

1 **Deciphering cellular transcriptional alterations in Alzheimer's disease brains**

2

3 Xue Wang^{1*}, Mariet Allen², Shaoyu Li³, Zachary S. Quicksall¹, Tulsi A. Patel², Troy P.

4 Carnwath², Joseph S. Reddy¹, Minerva M. Carrasquillo², Sarah J. Lincoln², Thuy T. Nguyen²,

5 Kimberly G. Malphrus², Dennis W. Dickson², Julia E. Crook¹, Yan W. Asmann¹, Nilüfer

6 Ertekin-Taner^{2,4*}

7

8

9

10

11 ¹ Department of Health Sciences Research, Mayo Clinic Florida, Jacksonville, FL, USA.

12 ² Department of Neuroscience, Mayo Clinic Florida, Jacksonville, FL, USA.

13 ³ Department of Mathematics and Statistics, University of North Carolina at Charlotte, Charlotte,
14 NC, USA.

15 ⁴ Department of Neurology, Mayo Clinic Florida, Jacksonville, FL, USA.

16

17 * Corresponding Authors

18 E-mails: Wang.Xue@mayo.edu (XW) and Taner.Nilufer@mayo.edu (NET)

19

20

21

22

23 **ABSTRACT**

24 Large-scale brain bulk-RNAseq studies identified molecular pathways implicated in Alzheimer's
25 disease (AD), however these findings can be confounded by cellular composition changes in
26 bulk-tissue. To identify cell intrinsic gene expression alterations of individual cell types, we
27 designed a bioinformatics pipeline and analyzed three AD and control bulk-RNAseq datasets of
28 temporal and dorsolateral prefrontal cortex from 685 brain samples. We detected cell-proportion
29 changes in AD brains that are robustly replicable across the three independently assessed
30 cohorts. We applied three different algorithms including our in-house algorithm to identify cell
31 intrinsic differentially expressed genes in individual cell types (CI-DEGs). We assessed the
32 performance of all algorithms by comparison to single nucleus RNAseq data. We identified
33 consensus CI-DEGs that are common to multiple brain regions. Despite significant overlap
34 between consensus CI-DEGs and bulk-DEGs, many CI-DEGs were absent from bulk-DEGs.
35 Consensus CI-DEGs and their enriched GO terms include genes and pathways previously
36 implicated in AD or neurodegeneration, as well as novel ones. We demonstrated that the
37 detection of CI-DEGs through computational deconvolution methods is promising and highlight
38 remaining challenges. These findings provide novel insights into cell-intrinsic transcriptional
39 changes of individual cell types in AD and may refine discovery and modeling of molecular
40 targets that drive this complex disease.

41

42 **Introduction**

43 Alzheimer's disease (AD) is a neurodegenerative disease that affects ~5.7 million
44 patients with annual cost of more than \$230 billion in the US¹. Effective disease-modifying
45 drugs are still elusive despite the urgent need and increasing global burden^{2,3}. Pathologically, AD
46 is marked by amyloid-beta plaques and neurofibrillary tangles, along with neuronal loss and
47 gliosis in the affected brain regions. Transcriptome-wide expression levels have been analyzed
48 from bulk brain tissue of hundreds of AD patients and neuropathologically normal controls⁴⁻⁸ to
49 discover genes and biological pathways that are perturbed in and/or lead to AD. Systems biology
50 and bioinformatics analysis of these data have implicated altered pathways in AD including
51 immune response⁶ and myelin metabolism^{4,5}. However, a fundamental knowledge gap remains
52 concerning whether disease-associated changes in brain gene expression are due to changes in
53 cellular composition of the AD samples secondary to disease neuropathology, or due to changes
54 in the intrinsic regulation/activity of genes in the central nervous system (CNS) cells. From a
55 clinical perspective, it is difficult to target changes in cellular composition secondary to
56 neuropathology, whereas intrinsic changes in gene expression that may drive AD progression are
57 potentially "druggable".

58 We expect that identifying cell-intrinsic differentially expressed genes (CI-DEGs) of
59 individual CNS cell types will reveal novel insights into the genes and pathways that could
60 potentially identify drug targets and lead to AD therapeutics. This approach circumvents the
61 influence of cell-composition changes that can impact disease associated DEGs obtained from
62 bulk tissue transcriptome analysis. RNA sequencing (RNAseq) studies from single cell, single
63 nucleus or purified adult human CNS cells⁹⁻¹¹ can be used to identify CI-DEGs. Even though
64 such single cell-type RNAseq data can effectively serve as a reference to annotate bulk-tissue

65 transcriptome data⁴, such approaches remain costly, cumbersome and limited in sample sizes. On
66 the other hand, there exist large-scale bulk brain RNAseq datasets^{5,8,12}, which can be leveraged to
67 identify CI-DEGs through analytic deconvolution of bulk tissue expression into signals of
68 individual cell types by using cell proportions or their proxies^{13,14}.

69 In this study, we describe the design and application of a bioinformatics pipeline that uses
70 cell type marker genes to estimate cell proportion^{15,16} to deconvolute bulk tissue transcriptome
71 data using three computational approaches^{13,14,17} and to subsequently identify CI-DEGs. We
72 applied our pipeline to the analysis of three post-mortem brain datasets, one from dorsolateral
73 prefrontal cortex (DLPFC)⁸ and two from temporal cortex (TCX)^{4,12,18} regions, comprised of 685
74 unique samples. Consensus CI-DEGs common to both TCX and DLPFC regions were analyzed
75 for enrichment of gene ontology (GO) terms. We compared the results of consensus CI-DEGs to
76 consensus bulk-DEGs. In addition, for the DLPFC⁸ dataset that had both bulk and single nucleus
77 RNAseq¹⁹ (snRNAseq) data, we compared the CI-DEGs from the computational deconvolution
78 to CI-DEGs obtained from snRNAseq¹⁹.

79 To our knowledge, this is the first study to detect consensus CI-DEGs and their enriched
80 gene ontology (GO) terms from multiple brain regions using multiple computational
81 deconvolution algorithms for AD and control RNAseq samples. Our study illustrates the cell
82 proportion landscape of AD and control brain regions assessed in three independent RNAseq
83 studies^{4,7,8,12}. We identify consensus CI-DEGs many of which are not observed in bulk-DEG
84 analysis and characterize their cell-type specificity. GO terms that are enriched for CI-DEGs
85 implicate cell intrinsic transcriptional alterations that may influence AD, rather than be a result
86 of cell-proportion changes in this disease. These CI-DEGs and their biological pathways may
87 serve as refined molecular targets for therapeutic discoveries and disease modeling in AD. Our

88 study also demonstrates that detection of CI-DEGs through computational deconvolution
89 methods is promising while some challenges remain.

90

91 **Results**

92 **Cellular composition in three brain cohorts from two brain regions**

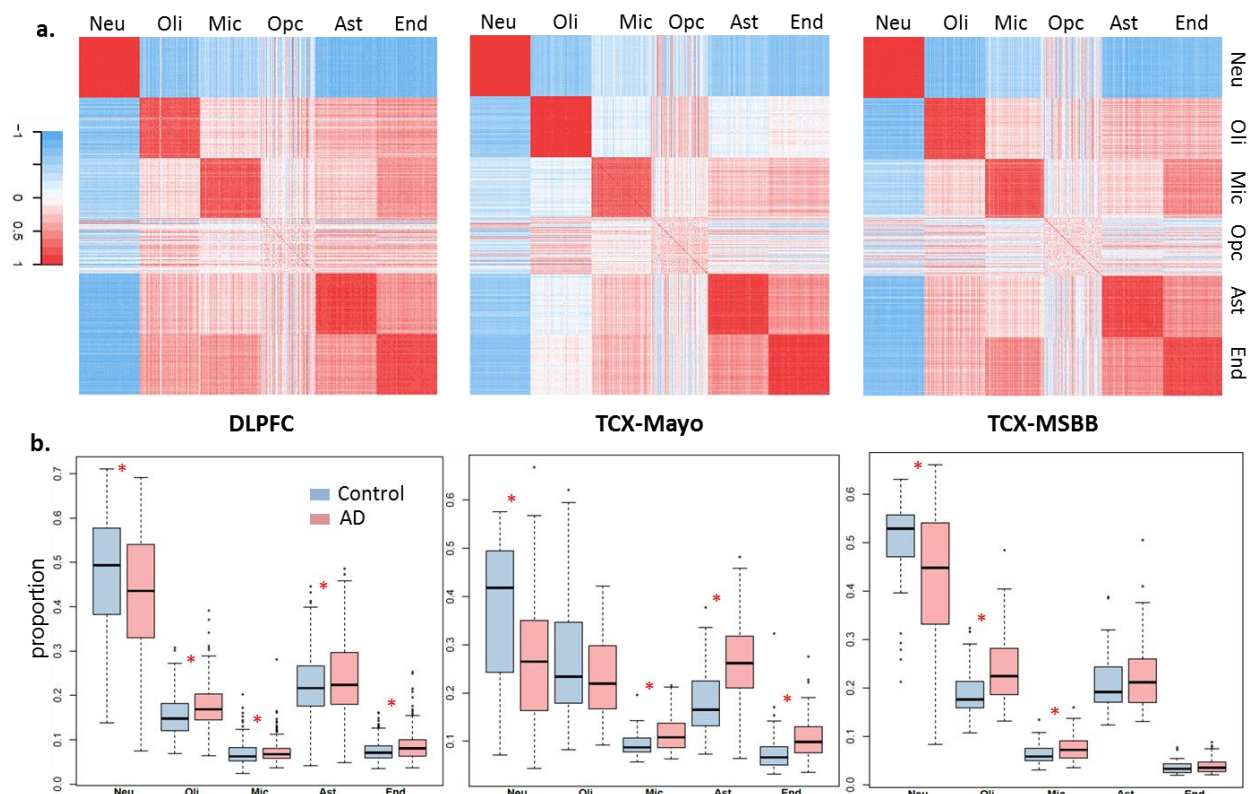
93 We analyzed three cohorts each consisting of post-mortem brains from AD and control
94 subjects (**Table S1**), namely the Rush Religious Orders Study and Memory and Aging Project
95 dorsolateral prefrontal cortex (DLPFC)^{7,8}, Mayo Clinic temporal cortex (TCX-Mayo)^{4,12}, and
96 Mount Sinai VA Medical Center Brain Bank temporal cortex (TCX-MSBB)¹⁸. We generated the
97 TCX-Mayo RNAseq dataset, and downloaded DLPFC and TCX-MSBB RNAseq datasets from
98 the AMP-AD knowledge portal on Synapse (www.synapse.org).

99 Cell proportions (**Table S2**) were estimated for DLPFC, TCX-Mayo and TCX-MSBB
100 datasets independently using the digital sorting algorithm (DSA) method¹⁶ and the top 100
101 marker genes (**Table S3**) obtained from R package BREGITEA¹⁵ for each of the following cell
102 types – neuron, oligodendrocyte, microglia, oligodendrocyte progenitor cell (OPC), astrocyte and
103 endothelial cell.

104 An inspection about the pairwise correlation between marker genes (**Fig 1a**) revealed that
105 markers of OPC have poor median pairwise Pearson correlation values of 0.12 in DLPFC, 0.11
106 in TCX-Mayo and 0.06 in TCX-MSBB respectively, whereas among the other five cell types
107 neuronal markers have the highest median correlation (0.68 in DLPFC, 0.78 in TCX-Mayo and
108 0.67 in TCX-MSBB), and microglia markers have the lowest correlation (0.37 in DLPFC, 0.42
109 in TCX-Mayo and 0.44 in TCX-MSBB). In addition, a computer simulation study (**Fig S1**)
110 demonstrated that the estimated proportions of OPC were not robust upon using different

111 selection of marker genes. Therefore, we did not include OPC in downstream analyses in this
112 study.

113 In all three datasets, neuronal cell proportion estimates were significantly lower in AD
114 compared to controls (**Fig 1b**). The magnitude of this decrease was the greatest for TCX-Mayo
115 (AD mean proportion = 28.0%, Control = 35.7%; ratio of AD:control cell proportions=0.78),
116 followed by TCX-MSSM (AD = 42.3%, control = 49.3%; ratio=0.87) and DLPFC (AD = 42.4%,
117 control = 47.4%; ratio=0.89). The estimated proportions of microglia were significantly higher in
118 AD vs. controls for all datasets, with higher magnitude in TCX-Mayo (AD:control ratio=1.19)
119 and TCX-MSBB (AD:control ratio=1.19) than for DLPFC (AD:control ratio=1.06). The
120 estimated proportions of astrocytes and endothelial cells were significantly higher in AD vs.
121 controls for DLPFC and TCX-Mayo datasets, although the magnitude was greater in TCX-Mayo
122 (1.40 and 1.30 respectively) than in DLPFC (1.07 and 1.14 respectively) for both cell types.
123 Oligodendrocyte proportion is significantly higher in AD in DLPFC with AD:control ratio 1.14
124 and TCX-MSBB with AD:control ratio 1.27, although remains unchanged in TCX-Mayo with
125 the ratio 0.94. Collectively, these findings demonstrate that the proportions of CNS cell types are
126 different in post-mortem AD vs. control brains for most cell types. Although these proportional
127 changes with AD are mostly consistent across the different studies, their extent varies across
128 brain regions, with TCX tending towards higher magnitude of neuronal loss and microglia
129 proliferation than DLPFC.



130

131

132 **Fig.1: a)** Pearson correlation between marker gene expressions in six cell types. Marker genes
133 are from literature. **b)** Estimated cell proportions in DLPFC, TCX-Mayo and TCX-MSBB
134 datasets in five cell types. Red asterisk indicates differences between cell proportions in AD and
135 control groups at nominal p value 0.05 from Wilcoxon rank sum test.

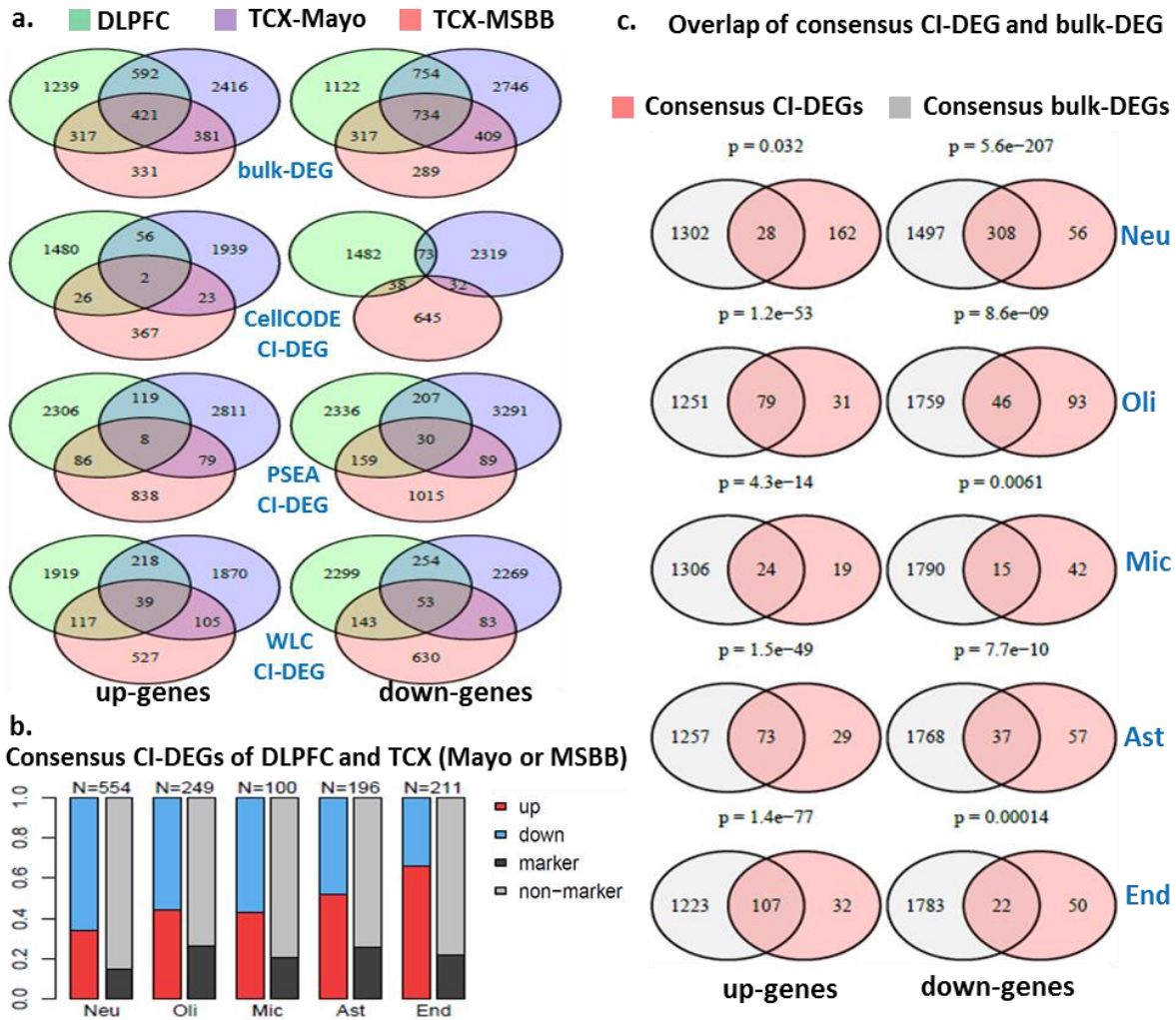
136

137 Differential Expression Analyses

138 In this study, three computational approaches were applied to identify cell intrinsic
139 differential expression in individual cell types (CI-DEGs, **Table S4-S6**), namely CellCODE¹⁴,
140 PSEA¹³ and our method WLC. Differentially expressed genes from bulk brain tissue (bulk-
141 DEGs) were identified through linear regression without adjusting for cellular composition
142 (**Table S7**). For the DLPFC, TCX-Mayo and TCX-MSBB datasets, we obtained bulk-DEGs and

143 CI-DEGs from the three computer algorithms for neuronal, oligodendrocytic, microglial,
144 astrocytic and endothelial cell types respectively.

145 We compared bulk-DEGs across the three datasets (**Fig 2a**, top panel). Similarly, CI-
146 DEGs from CellCODE, PSEA and WLC are compared across datasets (**Fig 2a**, lower panels),
147 such that CI-DEGs shared between datasets are required to be consistent in the designated cell
148 type. All DEGs are identified at nominal p-value cutoff 0.05 and shared CI-DEGs have the same
149 direction of change in the compared datasets. The ratio of overlap between any two datasets over
150 all DEGs, i.e. the number in overlapping areas of the Venn diagram over the total number (**Fig**
151 **2a**, top panel), is 30.0% or 1711/5697 in up-regulated bulk-DEGs, and 34.8% or 2214/6371 in
152 down-regulated bulk-DEGs. This ratio of overlap in bulk-DEGs is much higher than that in CI-
153 DEGs (2.7%, 4.7% and 10.0% in up-regulated genes from CellCODE, PSEA and WLC
154 respectively; 3.1%, 6.8% and 9.3% in down-regulated genes from CellCODE, PSEA and WLC
155 respectively).



156

157

158 **Fig.2: a)** Overlap across three independent RNAseq datasets of bulk-DEGs (upper panel) and CI-
 159 DEGs (lower panels) from three computational approaches. **b)** Consensus CI-DEGs between
 160 DLPFC and TCX brain regions, which consist of consensus CI-DEGs between DLPFC and
 161 TCX-Mayo, or between DLPFC and TCX-MSBB. **c)** Overlap between consensus CI-DEGs and
 162 consensus bulk-DEGs, per cell type. The p-values of overlap are from hypergeometric tests.

163

164 **Consensus CI-DEGs between DLPFC and TCX**

165 To obtain the consensus CI-DEGs that are shared between DLPFC and TCX brain
166 regions, we selected those CI-DEGs that are detected in “DLPFC and TCX-Mayo” or in
167 “DLPFC and TCX-MSBB” under any of the three algorithms (**Fig S2**). We combined all such
168 genes, which collectively comprised the consensus CI-DEGs for each cell type (**Fig 2b**).
169 Similarly, consensus bulk-DEGs were the combined set of bulk-DEGs shared between “DLPFC
170 and TCX-Mayo” or “DLPFC and TCX-MSBB”.

171 Most consensus CI-DEGs are from neuronal cells (N=554), followed by
172 oligodendrocytes (N=249), whereas microglia contributed the least number (N=100). The
173 majority (65.7% or 364/554) of neuronal CI-DEGs is down-regulated in AD, and the majority
174 (65.9% or 139/211) of endothelial CI-DEGs is up-regulated in AD, with other cell types lying in
175 between. Some of these CI-DEGs are also among the 1000 marker genes of the corresponding
176 cell type from BRETIGEA¹⁵; 14.8% or 82/554 of neuronal CI-DEGs are also neuronal markers,
177 26.5% or 66/249 of oligodendrocyte CI-DEGs are also oligodendrocyte markers, and other cell
178 types lie in between.

179 With regards to consensus bulk-DEGs (**Fig S3**), 28.2% of them (885/3135) are cell type
180 markers; 10.4% neuronal markers, 5.6% oligodendrocyte, 3.4% microglia, 4.8% astrocyte and
181 4.0% endothelial markers. The above observations indicate that computational deconvolution
182 algorithms could identify CI-DEGs for both marker genes and non-marker genes. Importantly,
183 the proportion of non-marker CI-DEGs is greater than that in bulk-DEGs. This suggests that
184 compared to bulk-DEGs, CI-DEGs may be capturing a greater proportion of expression changes
185 that are not due to mere cell population changes.

186 We also compared the consensus bulk-DEGs with consensus CI-DEGs (**Fig S4**). We
187 determined that only a small fraction (14.7% or 28/190) of the up-regulated neuronal CI-DEGs

188 was also present in up bulk-DEGs although the overlap is still significant (**Fig 2c**). In
189 comparison, most of the up-regulated CI-DEGs of the other four cell types were included in up
190 bulk-DEGs. On the other hand, most (85.0% or 308/365) of the down-regulated neuronal CI-
191 DEGs were also present in down bulk-DEGs, whereas most of the down-regulated CI-DEGs of
192 the other four cell types were absent from this group. Since bulk-DEGs did not adjust for
193 neuronal loss and gliosis in AD (**Fig 1b**), its ability to identify up-regulated neuronal genes and
194 down-regulated glial genes is likely to be compromised. For the same reason, bulk-DEGs may
195 have a false inflation of detecting down-regulated neuronal and up-regulated glial genes.

196

197 **Enriched GO terms of consensus CI-DEGs between DLPCF and TCX**

198 To identify pathways implicated by CI-DEGs that are robust across brain regions, we
199 performed Gene Ontology (GO) enrichment analysis^{20,21} for the consensus CI-DEGs, assessing
200 separately those that are up vs. down in AD subjects (**Table S8-S17**). **Fig 3** illustrates the top
201 two enriched GO terms by enrichment p-values, after filtering out terms that encompass less than
202 four CI-DEGs or are cellular compartments.

203 Consensus CI-DEGs revealed biological pathways that are perturbed in AD in specific
204 brain cell types. Some of these pathways have previously been implicated in AD and others are
205 novel. Down-regulated neuronal CI-DEGs were enriched in neuropeptide hormone activity
206 (GO:0005184) and hormone activity (GO:0005179) pathways, which include *VGF* (a.k.a.
207 neuroendocrine regulatory peptide 1)²² and corticotropin releasing hormone (CRH)²³ (**Table**
208 **S13**). Consensus up-regulated neuronal CI-DEGs were significantly enriched in potassium
209 channel activity (GO:0005267) and regulation of ion transport (GO:0043269) pathways (**Table**
210 **S8**). The latter GO term encompasses most of the genes from the former, and also includes other

211 genes involved in neuronal functions such as the glutamate ionotropic receptor NMDA type
212 subunit 1, *GRINI*²⁴ and *SYTI*, which encodes the synaptic vesicle protein, synaptotagmin²⁵.

213 Many of the most significant GO terms are related to key functions of the respective cell
214 types for the glial CI-DEGs, as well. The top enriched pathway of down-regulated CI-DEGs in
215 oligodendrocytes is myelination (GO:0042552), including myelin basic protein (*MBP*)⁴,
216 plasmolipin (*PLLP*)^{4,5}, myelin and lymphocyte protein (*MAL*), and myelin-associated
217 glycoprotein (*MAG*)²⁶ (**Table S14**). Up-regulated CI-DEGs of oligodendrocytes are enriched in
218 ceramide biosynthetic process (GO:0046513) including ceramide synthase 4 (*CERS4*) and UDP
219 glycosyltransferase 8 (*UGT8*)⁵ (**Table S9**). Ceramide is a constituent of sphingomyelin, a
220 sphingolipid which is particularly found in the myelin sheath; and also a multi-functional
221 signaling molecule^{27,28}. Hence, both the down-regulated and the up-regulated oligodendroglial
222 consensus CI-DEGs highlight different components of the myelin biology that are perturbed in
223 AD.

224 Similarly, microglial, astrocytic and endothelial CI-DEGs also highlight processes
225 pertinent to the functions of these cell types. Microglial up-regulated CI-DEGs are enriched in
226 inflammatory response (GO:000695) and leukocyte activation (GO:0002696), which includes
227 complement C3a receptor 1 (*C3ARI*)²⁹, interleukin 18 (*IL18*)^{30,31} and CCAAT enhancer binding
228 protein alpha (*CEBPA*)³² genes (**Table S10**).

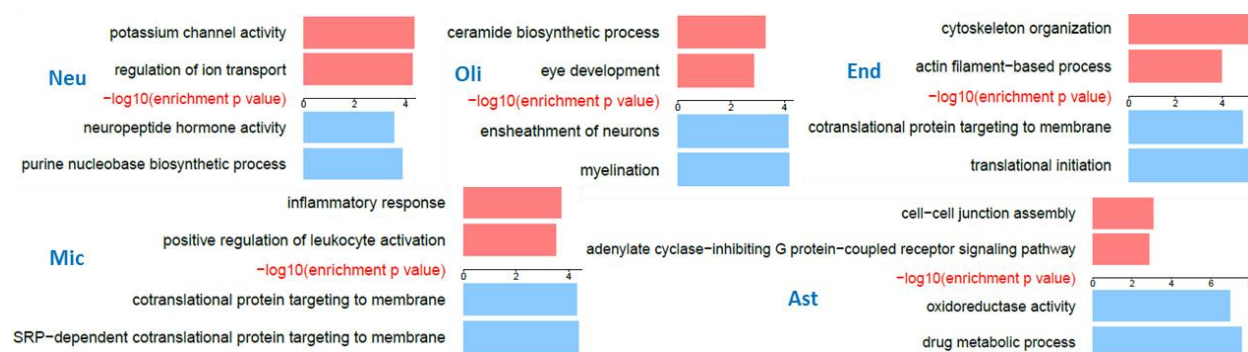
229 Astrocytes, a cell type that plays a critical role in maintaining brain energy dynamics³³
230 and metabolism³⁴, show enrichment of oxidoreductase activity (GO:0016491) and drug
231 metabolic process (GO:0017144) in down-regulated CI-DEGs which includes genes glutathione
232 S-transferase mu 2 (*GSTM2*)³⁵ and thioredoxin2 (*TXN2*)³⁶ (**Table S16**). Astrocytic up-regulated
233 consensus CI-DEGs are enriched for cell-cell junction assembly (GO:0007043) process (**Table**

234 **S11**), including the astrocytic gap junction protein connexin43 (*GJA1*)³⁷, which was identified as
235 a key regulator associated with AD related outcomes. The other top GO process for astrocytic
236 up-regulated consensus CI-DEGs is adenylate cyclase-inhibiting G protein-coupled receptor
237 signaling pathway (GO:0007193), which harbors adenylate cyclase 8 (*ADCY8*), involved in
238 memory functions³⁸.

239 Finally, endothelial cells, which are crucial in maintaining blood-brain barrier
240 integrity^{39,40}, show enrichment of up-regulated DEGs in cytoskeleton organization
241 (GO:0007010) and actin filament-based process (GO:0030029) (**Table S12**).

242 Importantly, some CI-DEGs highlight protein translation as a top perturbed biological
243 pathway. Down-regulated microglial consensus CI-DEGs show enrichment in processes
244 involved in protein translation (GO:0006614 and GO:0006613), which include ribosomal protein
245 encoding genes⁴¹⁻⁴³ (**Table S15**). Similarly, down-regulated endothelial consensus CI-DEGs also
246 harbor ribosomal protein encoding genes, with enrichment in protein translation related GO
247 processes (GO:0006413 and GO:0006613) (**Table S17**).

248



249

250

251 **Fig 3:** Top two enriched GO terms in up (red) or down-regulated (blue) consensus CI-DEGs
252 between DLPFC and TCX regions, per cell type.

253

254 **Comparison of CI-DEGs from computational deconvolution vs. snRNAseq**

255 We determined the extent to which each of the three computational deconvolution
256 algorithms could detect CI-DEGs from bulk tissue by comparison of their results with those
257 obtained in a published snRNAseq study¹⁹. The ROSMAP dataset utilized in our study has both
258 bulk RNAseq from DLPFC (bulk-DLPFC) as well as snRNAseq (snDLPFC) in a subset of its
259 participants¹⁹. We compared the bulk-DLPFC data deconvoluted with three different algorithms
260 with the published snDLPFC¹⁹ data. Endothelial CI-DEGs were not available from the
261 snRNAseq study, therefore overlap of results could be assessed only for four cell types.

262 We tested the overlap between the top CI-DEGs for each cell type obtained from
263 deconvoluted bulk-DLPFC and those from snDLPFC ranked by their p values (**Fig 4a**). We
264 evaluated the overlap for a range of top CI-DEGs up to top 1,000 genes. Overlap for CI-DEGs
265 that are either up (**Fig 4a, upper panel**) or down (**Fig 4a, lower panel**) in AD were assessed
266 separately. Hence, overlapping genes had both similar ranks and direction of effect in both
267 deconvoluted bulk-DLPFC and snDLPFC analyses. We established the significance of overlap
268 using simulations for a range of top ranked genes (N=200, 600 and 1,000) (**Table S18**).

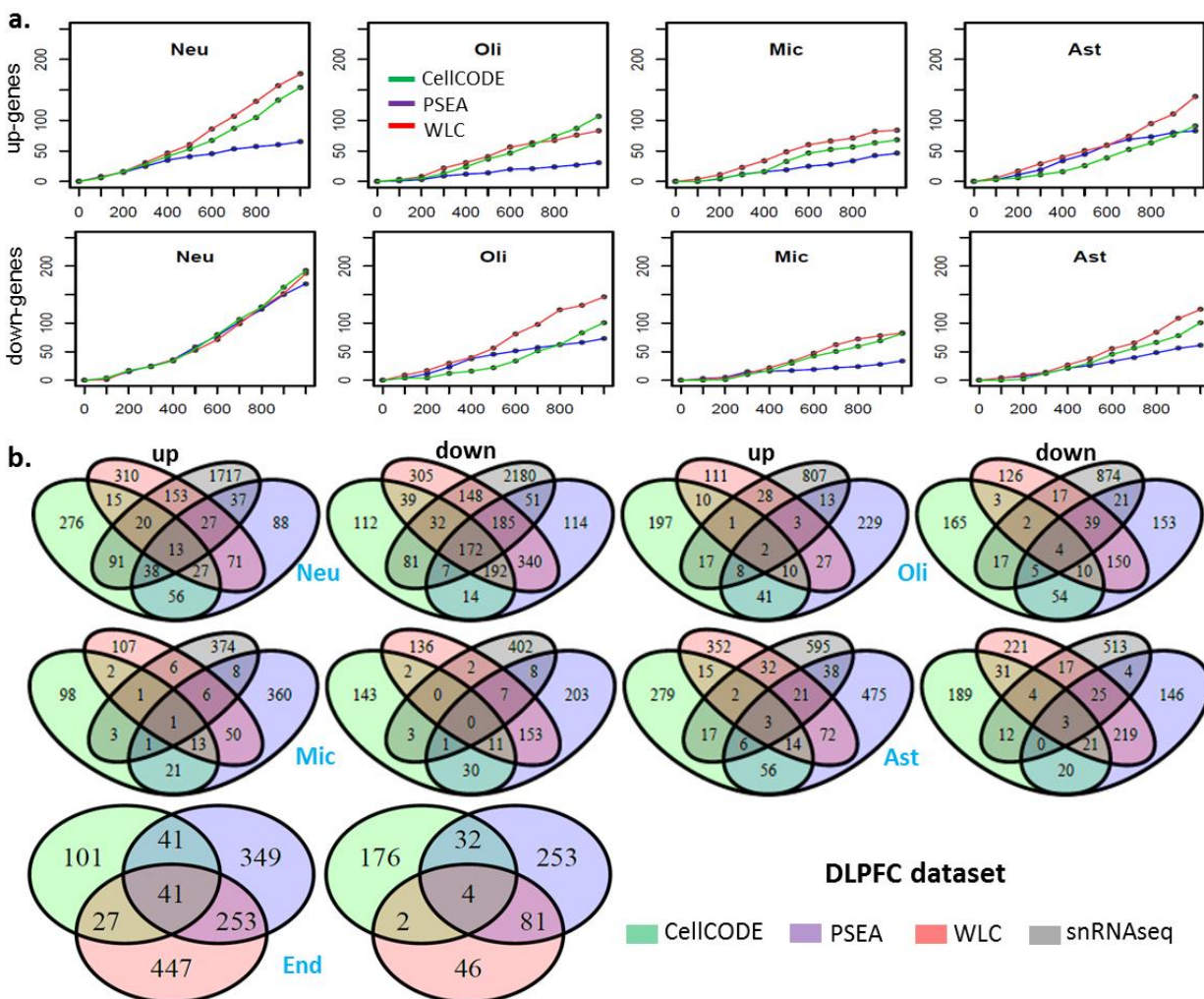
269 Neuronal CI-DEGs retained their significance of overlap across all comparisons and for
270 all algorithms, except for the top 1,000 up-regulated neuronal CI-DEGs deconvoluted with
271 PSEA. Microglial CI-DEGs had the least numbers of significant overlap for their top ranked
272 genes. Astrocytic and oligodendrocytic top ranked CI-DEGs had significance of overlap between
273 the neuronal and microglial results (**Table S18**). These findings are reflective of the abundance
274 of these cell types, with the most abundant neurons having the most overlap for the top ranked
275 CI-DEGs between deconvoluted bulk-DLPFC and snDLPFC.

276 Amongst these comparisons, we determined that the significance for overlap was best for
277 all algorithms for the top ranked 600 genes. Using WLC deconvoluted results, the overlap for the
278 top 600 CI-DEGs from bulk-DLPFC and snDNPFC are statistically significant for all eight
279 comparisons (**Fig 4a**). For the top 600 genes, overlap with CellCODE results is significant for all
280 except down-regulated oligodendrocyte and up-regulated astrocyte CI-DEGs. For PSEA, none of
281 the microglia CI-DEGs had significant overlap. PSEA results for the top 600 genes were
282 otherwise significant for all but up-regulated oligodendrocyte and down-regulated astrocyte
283 genes.

284 We also performed a comparison of CI-DEGs identified at nominal significance (p-value
285 <0.05) with each algorithm from bulk-DLPFC to nominally significant snDLPFC results (**Fig**
286 **4b, Table S19**). As with the above comparison, genes that are either up or down in both
287 deconvoluted bulk-DLPFC and snDLPFC data were analyzed separately for each cell type.

288 Not surprisingly, down-regulated neuronal CI-DEGs have the greatest overlap (537/3732
289 or 14.4% for WLC, 292/3213 or 9.1% for CellCODE, 415/3516 or 11.8% for PSEA). These
290 overlaps are significant for all three algorithms (**Table S19**). Down-regulated CI-DEGs in
291 microglia show the least proportion of overlap (9/723 or 1.2% for WLC, 4/609 or 0.66% for
292 CellCODE, 16/820 or 2.0% for PSEA) (empirical p-value > 0.05). Significant overlap detected
293 with WLC (all but down-regulated microglia) and PSEA (all but microglial results and up-
294 regulated oligodendrocytes) were similar, whereas CellCODE results had significant overlaps
295 only for the neuronal CI-DEGs (**Table S19**).

296



297

298 **Fig 4:** Comparison of CI-DEGs from computational deconvolution with CI-DEGs from
 299 snRNAseq on DLPFC dataset. **a.** Upper panel: number of overlapping genes (y-axis) between
 300 the top N (x-axis) up-regulated genes in snDLPFC and top N up-regulated genes from bulk-
 301 DLPFC deconvoluted with PSEA, WLC and CellCODE, respectively. Lower panel: number of
 302 overlapping genes (y-axis) between the top N (x-axis) down-regulated genes in snDLPFC and
 303 top N down-regulated genes from deconvoluted bulk-DLPFC. **b.** Venn diagram of CI-DEGs
 304 from computational deconvolution methods and those from snRNAseq. Overlap is evaluated for
 305 both bulk-DLPFC and snDLPFC CI-DEGs detected at nominal p value ≤ 0.05 .

306

307 **Discussion**

308 There is an increasing number of large scale RNAseq-based transcriptome datasets
309 generated in bulk tissue for many diseases, including brain tissue from patients with AD, other
310 neurodegenerative diseases and controls^{5-8,12}. Comparison of such transcriptome data from
311 patient and control individuals has been instrumental in the identification of genes and co-
312 expression networks that are altered in and may therefore be risk factors for these diseases^{4-7,44}.
313 The discovery that many of these transcriptional networks harbor genes with disease risk variants
314 provides support for the utility of this bulk-transcriptome approach in deciphering molecules and
315 pathways that are risk factors for these conditions. Nevertheless, there is clear evidence for
316 abundant transcriptome alterations in bulk tissue from affected organs of patients with disease,
317 which appears to be due to changes in cellular composition of that tissue as a consequence of the
318 disease processes^{4,45}. Given this, there is a strong need to detect cell-intrinsic transcriptional
319 alterations to be able to distinguish gene expression changes that are upstream of and may
320 therefore be causative factors for disease from those that are merely a result of cell proportion
321 changes that occur due to disease pathology. The discovery of molecules and pathways that are
322 upstream of and risk factors for disease pathology is of paramount importance for development
323 of targeted therapies. This information can also aid in the identification of refined disease
324 biomarkers reflective of disease-causative expression alterations in these conditions. Detection of
325 cell-specific transcriptional changes can also help develop more accurate disease models
326 harboring these cellular alterations. Further, discovery of cell-specific transcriptional alterations
327 in disease may uncover expression changes, particularly in less abundant cell types, which may
328 be missed by the analysis of bulk transcriptome. Thus, there is a growing effort to identify cell-
329 specific expression alterations in human diseases^{9-11,14,15,17,46-48}.

330 There are two general approaches to decipher cell-specific transcriptional changes in AD.
331 One approach is to conduct single nucleus (snRNAseq), single cell (scRNAseq) or purified cell
332 RNAseq experimentally, followed by data analyses. The alternative approach is to design
333 relatively complex bioinformatics pipelines to decipher cell-intrinsic information of individual
334 cell types from bulk tissue microarray or RNAseq data. The first approach is limited due to
335 significantly higher cost and experimental challenges. Additionally, these approaches may have
336 the drawback that the procedures of dissociating cells break cell-to-cell communication and thus
337 may not reflect the true expression signal *in vivo*. The alternative bioinformatics approach to
338 decipher cell-specific transcriptome alterations from bulk tissue has the potential to avoid the
339 above weaknesses of the experimental approach. Furthermore, a bioinformatics approach can
340 make use of the large amount of existing RNAseq data^{4-8,12} from fresh or frozen bulk brain
341 tissues with minimum added cost, and may better reflect the true situation in brain where
342 different cells are in communication.

343 In this study, applying three different algorithms^{13,14} including our own novel approach,
344 we estimated cell-intrinsic gene expression for deconvoluted cell types from three large bulk
345 RNAseq datasets^{4,8,12,18} from two brain regions. We identified consensus cell-intrinsic
346 transcriptional alterations (CI-DEGs) in AD, which are conserved across cohorts and brain
347 regions. We also performed an in-depth comparison of these CI-DEGs with bulk brain RNAseq
348 data obtained from the same datasets, collectively comprised of 685 unique brain samples. To
349 our knowledge, this is the first study to detect CI-DEGs and their enriched gene ontology (GO)
350 terms in computationally deconvoluted large-scale RNAseq data from AD and control brain
351 samples. Additionally, we conducted a detailed comparison of CI-DEGs deconvoluted from

352 bulk-DLPFC data with the three algorithms and those obtained from snRNAseq¹⁹ (snDLPFC) of
353 a subset of the samples from the same cohort^{7,8}.

354 The main findings from our study can be summarized as follows: **1)** The direction of
355 change in cellular proportions in AD is consistent across two brain regions and three datasets for
356 most cell types, although the magnitude of change seems to vary. Our findings revealed greater
357 neuronal loss and microgliosis in TCX compared to DLPFC. **2)** We identified CI-DEGs and bulk
358 tissue DEGs (bulk-DEGs) independently in two TCX and one DLPFC datasets. The overlap in
359 bulk-DEGs across datasets is greater than that for CI-DEGs. **3)** We performed an in-depth
360 comparison of the consensus CI-DEGs, common to both TCX and DLPFC against the consensus
361 bulk-DEGs detected in these same datasets. We identified significant overlap between consensus
362 CI-DEGs and consensus bulk-DEGs. The extent of overlap between consensus bulk-DEGs and
363 consensus CI-DEGs was greatest for *down-regulated neuronal* genes ($p=5.6E-207$). This was
364 followed by *up-regulated non-neuronal* genes (p ranging from $1.4E-77$ for endothelia to $4.3E-14$
365 for microglia). **4)** Despite the statistically significant overlap between consensus bulk vs. CI-
366 DEGs, the majority of the consensus CI-DEGs for *up-regulated neuronal* and *down-regulated*
367 *non-neuronal* genes were not detected in bulk tissue. This finding highlights the potential ability
368 of computational deconvolution approach to identify CI-DEGs that may be missed in bulk-DEGs
369 especially for genes that are not moving in the direction of cell proportion changes. **5)** We
370 identified GO-terms enriched for consensus CI-DEGs, and detected processes that have
371 previously been implicated in AD as well as novel ones. **6)** Using an snRNAseq¹⁹ dataset as
372 comparison, we assessed the performance of our CI-DEG detection algorithm (WLC), and the
373 published CellCODE¹⁴ and PSEA¹³ approaches. We determined that WLC had comparable or

374 superior performance in the detection CI-DEGs that had significant overlap with snDLPFC
375 results.

376 Our findings highlight the consistency and reproducibility of our findings across two
377 different brain regions from three different studies conducted separately. We identified similar
378 directions of change in AD:Control cell proportions in TCX and DLPFC. As expected from
379 known AD neuropathology, neuronal populations are significantly lower, and microglial
380 populations are significantly higher in AD vs. control brains in all datasets. Consistent with this
381 pattern of reproducibility, we also found significant overlap of consensus CI-DEGs detected in
382 TCX and DLPFC for all cell types and for both directions of change, i.e. up or down in AD, with
383 consensus bulk-DEGs.

384 Using consensus CI-DEGs, we identified GO terms, which include processes and genes
385 that have previously been implicated in AD, thereby providing further validation of our
386 approach. Detailed discussion of all of the pathways identified with the consensus CI-DEGs is
387 beyond the scope of this study. Instead, we herein highlight some of the CI-DEG enriched
388 pathways.

389 Down-regulated neuronal CI-DEGs were enriched in neuropeptide hormone activity
390 (GO:0005184) pathway. These terms include *VGF* (a.k.a. neuroendocrine regulatory peptide 1),
391 which is selectively expressed in some neurons and shown to be reduced in AD and Parkinson's
392 disease²². Corticotropin releasing hormone (CRH), which is a neuronally expressed peptide that
393 mediates stress in the hypothalamic-pituitary-adrenal axis, is also a member of the same GO
394 term. CRH has been implicated in both adverse outcomes related to AD pathology in model
395 systems and epidemiology studies, as well as having an important role in learning and memory²³.
396 Neuronal reduction of *CRH* may either be a potentially neuro-protective response in AD brains

397 or lead to further negative impact on memory. Although the biological implications of our
398 finding remain to be uncovered, our results are aligned with the potential importance of the
399 neuroendocrine system in AD.

400 Interestingly, CI-DEGs also implicate biological processes that are enriched for neuronal-
401 DEGs that are *up* in AD, despite reductions in neuronal cell populations in this condition for both
402 TCX and DLPFC. This suggests that our cell type specific transcriptome deconvolution
403 successfully captures transcript changes that occur in the direction opposite to that of cell-
404 population changes, and that may therefore be missed in bulk-DEG approaches.

405 Indeed, the significant GO terms potassium channel activity (GO:0005267) and
406 regulation of ion transport (GO:0043269) harbor many potassium channels, which are up in AD
407 for neuronal CI-DEGs in both TCX and DLPFC, but do not have consistent results in bulk-DEGs
408 from the same cohorts. These findings highlight the potential utility of cell-specific transcript
409 deconvolution approaches in reducing noise from cell population changes, thereby resulting in
410 consistent detection of true signal. Notably, potassium channels have been reported to be up in
411 AD brains and mouse models of AD⁴⁹⁻⁵¹, leading to their suggestion as a potential therapeutic
412 avenue in this condition.

413 Some consensus CI-DEGs point to AD-related perturbations of key cellular functions for
414 the specific cell type. An example of this is consensus oligodendrocyte CI-DEGs. The *down-*
415 *regulated* oligodendrocyte CI-DEGs are enriched for the myelination GO term (GO:0042552)
416 and those that are *up* in this cell type are enriched in ceramide biosynthetic process
417 (GO:0046513).

418 Down-regulated oligodendrocyte CI-DEGs include myelin basic protein (*MBP*)⁴,
419 plasmolipin (*PLLP*)^{4,5}, myelin and lymphocyte protein (*MAL*), and myelin-associated

420 glycoprotein (*MAG*)²⁶, even though bulk-DEGs for these genes did not show consistent changes.
421 We⁴ and others⁵ demonstrated lower levels in AD of co-expression networks of genes implicated
422 in myelination, which is consistent with the present findings from oligodendrocyte CI-DEGs.

423 Ceramide dysregulation has been implicated in AD^{52,53}. Increased ceramide species were
424 observed in AD and other neuropathological disorders compared to controls⁵³, and the activation
425 of the neutral sphingomyelinase–ceramide pathway induces oligodendrocyte death⁵⁴. Our present
426 observation from oligodendrocyte CI-DEGs are consistent with these findings.

427 Despite the biological insights gained from computationally deconvoluted CI-DEGs, they
428 also have some shortcomings. Compared to bulk-DEGs, CI-DEGs between different datasets
429 show less degree of overlap, regardless of the deconvolution algorithm utilized. This highlights
430 the challenge in the field for the ultimate goal of minimizing detection of transcriptional
431 perturbations due to cell proportion changes while maximizing discovery of those that lead to
432 disease. Put differently, CI-DEGs may enhance discovery of true cell-specific transcriptional
433 changes but this may come at the expense of increased false negative findings. In contrast, bulk-
434 DEGs may capture a larger number of perturbed transcripts, but some may be merely due to cell
435 population differences between diseased and healthy tissue. Ultimately, findings from both
436 approaches may provide the greatest utility in detecting true positives.

437 Comparison of CI-DEGs deconvoluted from bulk-DLPFC and those detected using an
438 orthogonal approach of snRNAseq¹⁹ from the same cohort (ROSMAP) demonstrates the ability
439 of these computational approaches to identify true cell-specific expression changes in AD. Using
440 our in-house WLC algorithm, there was significant overlap with the results from snRNAseq and
441 CI-DEGs of most cell types (**Fig 4, Tables S18-S19**). CellCODE and PSEA also identified
442 significant overlaps, but to a lesser extent, especially for rarer cell types such as microglia.

443 Hence, our WLC algorithm demonstrates at least comparable performance to these two
444 algorithms^{13,14}. This is also supported by the higher degree of overlap among different cohorts
445 for CI-DEGs detected by WLC than the other two algorithms (**Fig 2a**). Due to the challenges in
446 deconvoluting noisy data from human series, different computational approaches may be utilized
447 and combined, and that calls for a more devoted effort in developing such algorithms.

448 In summary, using three distinct computational approaches, we deconvoluted brain bulk-
449 RNAseq data from three large and independent cohorts^{8,12,18}. We detected cell population
450 changes that are observed consistently across cohorts, and congruent with the known disease
451 pathology. Although there is significant overlap between consensus CI-DEGs and consensus
452 bulk-DEGs, there are more unique CI-DEGs that change in the direction opposite to that of cell
453 population changes. This suggests that CI-DEGs may have utility in detecting disease-related
454 transcriptional changes above and beyond those due to cell proportion changes. Consensus CI-
455 DEGs identify GO terms, including those for hormone activity, myelin biology and channel
456 activity. The enriched CI-DEGs include genes previously implicated in AD or
457 neurodegeneration, such as *VGF*, *CRH*, *MOBP* and *MBP*, and other novel genes.

458 This study demonstrates the utility of our analytic approach in deciphering cell-specific
459 transcriptional alterations using bulk tissue in a complex disease, provides a comprehensive
460 comparison of our pipeline to existing ones, identifies patterns of cell proportions in AD and
461 control samples across brain regions, discovers novel CI-DEGs with replication across
462 independent cohorts and highlights biological processes with cell-specific expression changes in
463 AD. These findings are expected to refine discovery of molecular therapeutic targets, bio markers
464 that reflect cellular transcriptional alterations in AD and accelerate generation of more accurate
465 disease models.

466

467 **Methods**

468 **Analysis datasets**

469 We generated the TCX-Mayo data, which consists of temporal cortex RNAseq
470 measurement of 80 AD patients and 28 controls diagnosed according to neuropathologic
471 criteria^{4,12}. RNAseq data were processed and quality control (QC) was conducted as
472 described^{4,12}. ROSMAP DLPFC^{7,8} and TCX-MSBB¹⁸ datasets were downloaded from the AMP-
473 AD Knowledge Portal on Synapse (syn8691134 and syn8691099). We further filtered out non-
474 Caucasian samples and those that had incongruent sex based on provided covariate vs.
475 transcriptome data. All samples were classified as AD or control based on neuropathological
476 data. All TCX-Mayo AD samples had Braak neurofibrillary tangle (NFT) score ≥ 4 . TCX-Mayo
477 controls had Braak score ≤ 3 and were without any neurodegenerative disease diagnoses. TCX-
478 MSBB AD samples had Braak ≥ 4 and CERAD ≥ 2 ; and controls had Braak ≤ 3 and CERAD ≤ 1 .
479 DLPFC AD samples from ROSMAP had Braak score ≥ 4 and CERAD neuritic plaque score ≤ 2 .
480 ROSMAP control samples had Braak ≤ 3 and CERAD neuritic plaque score ≥ 3 .

481 It should be noted that the CERAD⁵⁵ neuritic plaque score as applied by the ROSMAP
482 study is defined such that high CERAD indicates lower neuritic plaque burden and decreased
483 probability of AD. In the MSBB study, higher CERAD indicates higher plaque burden.

484 Raw RNA read counts were normalized using conditional quantile normalization (CQN)
485 method⁵⁶ implemented in R cqn package, as previously described⁴. This normalization takes into
486 consideration library size, gene length and GC content. It also performs a log2 transformation so
487 that the resulting distribution for each gene is Gaussian-like. We also determined covariates that
488 contributed significantly to the variation of gene expression in these RNAseq cohorts (**Fig S5**)
489 for adjustment in the analyses.

490

491 **Cell proportion estimation**

492 Digital sorting algorithm (DSA)¹⁶ was applied to estimate cell proportions through R
493 DSA package, function DSA. For each cell type, DSA first computes the average marker gene
494 expression in the analysis dataset, the purpose of which is to construct a variable that better
495 reflects cell proportion variation among subjects. To reduce the effect of outlier expression that
496 is occasionally seen in RNAseq data, we modified the original DSA so that the median instead of
497 mean expression was computed.

498

499 **CI-DEG analysis for individual cell types**

500 In this study we identified CI-DEGs from deconvoluted bulk RNAseq data using three
501 different algorithms, namely PSEA¹³, CellCODE¹⁴ and our in-house algorithm WLC. All
502 analyses adjusted for the following variables: Sex, RIN, age at death and batch for DLPFC and
503 TCX-MSBB datasets, and sex, RIN and age at death for TCX-Mayo dataset (**Fig S5**).

504 PSEA¹³ applies model selection procedures to select cell type(s) that should be included
505 in baseline (control) or AD condition, and then estimate differential expression in specific cell
506 types (CI-DEGs). We used the R package PSEA to obtain CI-DEG results of PSEA approach,
507 through functions `em_quantvg` (to generate candidate models) and `lmfitst` (to fit all candidate
508 models and pick the best one). Expression values used in PSEA are in linear scale (non log-
509 transformed).

510 CellCODE¹⁴ assesses overall gene expression difference between AD and control groups
511 with adjustment of relative cell proportion, followed by assigning the cell type that correlates
512 best with the difference (CI-DEGs). R package CellCODE was used to obtain DEG results of

513 CellCODE approach, through functions `getAllSPVs` (to construct surrogate variable using
514 marker genes through singular value decomposition) and `lm.coef` (to estimate difference between
515 AD and control groups). Strictly speaking, CellCODE measures the overall differential
516 expression rather than CI-DEGs but identifies the cell type that is most correlated with the
517 observed difference between AD and control using `cellPopT` function. However, for simplicity,
518 we refer to the DEGs from CellCODE as CI-DEGs in this study. Expression values used in
519 CellCODE are in log scale.

520 Our in-house method WLC, described in the method section, applies weighted linear
521 regression with constraints and model selection procedures, which also estimates differential
522 expression in specific cell types (CI-DEGs). It guarantees the fitted relative gene expression is
523 non-negative. By weighing the expression, it smooths out the extreme data points. The
524 procedures of this algorithm is illustrated by the following high level pseudo code.

525

526

527 *Assume cell type 1,2,3,4,5 are neu, oli, mic, ast and end respectively*

528 *Fit Eq.(1) to identify a set of cell type T in which the gene is significantly expressed*

529 *If the set T is not empty:*

530 *Fit Eq.(2) to identify a set of cell type $\Phi \subseteq T$ in which the gene is differentially expressed*

531 *Let $\Omega = \{\text{each cell type in } T, \Phi, T\}$*

532 *For each element $\theta \in \Omega$*

533 *Fit Eq.(3) with adjustment for other covariates $C_k (1 \leq k \leq m)$*

534 *Keep Akaike information criterion (AIC) of this model fitting*

535 *Identify θ_{best} that gives the best AIC.*

536 *Use the estimated values from the best model*

537

$$y_i \sim \alpha_0 + \sum_{t \in \{1,2,3,4,5\}} \alpha_t x_{i,t} + \epsilon, \quad \text{s.t. } \alpha_t \geq 0 \quad (1 \leq t \leq 5) \quad (1)$$

$$y_i \sim \beta_0 + \sum_{t \in T} \beta_t x_{i,t} + d_i \sum_{t \in T} \beta_t^\Delta x_{i,t} + \epsilon, \quad \text{s.t. } \beta_t \geq 0, \beta_t + \beta_t^\Delta \geq 0 \quad (t \in T) \quad (2)$$

$$y_i \sim \gamma_0 + \sum_{t \in T} \gamma_t x_{i,t} + d_i \sum_{t \in \Theta} \gamma_t^\Delta x_{i,t} + \sum_{k=1}^m \lambda_k C_{i,k} + \epsilon, \quad \text{s.t. } \gamma_t \geq 0 \quad (t \in T), \gamma_t + \gamma_t^\Delta \geq 0 \quad (t \in \Theta) \quad (3)$$

538

539 In Eqs.1-3, y_i is the observed expression of a gene in subject i ; $x_{i,t}$ is the median marker gene
540 expression of cell type t in subject i ; $C_{i,k}$ is covariate k in subject i . In Eq.1, α_t is the overall
541 relative expression in cell type t . In Eq.2, β_t is relative expression at the baseline condition in cell
542 type t ; $d_i = 0$ if subject i is in control group, and $d_i = 1$ if subject i is in AD group; therefore,
543 β_t^Δ is the difference of relative gene expression between baseline condition and AD condition in
544 cell type t . Of note, due to the biological meaning these coefficients, they must satisfy constraints
545 such that $\alpha_t, \beta_t, \beta_t + \beta_t^\Delta$ be non-negative. In addition, y_i is in linear scale rather than in log
546 scale⁵⁷, and these non-log-transformed expression values tend to have extreme data points that
547 need to be weighted down. Based on the above considerations, Eq. 1 was fitted by weighted least
548 square linear regression with constraints, which is implemented in lsei function in R package
549 limSolve. The weight of each observation (w_i) is determined by formulae Eq.3 and Eq.4 below.
550 Intuitively, if the expression of a gene in a sample is extremely distant from the median
551 expression of all samples in the same diagnosis group, the weight of that sample is smaller than 1
552 for that gene.

Wang et al.

Main Text

$$w_i = 1 / \left(1 + \frac{|y_i - \text{median}(y_{AD})|}{1 + \text{median}(y_{AD})^2} \right) \quad \text{sample } i \text{ is AD} \quad (4)$$

$$w_i = 1 / \left(1 + \frac{|y_i - \text{median}(y_{control})|}{1 + \text{median}(y_{control})^2} \right) \quad \text{sample } i \text{ is control} \quad (5)$$

553

554 **GO enrichment analysis**

555 Using genes included in the CI-DEG analysis as background genes, p-value for GO term
556 enrichment with consensus CI-DEGs was calculated by “enrichmentAnalysis” function from
557 WGCNA R package³⁵.

558

559 **Comparison of CI-DEGs from computational deconvolution vs. snRNAseq:**

560 To determine if computational bulk-tissue RNAseq could reveal true CI-DEGs, we
561 downloaded and utilized a published snRNAseq study¹⁹ from frozen DLPFC tissues (snDLPFC)
562 which compared gene expression from 24 individuals with AD-pathology to that from 24
563 individuals without AD-pathology in six cell types - excitatory neurons, inhibitory neurons,
564 oligodendrocyte, microglia, oligodendrocyte precursor cells and astrocytes. These snDLPFC
565 samples are from the ROSMAP project, of which we analyzed 474 bulk-DLPFC RNAseq data in
566 our current study. Twenty-four (9 AD cases, 15 controls) of the 48 individuals in snDLPFC
567 study are also included in the bulk-DLPFC dataset. Hence both the snDLPFC and bulk-DLPFC
568 are from the same cohort with significant overlap. Three deconvolution methods were included
569 in this comparison – CellCODE¹⁴, PSEA¹³ and WLC, our in-house method.

570 Two types of comparisons were made between the deconvoluted bulk-DLPFC and
571 snDLPFC results. In the first comparison, we ranked the genes by their p-values of differential
572 expression between AD and control subjects, per cell type. We compared the top N up- or down-
573 genes from the snDLPFC study to those identified by each deconvolution algorithm, per cell

574 type. Genes common to both the deconvoluted bulk-DLPFC and snDLPFC were counted and
575 plotted.

576 In the second set of comparisons, CI-DEGs identified from deconvolution methods at
577 nominal significance (p-value <0.05) were compared to those identified in the snRNAseq data
578 (p-value <0.05).

579 To assess if the observed overlap from each set of analyses is significant with regards to
580 overlap between random selections, we obtained empirical p values from computer simulations
581 described below.

582

583 **Empirical p-value for the number of overlapping genes**

584 The empirical p-values for the number of overlapping genes between snDLPFC and bulk-
585 DLPFC was obtained using a computer simulation as follows. **(A)** Let S_{sn} stand for all genes in
586 snDLFC and S_{bulk} for all genes in bulk-DLFC; **(B)** randomly assign up-regulation on each S_{sn}
587 gene at probability 0.5, and on each S_{bulk} gene at probability 0.5; **(C)** randomly pick N genes
588 from up-genes of S_{sn} , randomly pick N genes from up-genes of S_{bulk} , and count the number of
589 overlapping genes; **(D)** Steps B-C were repeated 10000 times, and the numbers of overlaps
590 $(M_1, M_2, \dots, M_{10000})$ were obtained; **(E)** Let M be the number of observed overlapping genes, and
591 the empirical p-value is $(1 + \text{number of occurrences that } M_i \geq M) / 10001$.

592

593 **Declarations:**

594 **Ethics, consent and permission:**

595 All data were generated from deceased individual's autopsied brain tissue. This study
596 was approved by Mayo Clinic Institutional Review Board.

597

598 **Consent for publication:**

599 All authors reviewed and approved the final manuscript.

600

601 **Availability of Data and Materials:**

602 The data used in this manuscript is available to the research community through the

603 AMP-AD knowledge portal on Sage Synapse as follows:

604 Mayo Clinic RNAseq: <https://www.synapse.org/#!/Synapse:syn8690799>

605 ROSMAP RNAseq: <https://www.synapse.org/#!/Synapse:syn8691134>

606 MSBB RNAseq: <https://www.synapse.org/#!/Synapse:syn8691099>

607 Scripts to perform the analysis and results reported here will be made available upon
608 acceptance of the manuscript for publication.

609

610 **Competing interests:**

611 The authors declare that they have no competing interests.

612

613 **Funding and Acknowledgements:**

614 We thank the patients and their families for their participation, without them these studies
615 would not have been possible. The Mayo RNAseq study data was led by Dr. Nilüfer Ertekin-
616 Taner, Mayo Clinic, Jacksonville, FL as part of the multi-PI U01 AG046139 (MPIs Golde,
617 Ertekin-Taner, Younkin, Price) using samples from The Mayo Clinic Brain Bank. Data
618 collection was supported through funding by NIA grants P50 AG016574, R01 AG032990, U01
619 AG046139, R01 AG018023, U01 AG006576, U01 AG006786, R01 AG025711, R01
620 AG017216, R01 AG003949, NINDS grant R01 NS080820, CurePSP Foundation, and support
621 from Mayo Foundation. This study was in part funded by NIH RF1 AG051504 and R01
622 AG061796 (NET). We thank the Mayo Clinic Advanced Genomic Technology Center staff and
623 bioinformatics core facility for all gene expression measurements. XW thanks Dr. E. Aubrey
624 Thompson for text editing and Jeremy Burgess for provide insightful inputs. XW was funded in
625 part by a pilot grant from the Mayo Clinic ADRC (P50 AG016574).

626 The results published here are in whole or in part based on data obtained from the AMP-
627 AD Knowledge Portal (<https://adknowledgeportal.synapse.org>). Study data were provided by the
628 Rush Alzheimer's Disease Center, Rush University Medical Center, Chicago. Data collection
629 was supported through funding by NIA grants P30AG10161 (ROS), R01AG15819 (ROSMAP;
630 genomics and RNAseq), R01AG17917 (MAP), R01AG30146, R01AG36042 (5hC methylation,
631 ATACseq), RC2AG036547 (H3K9Ac), R01AG36836 (RNAseq), R01AG48015 (monocyte
632 RNAseq) RF1AG57473 (single nucleus RNAseq), U01AG32984 (genomic and whole exome
633 sequencing), U01AG46152 (ROSMAP AMP-AD, targeted proteomics), U01AG46161(TMT
634 proteomics), U01AG61356 (whole genome sequencing, targeted proteomics, ROSMAP AMP-

635 AD), the Illinois Department of Public Health (ROSMAP), and the Translational Genomics
636 Research Institute (genomic).

637 The results published here are in whole or in part based on data obtained from the AMP-
638 AD Knowledge Portal (<https://adknowledgeportal.synapse.org/>). These data were generated from
639 postmortem brain tissue collected through the Mount Sinai VA Medical Center Brain Bank and
640 were provided by Dr. Eric Schadt from Mount Sinai School of Medicine.

641

642 **Author contributions:**

643 X.W., S.L., M.A. and N.E-T conceived the idea. X.W., M.A. and N.E-T developed and designed
644 the study, and wrote the manuscript. X.W. and S.L. developed the mathematical theory and X.W.
645 performed the computations. D.D. provided tissue from the Mayo Clinic Brain Bank and
646 neuropathologically characterized TCX-Mayo brain samples. T.N., K.M., S.L., M.A. and M.C.
647 performed sample selection and library preparation of the TCX-Mayo dataset. M.A., Z.Q., T.C.,
648 T.P. and J.R. performed quality control of DLPFC and TCX-MSBB datasets. J.C. and Y.A.
649 supervised the analytical aspects of the project. N.E-T. provided funding, supervision and
650 direction for the whole study. All authors provided critical feedback and approved the final
651 manuscript.

652

653

654 **References:**

- 655 1. 2018 Alzheimer's disease facts and figures. *Alzheimer's & Dementia* **14**, 367-429 (2018).
656 2. Mehta, D., Jackson, R., Paul, G., Shi, J. & Sabbagh, M. Why do trials for Alzheimer's disease drugs
657 keep failing? A discontinued drug perspective for 2010-2015. *Expert opinion on investigational*
658 *drugs* **26**, 735-739 (2017).
659 3. Anderson, R.M., Hadjichrysanthou, C., Evans, S. & Wong, M.M. Why do so many clinical trials of
660 therapies for Alzheimer's disease fail? *The Lancet* **390**, 2327-2329 (2017).
661 4. Allen, M. *et al.* Conserved brain myelination networks are altered in Alzheimer's and other
662 neurodegenerative diseases. *Alzheimer's & Dementia* (2017).
663 5. McKenzie, A.T. *et al.* Multiscale network modeling of oligodendrocytes reveals molecular
664 components of myelin dysregulation in Alzheimer's disease. *Molecular Neurodegeneration* **12**,
665 82 (2017).
666 6. Zhang, B. *et al.* Integrated systems approach identifies genetic nodes and networks in late-onset
667 Alzheimer's disease. *Cell* **153**, 707-720 (2013).
668 7. Mostafavi, S. *et al.* A molecular network of the aging human brain provides insights into the
669 pathology and cognitive decline of Alzheimer's disease. *Nature Neuroscience* **21**, 811-819
670 (2018).
671 8. De Jager, P.L. *et al.* A multi-omic atlas of the human frontal cortex for aging and Alzheimer's
672 disease research. *Scientific Data* **5**, 180142 (2018).
673 9. Zhang, Y. *et al.* Purification and Characterization of Progenitor and Mature Human Astrocytes
674 Reveals Transcriptional and Functional Differences with Mouse. *Neuron* **89**, 37-53 (2016).
675 10. Darmanis, S. *et al.* A survey of human brain transcriptome diversity at the single cell level.
676 *Proceedings of the National Academy of Sciences* **112**, 7285-7290 (2015).
677 11. Lake, B.B. *et al.* Integrative single-cell analysis of transcriptional and epigenetic states in the
678 human adult brain. *Nature biotechnology* **36**, 70-80 (2018).
679 12. Allen, M. *et al.* Human whole genome genotype and transcriptome data for Alzheimer's and
680 other neurodegenerative diseases. *Scientific Data* **3**, 160089 (2016).
681 13. Kuhn, A., Thu, D., Waldvogel, H.J., Faull, R.L.M. & Luthi-Carter, R. Population-specific expression
682 analysis (PSEA) reveals molecular changes in diseased brain. *Nat Meth* **8**, 945-947 (2011).
683 14. Chikina, M., Zaslavsky, E. & Sealfon, S.C. CellCODE: a robust latent variable approach to
684 differential expression analysis for heterogeneous cell populations. *Bioinformatics* **31**, 1584-
685 1591 (2015).
686 15. McKenzie, A.T. *et al.* Brain Cell Type Specific Gene Expression and Co-expression Network
687 Architectures. *Scientific Reports* **8**, 8868 (2018).
688 16. Zhong, Y., Wan, Y.-W., Pang, K., Chow, L.M. & Liu, Z. Digital sorting of complex tissues for cell
689 type-specific gene expression profiles. *BMC Bioinformatics* **14**, 89 (2013).
690 17. Kuhn, A. *et al.* Cell population-specific expression analysis of human cerebellum. *BMC genomics*
691 **13**, 610 (2012).
692 18. Wang, M. *et al.* The Mount Sinai cohort of large-scale genomic, transcriptomic and proteomic
693 data in Alzheimer's disease. *Sci Data* **5**, 180185 (2018).
694 19. Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer's disease. *Nature* (2019).
695 20. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology
696 Consortium. *Nature genetics* **25**, 25-29 (2000).
697 21. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network analysis.
698 *BMC Bioinformatics* **9**, 559 (2008).

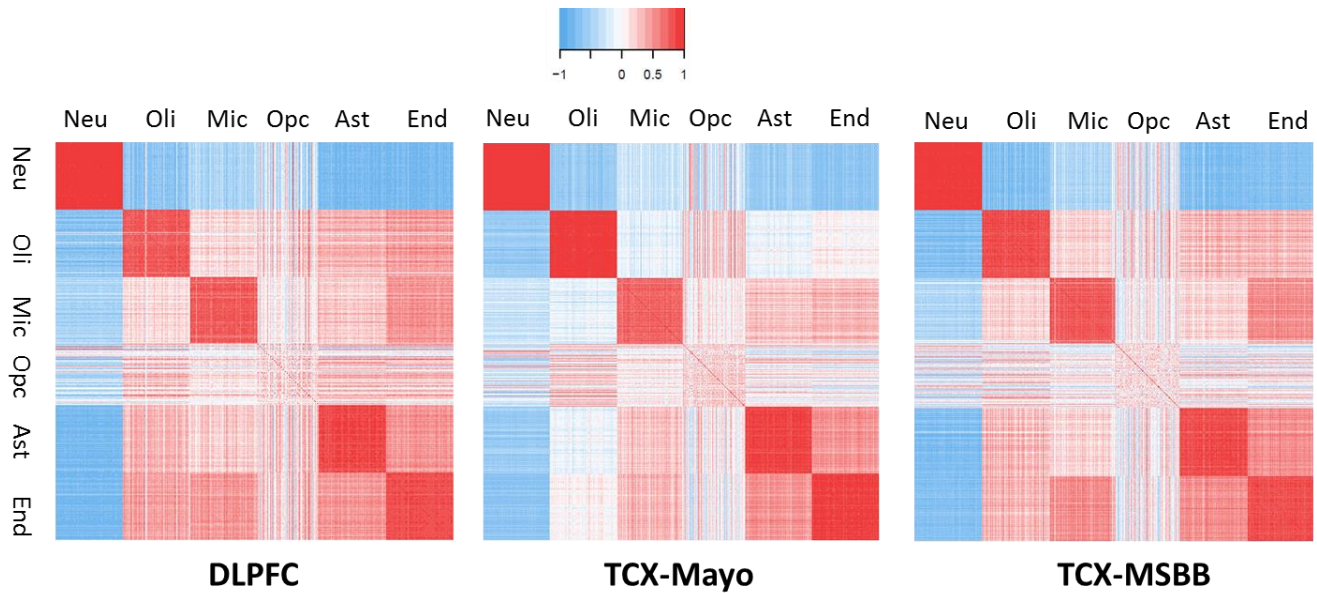
- 699 22. Cocco, C. *et al.* Distribution of VGF peptides in the human cortex and their selective changes in
700 Parkinson's and Alzheimer's diseases. *Journal of anatomy* **217**, 683-93 (2010).
- 701 23. Futch, H.S., Croft, C.L., Truong, V.Q., Krause, E.G. & Golde, T.E. Targeting psychologic stress
702 signaling pathways in Alzheimer's disease. *Molecular neurodegeneration* **12**, 49 (2017).
- 703 24. Wang, R. & Reddy, P.H. Role of Glutamate and NMDA Receptors in Alzheimer's Disease. *J*
704 *Alzheimers Dis* **57**, 1041-1048 (2017).
- 705 25. Berchtold, N.C. *et al.* Brain gene expression patterns differentiate mild cognitive impairment
706 from normal aged and Alzheimer's disease. *Neurobiol Aging* **35**, 1961-72 (2014).
- 707 26. McAleese, K.E. *et al.* Parietal white matter lesions in Alzheimer's disease are associated with
708 cortical neurodegenerative pathology, but not with small vessel disease. *Acta neuropathologica*
709 **134**, 459-473 (2017).
- 710 27. Crivelli, S.M. *et al.* Sphingolipids in Alzheimer's disease, how can we target them? *Adv Drug Deliv*
711 *Rev* (2020).
- 712 28. Olsen, A.S.B. & Faergeman, N.J. Sphingolipids: membrane microdomains in brain development,
713 function and neurological diseases. *Open Biol* **7**(2017).
- 714 29. El Gaamouch, F. *et al.* VGF-derived peptide TLQP-21 modulates microglial function through
715 C3aR1 signaling pathways and reduces neuropathology in 5xFAD mice. *Molecular*
716 *neurodegeneration* **15**, 4-4 (2020).
- 717 30. Tzeng, T.-C. *et al.* Inflammasome-derived cytokine IL18 suppresses amyloid-induced seizures in
718 Alzheimer-prone mice. *Proceedings of the National Academy of Sciences of the United States of*
719 *America* **115**, 9002-9007 (2018).
- 720 31. White, C.S., Lawrence, C.B., Brough, D. & Rivers-Auty, J. Inflammasomes as therapeutic targets
721 for Alzheimer's disease. *Brain Pathol* **27**, 223-234 (2017).
- 722 32. Gao, T. *et al.* Transcriptional regulation of homeostatic and disease-associated-microglial genes
723 by IRF1, LXR β , and CEBP α . *Glia* **67**, 1958-1975 (2019).
- 724 33. Deitmer, J.W., Theparambil, S.M., Ruminot, I., Noor, S.I. & Becker, H.M. Energy Dynamics in the
725 Brain: Contributions of Astrocytes to Metabolism and pH Homeostasis. *Frontiers in Neuroscience*
726 **13**(2019).
- 727 34. Bélanger, M., Allaman, I. & Magistretti, Pierre J. Brain Energy Metabolism: Focus on Astrocyte-
728 Neuron Metabolic Cooperation. *Cell Metabolism* **14**, 724-738 (2011).
- 729 35. Zhang, Q. *et al.* Integrated proteomics and network analysis identifies protein hubs and network
730 alterations in Alzheimer's disease. *Acta neuropathologica communications* **6**, 19-19 (2018).
- 731 36. Lovell, M.A., Xie, C., Gabbita, S.P. & Markesbery, W.R. Decreased thioredoxin and increased
732 thioredoxin reductase levels in alzheimer's disease brain. *Free Radical Biology and Medicine* **28**,
733 418-427 (2000).
- 734 37. Kajiwara, Y. *et al.* GJA1 (connexin43) is a key regulator of Alzheimer's disease pathogenesis. *Acta*
735 *Neuropathol Commun* **6**, 144 (2018).
- 736 38. de Quervain, D.J. & Papassotiropoulos, A. Identification of a genetic cluster influencing memory
737 performance and hippocampal activity in humans. *Proc Natl Acad Sci U S A* **103**, 4270-4 (2006).
- 738 39. Shi, Y. *et al.* Rapid endothelial cytoskeletal reorganization enables early blood-brain barrier
739 disruption and long-term ischaemic reperfusion brain injury. *Nature Communications* **7**, 10523
740 (2016).
- 741 40. Stamatovic, S.M., Keep, R.F. & Andjelkovic, A.V. Brain endothelial cell-cell junctions: how to
742 "open" the blood brain barrier. *Current neuropharmacology* **6**, 179-192 (2008).
- 743 41. Evans, H.T., Benetatos, J., van Roijen, M., Bodea, L.G. & Gotz, J. Decreased synthesis of
744 ribosomal proteins in tauopathy revealed by non-canonical amino acid labelling. *EMBO J* **38**,
745 e101174 (2019).

- 746 42. Garcia-Esparcia, P. *et al.* Altered machinery of protein synthesis is region- and stage-dependent
747 and is associated with alpha-synuclein oligomers in Parkinson's disease. *Acta Neuropathol*
748 *Commun* **3**, 76 (2015).
- 749 43. Koren, S.A. *et al.* Tau drives translational selectivity by interacting with ribosomal proteins. *Acta*
750 *Neuropathol* **137**, 571-583 (2019).
- 751 44. Conway, O.J. *et al.* ABI3 and PLCG2 missense variants as risk factors for neurodegenerative
752 diseases in Caucasians and African Americans. *Molecular neurodegeneration* **13**, 53 (2018).
- 753 45. Srinivasan, K. *et al.* Untangling the brain's neuroinflammatory and neurodegenerative
754 transcriptional responses. *Nature communications* **7**, 11295 (2016).
- 755 46. Lin, M.-Y. *et al.* Releasing Syntaphilin Removes Stressed Mitochondria from Axons Independent
756 of Mitophagy under Pathophysiological Conditions. *Neuron* **94**, 595-610.e6 (2017).
- 757 47. Newberg, L.A., Chen, X., Kodira, C.D. & Zavodszky, M.I. Computational de novo discovery of
758 distinguishing genes for biological processes and cell types in complex tissues. *PloS one* **13**,
759 e0193067 (2018).
- 760 48. Li, Z. *et al.* Genetic variants associated with Alzheimer's disease confer different cerebral cortex
761 cell-type population structure. *Genome medicine* **10**, 43 (2018).
- 762 49. Angulo, E. *et al.* Up-regulation of the Kv3.4 potassium channel subunit in early stages of
763 Alzheimer's disease. *J Neurochem* **91**, 547-57 (2004).
- 764 50. Maezawa, I. *et al.* Kv1.3 inhibition as a potential microglia-targeted therapy for Alzheimer's
765 disease: preclinical proof of concept. *Brain* **141**, 596-612 (2018).
- 766 51. Yi, M. *et al.* KCa3.1 constitutes a pharmacological target for astrogliosis associated with
767 Alzheimer's disease. *Mol Cell Neurosci* **76**, 21-32 (2016).
- 768 52. Czubowicz, K., Jęśko, H., Wencel, P., Lukiw, W.J. & Strosznajder, R.P. The Role of Ceramide and
769 Sphingosine-1-Phosphate in Alzheimer's Disease and Other Neurodegenerative Disorders.
770 *Molecular neurobiology* **56**, 5436-5455 (2019).
- 771 53. Filippov, V. *et al.* Increased ceramide in brains with Alzheimer's and other neurodegenerative
772 diseases. *Journal of Alzheimer's disease : JAD* **29**, 537-547 (2012).
- 773 54. Lee, J.-T. *et al.* Amyloid-beta peptide induces oligodendrocyte death by activating the neutral
774 sphingomyelinase-ceramide pathway. *The Journal of cell biology* **164**, 123-131 (2004).
- 775 55. Mirra, S.S. *et al.* The Consortium to Establish a Registry for Alzheimer's Disease (CERAD). Part II.
776 Standardization of the neuropathologic assessment of Alzheimer's disease. *Neurology* **41**, 479-
777 86 (1991).
- 778 56. Hansen, K.D., Irizarry, R.A. & Wu, Z. Removing technical variability in RNA-seq data using
779 conditional quantile normalization. *Biostatistics* **13**, 204-216 (2012).
- 780 57. Zhong, Y. & Liu, Z. Gene expression deconvolution in linear space. *Nat Meth* **9**, 8-9 (2012).

781

1

2



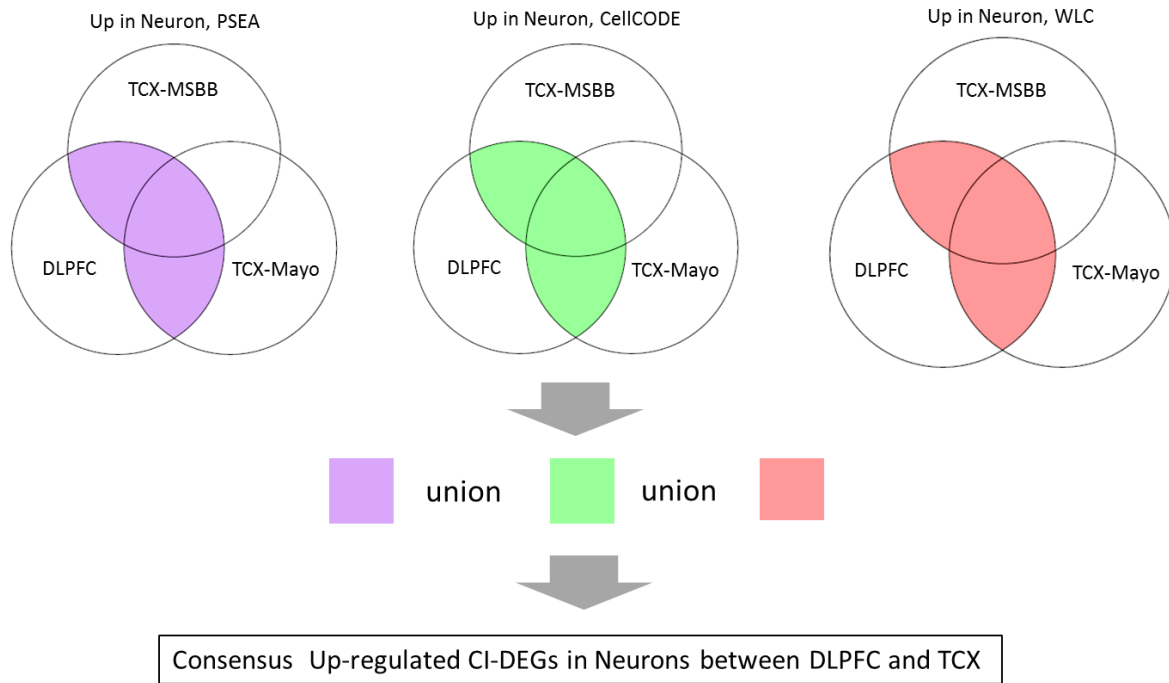
3

4 **Figure S1:** Pearson correlation of estimated cell proportions using different sets of marker genes
5 randomly selected. The steps of this simulation study are as follows. (1) Obtain candidate markers from
6 R package BRETIGEA¹ for neurons (top 500 out of 1000 listed), oligodendrocyte (top 500 out of 1000
7 listed), microglia (top 500 out of 1000 listed), OPC (top 250 out of 500 listed), astrocyte (top 500 out of
8 1000 listed) and endothelial (top 500 out of 1000 listed). (2) From candidate markers, randomly select
9 1/10 genes for each cell type, that is 50 selected genes for all cell types except OPC which has 25
10 selected genes. (3) Estimate cell proportion using DSA algorithm². (4) Repeat (2)-(3) 100 times, and
11 compute Pearson correlation of estimated cell proportion between different runs.

12

13

14



15

16

17 **Figure S2:** Illustration of obtaining consensus CI-DEGs between DLPFC and TCX regions,
18 using up-regulated CI-DEGs in neuronal cells as an example.

19

20

21

22

23

24

25

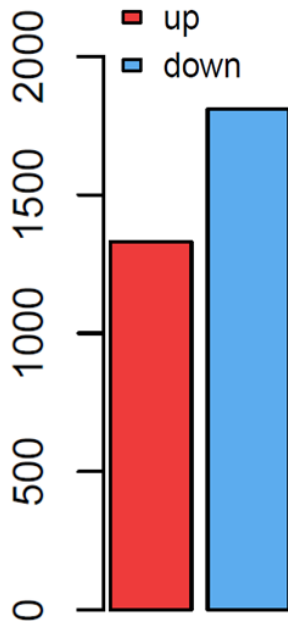
26

27

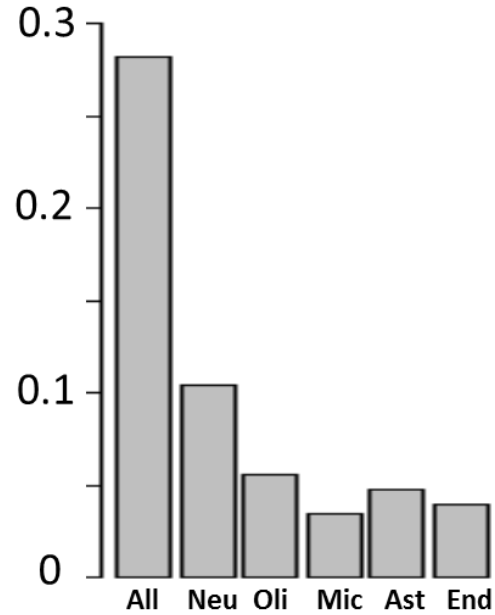
28

29

Consensus bulk-DEGs



% Consensus bulk-DEGs that are cell type markers



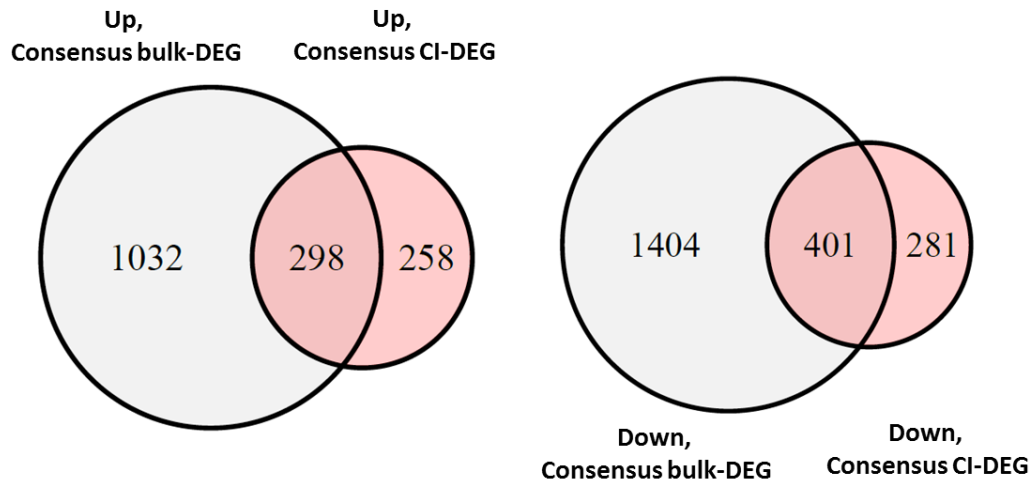
30

31 **Figure S3:** Left panel: the number of up-regulated and down-regulated consensus bulk-DEGs between
32 DLPFC and TCX-Mayo, or between DLPFC and TCX-MSSM. Right panel: percent of consensus CI-
33 DEGs that are also cell type marker genes. Cell type markers are from BRETIGEA¹, containing 1000
34 markers for each of the five cell types.

35

36

37



38

39 **Figure S4:** Venn diagram for all consensus bulk-DEGs and consensus CI-DEGs.

40

41

42

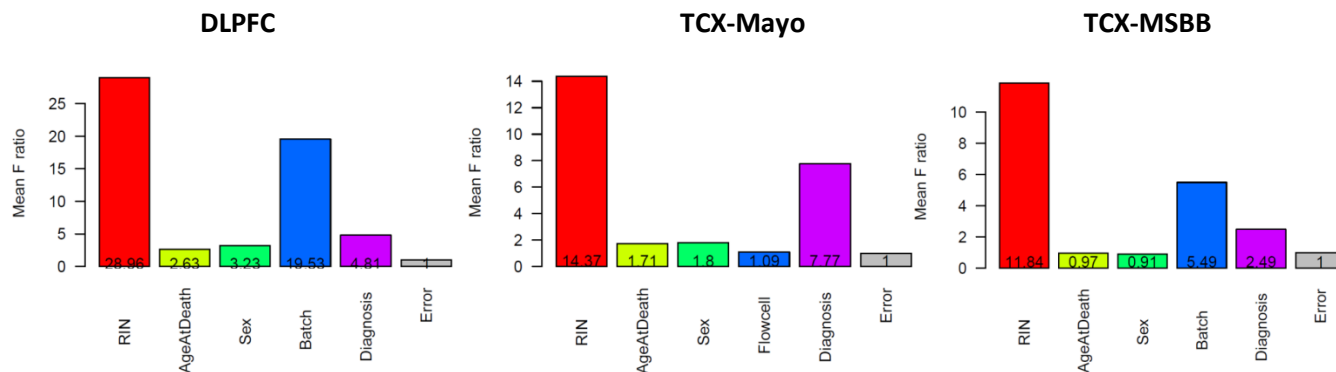
43

44

45

46

47



48

49

50 **Figure S5:** Source of variance analysis of in the DLPFC RNaseq dataset (left panel), TCX-Mayo
51 (middle panel) and TCX-MSBB (right panel). For each gene, a full model was fitted in which cqn
52 normalized gene expression with gene expression as dependent variable and RIN, age at death, sex,
53 batch, and diagnosis group as independent variables (for DLPFC); RIN, age at death, sex, flowcell, and
54 diagnosis group as independent variables (for TCX-Mayo); RIN, age at death, sex, batch, and diagnosis
55 group as independent variables (for TCX-MSBB). Partial models were fitted using the same dependent
56 variable and all but one independent variable. F statistics were obtained by comparing the full model and
57 partial model for each independent variable. Y-axis is the mean of values of F statistics over all genes. In
58 DLPFC and TCX-MSBB, diagnosis, age at death, sex, RIN and batch contributed more than random
59 errors to the variation of gene expression, whereas in TCX-Mayo diagnosis, age at death, sex, and RIN
60 contributed more than random errors to the variation of gene expression.

61

62

63 **References:**

- 64 1. McKenzie, A.T. *et al.* Brain Cell Type Specific Gene Expression and Co-expression Network
65 Architectures. *Scientific Reports* **8**, 8868 (2018).
66 2. Zhong, Y., Wan, Y.-W., Pang, K., Chow, L.M. & Liu, Z. Digital sorting of complex tissues for
67 cell type-specific gene expression profiles. *BMC Bioinformatics* **14**, 89 (2013).

68