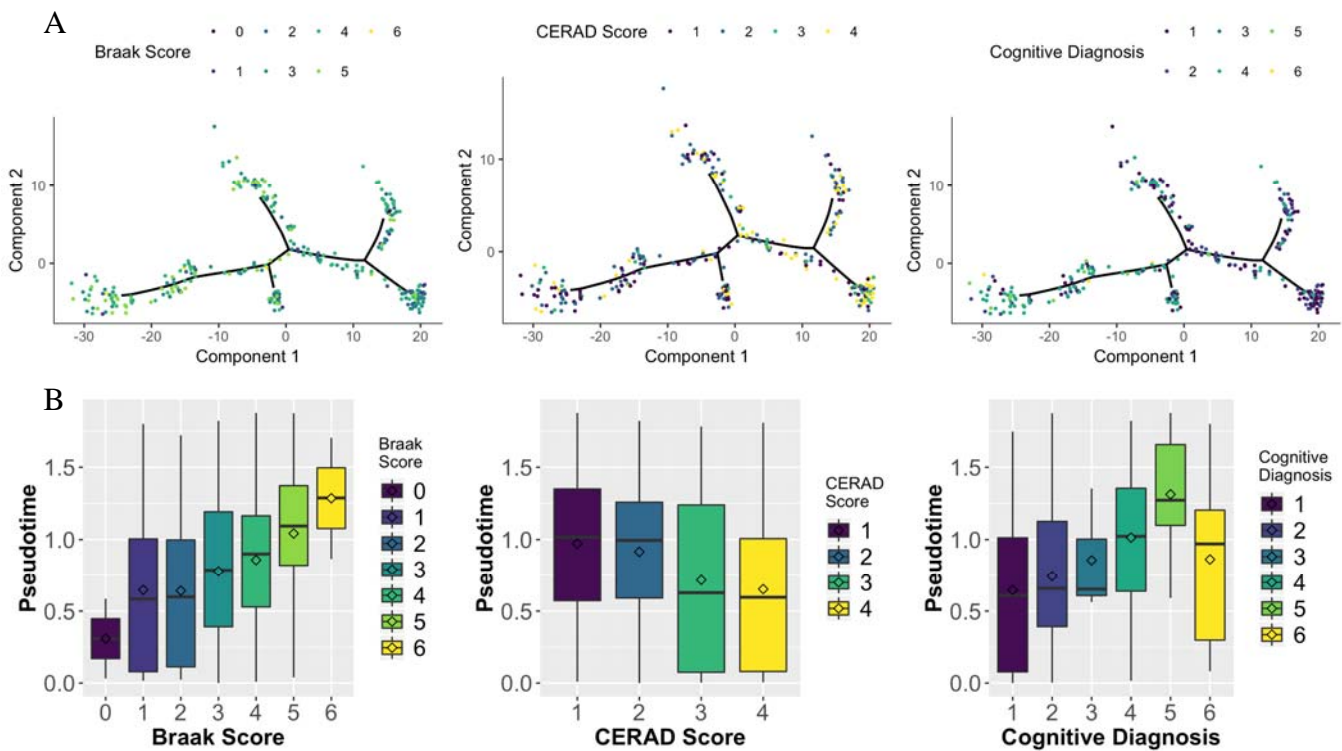


611

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



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



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

