

## **Variation among intact tissue samples reveals the core transcriptional features of human CNS cell classes**

Kevin W. Kelley<sup>1-4</sup>, Hiromi Inoue<sup>2,3</sup>, Anna V. Molofsky<sup>2,3</sup>, Michael C. Oldham<sup>1,2\*</sup>

<sup>1</sup>Dept. of Neurological Surgery, UCSF

<sup>2</sup>Eli and Edythe Broad Center of Regeneration Medicine and Stem Cell Research, UCSF

<sup>3</sup>Dept. of Psychiatry, UCSF

<sup>4</sup>Medical Scientist Training Program and Neuroscience Graduate Program, UCSF

\*Correspondence: [Michael.Oldham@ucsf.edu](mailto:Michael.Oldham@ucsf.edu)

## 1 **ABSTRACT**

2 It is widely assumed that cells must be physically isolated to study their molecular profiles. However,  
3 intact tissue samples naturally exhibit variation in cellular composition, which drives covariation of cell-  
4 class-specific molecular features. By analyzing transcriptional covariation in 7221 intact CNS samples  
5 from 840 individuals representing billions of cells, we reveal the core transcriptional identities of major  
6 CNS cell classes in humans. By modeling intact CNS transcriptomes as a function of variation in cellu-  
7 lar composition, we identify cell-class-specific transcriptional differences in Alzheimer's disease, among  
8 brain regions, and between species. Among these, we show that *PMP2* is expressed by human but not  
9 mouse astrocytes and significantly increases mouse astrocyte size upon ectopic expression *in vivo*, caus-  
10 ing them to more closely resemble their human counterparts. Our work is available as an online resource  
11 (<http://oldhamlab.ctec.ucsf.edu>) and provides a generalizable strategy for determining the core molecu-  
12 lar features of cellular identity in intact biological systems.

## 13 **INTRODUCTION**

14 Identifying the molecular features that define cellular identities is a fundamental goal of biological re-  
15 search. Consequently, a number of methods have been developed to isolate cells for molecular profiling,  
16 including fluorescence-activated cell sorting (FACS), immunopanning (IP), and sorting of single cells  
17 (SC) or nuclei (SN). Although these methods are readily applied to many biological systems, their ap-  
18 plicability to the adult human CNS has been limited by technical factors and practical considerations.  
19 For example, FACS, IP, and SC typically require fresh tissue and have therefore been mostly limited to  
20 surgical samples from a handful of CNS regions and individuals<sup>1-3</sup>. SN<sup>4,5</sup> is compatible with frozen tis-  
21 sue but, like SC, suffers from technical noise caused by tissue dissociation, nucleus/cell capture, cDNA  
22 preamplification, and stochastic transcript coverage<sup>6</sup>. Furthermore, there is a trade-off between sequenc-  
23 ing depth and the number of nuclei/cells that can be analyzed.

24 The adult human CNS is large, heterogeneous, and difficult to dissociate due to extensive mye-  
25 lin. It consists of ~170 billion cells, about half of which are neurons<sup>7</sup>. The remaining cells consist mostly  
26 of oligodendrocytes, astrocytes, and microglia, which are collectively referred to as glia. Identifying  
27 transcriptional differences among neuronal and glial subtypes is an important goal, since the extent of  
28 heterogeneity among major CNS cell classes is not fully understood. However, overlooked in the focus  
29 on heterogeneity is the equally important question of what CNS cell subtypes have in common. For ex-  
30 ample, is there a core set of genes whose expression is shared among all neurons? All astrocytes? Etc.  
31 Answering these questions will fill critical gaps in our understanding of CNS cell biology, produce nov-  
32 el experimental and analytical strategies, and provide important insights into the cellular origins of CNS  
33 pathologies. ‘Bottom-up’ methods such as SC/SN are poorly suited to address these questions, since  
34 they are difficult to apply to the adult human CNS at scale.

35 Most gene expression studies of the human CNS have analyzed intact postmortem tissue sam-  
36 ples. Because these samples are heterogeneous and cells must be destroyed to extract RNA, it is often  
37 assumed that these datasets contain no information about the cellular origins of gene expression. How-  
38 ever, it is axiomatic that intact tissue samples from any biological system will exhibit variation in cellu-  
39 lar composition. Therefore, when many intact tissue samples are analyzed, genes expressed with the  
40 greatest sensitivity and specificity in the same cell class should appear highly correlated, since their ex-  
41 pression levels depend primarily on the proportion of that cell class in each sample<sup>8</sup>. In support of this  
42 reasoning, we previously discovered highly reproducible gene coexpression modules in microarray data  
43 from intact human brain samples that were significantly enriched with markers of major CNS cell clas-  
44 ses<sup>9</sup>. These findings were replicated in studies of intact CNS transcriptomes from mice<sup>10</sup>, rats<sup>11</sup>, zebra  
45 finches<sup>12</sup>, macaques<sup>13</sup>, and humans<sup>14</sup>.

46 Gene coexpression modules corresponding to major cell classes are therefore robust and predict-  
47 able features of CNS transcriptomes derived from intact tissue samples. Furthermore, the same genes  
48 consistently show the strongest affinities for these modules, offering substantial information about the  
49 molecular correlates of cellular identity<sup>9</sup>. Over the past decade, thousands of intact, neurotypical human  
50 tissue samples from every major CNS region have been transcriptionally profiled with multiple technol-  
51 ogy platforms. These data provide an unprecedented opportunity to determine the core transcriptional  
52 features of cellular identity in the human CNS from the ‘top down’ by integrating cell-class-specific  
53 gene coexpression modules from a large number of independent datasets.

## 54 **RESULTS**

### 55 **Gene coexpression analysis of synthetic brain samples accurately predicts differential expression** 56 **among CNS cell classes**

57 To illustrate the premise of our approach, we aggregated single-cell RNA-seq data from the adult human  
58 brain<sup>1</sup> to create synthetic samples that mimic the heterogeneity of intact tissue (**Fig. 1A**). We performed  
59 unsupervised gene coexpression analysis of synthetic datasets and identified modules of coexpressed  
60 genes in each dataset that were maximally enriched with published markers<sup>15, 16</sup> of astrocytes, oligoden-  
61 drocytes, microglia, or neurons (‘cell-class modules’; e.g. **Fig. 1A**). Intuitively, the primary source of  
62 expression variation in a cell-class module is variation in the representation of that cell class in each  
63 sample. Mathematically, the vector that explains the greatest amount of expression variation in a coex-  
64 pression module is its first principal component, or module ‘eigengene’ (**Fig. 1A**)<sup>17</sup>. This line of reason-  
65 ing suggests that the eigengene of a cell-class module should approximate the relative abundance of that  
66 cell class in each sample. Because the precise cellular composition of each synthetic sample is known,  
67 we tested this hypothesis and found that actual cellular abundance was nearly indistinguishable from that  
68 predicted by cell-class module eigengenes (**Fig. S1A**).

69 To determine the affinity of each gene for each significant cell-class module, we calculated the  
70 Weighted Gene Coexpression Network Analysis measure of intramodular connectivity, or  $k_{ME}$ <sup>18</sup>.  $k_{ME}$  is  
71 defined as the Pearson correlation between the expression pattern of a gene and a module eigengene. In  
72 the special situation of a cell-class module,  $k_{ME}$  therefore quantifies the similarity between the expres-  
73 sion pattern of a gene and the relative abundance of that cell class in each sample. Because each sample  
74 is a heterogeneous mixture of cells, a high  $k_{ME}$  value for a cell-class module suggests that expression of  
75 the gene in that particular cell class is sensitive and specific. We tested this hypothesis by performing  
76 differential expression analysis of single-cell RNA-seq data for each cell class, restricting our analysis to  
77 exactly the same cells that were used to construct the synthetic brain samples. As shown in **Fig. 1B**, the



78 genes that are most significantly up-regulated in a given cell class also have the highest  $k_{ME}$  values for  
79 the corresponding cell-class module. We obtained nearly identical results by aggregating single-cell  
80 RNA-seq data from the adult mouse brain<sup>19</sup> (**Fig. S1B,C**). These findings demonstrate that gene coex-  
81 pression analysis of intact CNS samples can determine which genes are most differentially expressed  
82 among CNS cell classes. More generally, our results suggest that it is not always necessary to physically  
83 isolate cells in order to ascertain their defining transcriptional features.

84

## 85 **Integrative gene coexpression analysis of intact tissue samples reveals consensus transcriptional** 86 **profiles of major CNS cell classes in humans**

87 To determine consensus transcriptional profiles of human CNS cell classes, we analyzed 7221 CNS  
88 transcriptomes from 840 neurotypical adult humans by combining data from eight studies<sup>14, 20-26</sup> and one  
89 resource ([www.brainspan.org](http://www.brainspan.org)). These data were generated from intact postmortem tissue samples using  
90 diverse technology platforms (**Table S1**) and collectively represent billions of cells. Each sample was  
91 assigned to one of 19 broad neuroanatomical regions, resulting in 62 regional datasets (**Fig. 1C**). After  
92 data preprocessing and quality control, each dataset consisted of  $\geq 25$  samples (median: 76) (**Table S1**).  
93 For each dataset, we performed unsupervised gene coexpression analysis and identified the module that  
94 was maximally enriched with published markers<sup>15, 16</sup> of astrocytes, oligodendrocytes, microglia, or neu-  
95 rons (**Fig. 1D, Table S2**). PC1 of these modules was used to estimate the relative abundance of each cell  
96 class over all samples and calculate genome-wide  $k_{ME}$  values (**Fig. 1E,F**). Finally, we combined  $k_{ME}$   
97 values for significant cell-class modules from all 62 datasets, producing a single value (z-score) for each  
98 gene that quantifies its global expression *fidelity* for each cell class (**Fig. 1G**). Importantly, estimates of  
99 fidelity were highly robust to the choice of gene set used for enrichment analysis (especially for glia;  
100 **Fig. S2**). Canonical markers consistently had high fidelity for the expected cell class and low fidelity for  
101 other cell classes (**Fig. 2A-D**). High-fidelity genes were also significantly and specifically enriched with  
102 expected cell-class markers from multiple independent studies (**Fig. 2A-D**). Compared to glia, the distri-  
103 bution of expression fidelity for neurons was compressed (**Fig. 2A-D**), likely reflecting neuronal hetero-  
104 geneity among CNS regions. Genome-wide estimates of expression fidelity for major cell classes are  
105 provided in **Table S3** and on our web site (<http://oldhamlab.ctec.ucsf.edu>).

## 106 **High-fidelity genes reveal the core transcriptional identities of major CNS cell classes in humans**

107 The genes with the highest expression fidelity for major CNS cell classes are consistently coexpressed  
108 across regions and technology platforms (**Fig. S3**). This consistency suggests that high-fidelity genes  
109 can provide an unbiased view of the core transcriptional identities of major cell classes, thereby reveal-

110 ing novel cellular functions and biomarkers. We visualized the top 50 genes ranked by expression fidelity  
111 ty for each cell class to compare their expression levels, mutation intolerance, literature citations, cellular  
112 lar localization, and protein-protein interactions (PPI) (**Fig. 3A-D**). Overall, absolute expression levels  
113 of high-fidelity genes were highest for neurons and lowest for microglia (**Fig. 3A-D**, red tracks). How-  
114 ever, for each cell class there was a wide range of expression levels for high-fidelity genes, suggesting  
115 parallel regulatory mechanisms and/or differential transcript stability.

116 To assess the tolerance of high-fidelity genes to loss-of-function (LoF) mutations, we analyzed  
117 data from the Exome Aggregation Consortium (ExAC), which summarizes the prevalence of coding mu-  
118 tations in ~61K human exomes<sup>27</sup>. Unexpectedly, high-fidelity neuronal genes were significantly less tol-  
119 erant to LoF mutations than high-fidelity glial genes (**Fig. 3A-D**, black tracks). To determine whether  
120 high-fidelity genes have been studied in their respective cell classes, we searched PubMed for each gene  
121 symbol and the name of the cell class (**Fig. 3A-D**, green tracks). Interestingly, many searches returned  
122 no citations, highlighting critical gaps in our understanding of CNS cell biology. For example, the top  
123 microglial gene (amyloid beta precursor protein binding family B member 1 interacting protein, or  
124 *APBB1IP*) is unstudied in microglia.

125 We examined the cellular localization of proteins<sup>28</sup> encoded by high-fidelity genes and observed  
126 another distinction between neurons and glia. Among the proteins encoded by genes in **Fig. 3A-D**,  
127 membrane localization was reported for 33 in astrocytes, 22 in oligodendrocytes, and 30 in microglia,  
128 but only 13 in neurons (inside track). This result may reflect the homeostatic functions of glia as sensors  
129 and regulators of extracellular CNS environments. More generally, the non-random distributions of cel-  
130 lular localizations suggest that high-fidelity genes are expressed at the protein level in the corresponding  
131 cell classes. To further explore this topic, we examined PPI<sup>29</sup> among high-fidelity gene products for each  
132 cell class and observed significantly more interactions than expected by chance (**Fig. 3A-D**, interior  
133 lines).

134 Because high-fidelity genes should encode optimal biomarkers, we searched for high-fidelity  
135 genes in the Human Protein Atlas (<http://www.proteinatlas.org>) to identify novel reagents for labeling  
136 human CNS cell classes. We identified validated antibodies for PON2 (astrocytes), DBNDD2 (oli-  
137 godendrocytes), APBB1IP (microglia), and CELF2 (neurons) (**Fig. 3A-D**). Dual immunostaining with  
138 canonical markers revealed almost perfect concordance in human frontal cortex (**Fig. 3E-H**).

139 **Gene coexpression analysis of intact tissue samples reveals the core transcriptional features of di-**  
140 **verse CNS cell classes**

141 Variation among intact tissue samples can also reveal transcriptional features of less abundant cell clas-  
142 ses in the human CNS. Following the general strategy outlined in **Fig. 1**, we calculated genome-wide  
143 expression fidelity for human cholinergic neurons, midbrain dopaminergic neurons, endothelial cells,  
144 ependymal cells, choroid plexus cells, mural cells, oligodendrocyte progenitor cells, and Purkinje neu-  
145 rons (**Figs. 4, S4; Table S3**). This analysis correctly assigned high-fidelity scores for canonical markers  
146 of these cells. For example, choline acetyltransferase (*CHAT*), the high-affinity choline transporter  
147 (*SLC5A7*), and the vesicular acetylcholine transporter (*VACHT*) were all ranked within the top ~0.2% of  
148 all genes for cholinergic neuron expression fidelity, while claudin 5 (*CLDN5*), tyrosine kinase with im-  
149 munoglobulin like and EGF like domains 1 (*TIE1*), and platelet and endothelial cell adhesion molecule 1  
150 (*PECAMI*) were all ranked within the top ~0.3% of all genes for endothelial cell expression fidelity  
151 (**Table S3**). Comparisons with published gene sets revealed that high-fidelity genes were significantly  
152 and specifically enriched with expected markers of each cell class from multiple independent studies.  
153 Furthermore, novel markers predicted by our analysis were validated by *in situ* hybridization in the adult  
154 mouse brain<sup>30</sup> (**Figs. 4, S4**).

155

### 156 **High-fidelity genes enable predictive modeling of gene expression in transcriptomes from intact** 157 **tissue samples**

158 The reproducibility of gene coexpression modules corresponding to major cell classes (**Table S2, Fig.**  
159 **S3**) suggests that transcriptional variation among intact CNS samples can be modeled as a function of  
160 cellular abundance. We explored this topic systematically by performing multiple linear regression in 47  
161 CNS datasets with  $\geq 40$  samples to determine how much expression variation in a shared set of ~9600  
162 genes could be explained by variation in the abundance of neurons, astrocytes, oligodendrocytes, and  
163 microglia. To estimate the relative abundance of each cell class in each dataset, we summarized the ex-  
164 pression patterns of high-fidelity genes (**Fig. 5A**). To avoid circularity, we used a leave-one-out cross-  
165 validation strategy to redefine high-fidelity genes for each dataset by recalculating expression fidelity  
166 for each cell class using the remaining 46 datasets (as in **Fig. 1C-G**). Prior to modeling, each dataset was  
167 downsampled ( $n=40$ ) to facilitate comparisons of results; this process was performed iteratively to en-  
168 sure robustness (**Fig. 5A**).

169 Implementing this strategy, we obtained several important results (**Fig. 5B**). First, using only one  
170 gene (with the highest fidelity) as a surrogate for each cell class, our models explained 32.2% of total  
171 transcriptional variation averaged over all datasets and up to ~50% in some datasets (vs. ~0.1% for per-  
172 muted data). Second, increasing the number of gene surrogates/cell class (e.g. using the top 10 or top 50  
173 high-fidelity genes) provided only modest performance improvements (unless otherwise stated, subse-

174 quent models used the top 10 high-fidelity genes). Third, prediction accuracy depended strongly on  
175 technology platform ( $p < 10^{-7}$ , ANOVA) but not CNS region ( $p = 0.92$ , ANOVA). Among microarrays,  
176 older platforms fared substantially worse than newer platforms, while RNA-seq generally outperformed  
177 all microarrays.

178 Despite their simplicity, our models explained  $>50\%$  of expression variation, averaged over all  
179 datasets, for  $\sim 2000$  genes (**Fig. 5C**). Over all genes, the average amount of expression variation ex-  
180 plained by our models followed a sigmoid function (**Fig. 5C**). We benchmarked model performance  
181 against the maximal explanatory power of any 4 predictors by using PC1-4 from each dataset as covari-  
182 ates for multiple regression. On average, PC1-4 explained 49.6% of total gene expression variation over  
183 all datasets (**Fig. 5B**). Thus, modeling gene expression in the human CNS as a function of neuron, astro-  
184 cyte, oligodendrocyte, and microglia abundance achieved, on average, 72.0% of the maximal explanato-  
185 ry power for all datasets and 80.1% for RNA-seq datasets (**Fig. 5B**).

186 We reasoned that model performance for RNA-seq might exceed that for microarrays since the  
187 latter have many probes for transcripts that are unlikely to be expressed in the CNS. We therefore strati-  
188 fied genes by expression levels and examined model performance. As expected, predictive power de-  
189 creased at lower expression levels, with the sharpest decline between the 3rd and 4th quartiles (**Fig. 5D**).

190 We next explored how transcriptional variation related to variation in the abundance of individu-  
191 al cell classes, sex, and age. We found that neuronal abundance explained more transcriptional variation  
192 than glial abundance (**Fig. 5E**). After controlling for variation in the abundance of major cell classes,  
193 model performance did not substantially improve by including sex or age as covariates (**Fig. 5E**). We  
194 further explored this topic by correlating the estimated abundance of each cell class with age in 32 CNS  
195 datasets. We found that neuronal and oligodendroglial abundance were negatively correlated with age,  
196 while astrocytic and microglial abundance were positively correlated (**Fig. 5F**). These results suggest  
197 that age-related changes in gene expression in bulk CNS transcriptomes are primarily driven by age-  
198 related changes in cellular composition.

## 199 **Gene expression modeling applications**

200 The ability to predict gene expression in transcriptomes from intact CNS samples has substantial impli-  
201 cations for many areas of neurobiological inquiry. We illustrate the relevance of this approach through  
202 comparative analysis of gene expression models in disease, among CNS regions, and between species.

## 203 **Application #1: Contextualizing disease genes and modeling gene expression in pathological sam-** 204 **ples**

205 Using a curated database of results from genetic association studies<sup>31</sup>, we asked whether genes associat-  
206 ed with CNS diseases are enriched among genes primarily expressed by astrocytes, oligodendrocytes,  
207 microglia, or neurons (**Fig. 6A-B**). Clustering of select CNS diseases by enrichment p-values revealed  
208 several interesting findings. First, with the exception of ALS, genes associated with neurodegenerative  
209 disorders were most enriched among genes expressed by microglia and astrocytes. Second, genes asso-  
210 ciated with neurodevelopmental disorders, epilepsy, and psychiatric disorders were most enriched  
211 among genes expressed by astrocytes and neurons. Third, genes expressed by astrocytes consistently  
212 showed the greatest enrichment with candidate CNS disease genes.

213 Beyond broad associations between diseases and cell classes, gene expression modeling can also  
214 reveal which cell class is most likely to express a candidate disease gene. For example, we modeled  
215 gene expression for Alzheimer's diseases (AD) risk genes as a function of neuronal, oligodendroglial,  
216 astrocytic, and microglial abundance in transcriptomes from intact neurotypical adult human temporal  
217 cortex (**Fig. 6C**). Expression levels of early-onset AD risk genes *APP* and *PSEN1* were mostly ex-  
218 plained by variation in neuronal and oligodendroglial abundance, respectively. In contrast, expression  
219 levels of late-onset AD risk genes *APOE* and *TREM2* were mostly explained by variation in astrocytic  
220 and microglial abundance, respectively. These results were highly consistent across 47 CNS datasets  
221 (**Fig. 6D**).

222 Compared to control (CTRL) human brain samples, AD samples should contain fewer neurons  
223 and proportionately more glia. We tested this hypothesis by using expression patterns of high-fidelity  
224 genes to infer the relative abundance of neurons, astrocytes, microglia, and oligodendrocytes in 3 gene  
225 expression datasets from intact postmortem brain samples of CTRL and AD subjects<sup>32-34</sup>. We observed a  
226 highly significant decrease in neuronal abundance in AD in all datasets (**Figs. 6E, S5A-B**). In 2 out of 3  
227 datasets, there were significant increases in the relative abundance of astrocytes and microglia in AD,  
228 with similar trends in the third (**Figs. 6E, S5A-B**). Interestingly, there were no significant differences in  
229 oligodendrocyte abundance between CTRL and AD in any dataset (**Figs. 6E, S5A-B**). This strategy can  
230 help determine whether variable cellular composition is associated with diverse CNS disorders.

231 Because AD brain samples tend to have fewer neurons and proportionately more astro-  
232 cytes/microglia than CTRL, differential expression analysis of intact tissue samples will reveal down-  
233 regulation of neuronal transcripts and up-regulation of astrocytic/microglial transcripts. However, pre-  
234 dictive modeling can identify cell-intrinsic transcriptional differences between CTRL and AD that are  
235 independent of changes in cellular composition. This strategy is analogous to that of Kuhn et al.<sup>35</sup>, ex-  
236 cept here we use expression patterns of high-fidelity genes to estimate cellular abundance. Surprisingly,  
237 after controlling for differences in cellular composition between CTRL and AD, we identified many

238 genes that were consistently up-regulated in AD neurons (**Fig. 6F, Table S4**). These genes did not in-  
239 clude canonical AD risk genes (**Fig. S5C**), but rather genes involved in protein ubiquitination, catabo-  
240 lism, proteasome degradation, and mitochondrial function (**Fig. S5D**), suggesting efforts by AD neurons  
241 to mitigate the effects of misfolded protein aggregates. Examples are shown in **Figs. 6G, S5**.

## 242 **Application #2: Identifying transcriptional differences in major cell classes among CNS regions**

243 We recalculated expression fidelity separately for each CNS region with  $\geq 3$  datasets and performed hi-  
244 erarchical clustering for each cell class (**Fig. 7A-D**). Regional differences in expression fidelity were  
245 greatest for neurons, with a clear bifurcation between cortical/subcortical structures (**Fig. 7D-E**). In con-  
246 trast, expression fidelity for oligodendrocytes was very similar among brain regions (**Fig. 7B,E**). Com-  
247 paratively, microglia and astrocytes exhibited more regional variation in expression fidelity than oli-  
248 godendrocytes, but less than neurons (**Fig. 7A,C,E**).

249 We developed a conservative strategy to identify binary expression differences in major cell  
250 classes among human brain regions (**Fig. 7F-G, Table S5**). Using these criteria, many genes were pre-  
251 dicted to distinguish regional subpopulations of neurons (**Figs. 7H, S6**). Using the same criteria, we  
252 found no evidence for binary expression differences among regional subpopulations of microglia or oli-  
253 godendrocytes (**Fig. 7H**). However, we did predict binary expression differences among regional sub-  
254 populations of human astrocytes (**Fig. 7H-I**). For example, *CHRD11* was predicted to be expressed by  
255 astrocytes in frontal cortex and striatum, but not by astrocytes in diencephalon and midbrain (**Fig. 7I-K**).  
256 To validate this prediction, we performed single-molecule fluorescent *in situ* hybridization (FISH) for  
257 *Chrd11* and *Aldh11l1* in cortical and thalamic samples from mice. *Aldh11l1* is expressed ubiquitously by  
258 astrocytes<sup>15</sup> and was detected in mouse cortex and thalamus (**Fig. 7J-L**). Expression of *Chrd11* colocal-  
259 ized with *Aldh11l1* in mouse cortex but not thalamus (**Fig. 7L**), as predicted.

## 260 **Application #3: Identifying transcriptional differences in major CNS cell classes between species**

261 We analyzed 1346 mouse brain transcriptomes to determine genome-wide expression fidelity for astro-  
262 cytes, oligodendrocytes, microglia, and neurons (**Tables S1, S6; Fig. S7**). Over all homologous genes,  
263 expression fidelity was significantly correlated between mice and humans for each cell class, with the  
264 greatest similarity for neurons (**Fig. 8A**). We note that the strong conservation of neuronal expression  
265 fidelity relative to glia is mirrored at the protein level: high-fidelity neuronal genes are significantly less  
266 tolerant to LoF mutations than high-fidelity glial genes (**Fig. 3A-D**, black tracks). These findings may  
267 indicate that neurons are under greater evolutionary constraint than glia.



268 We applied stringent criteria and identified 50 genes predicted to be ‘on’ in human CNS cell  
269 classes and ‘off’ in the corresponding mouse CNS cell classes (**Fig. 8B, Table S7**). Only 6 genes were  
270 predicted with the opposite pattern (**Fig. 8B, Table S7**), which may reflect the smaller number of mouse  
271 transcriptomes analyzed. ~85% of predicted transcriptional differences between humans and mice were  
272 in glia (**Fig. 8B**). Because these differences could reflect an evolutionary gain of expression in one spe-  
273 cies or loss in the other, we analyzed 476 outgroup samples from chimpanzee and macaque brains (**Ta-**  
274 **ble S1**). Of the 50 genes predicted to be expressed in human but not mouse cell classes, 29 were signifi-  
275 cantly associated with the same cell class in at least one primate dataset; conversely, of the 6 genes with  
276 the opposite pattern, none was significantly associated with the same cell class in any primate dataset  
277 (**Table S7**). For example, expression variation of *MRVII* was largely explained by variation in astrocyte  
278 abundance in primates, but not mice (**Figs. 8B, S8A-B**). Conversely, expression variation of *PLA2G7*  
279 was largely explained by variation in astrocyte abundance in mice, but not primates (**Figs. 8B, S8A-B**).  
280 Single-molecule FISH in human and mouse cerebral cortex confirmed that expression of *MRVII* and  
281 *PLA2G7* is specific to human and mouse astrocytes, respectively (**Fig. S8C-D**).

282 To provide proof of concept for the ability of our analyses to deliver functional insights into the  
283 unique biology of human brains, we focused on a major unexplained cellular phenotype, which is the  
284 fact that human astrocytes are much larger than mouse astrocytes (as well as non-human primate astro-  
285 cytes)<sup>36</sup>. This phenotype has important implications for neuronal function, since the domain of one hu-  
286 man astrocyte can encompass up to ~2MM synapses vs. only ~100K synapses for one mouse astrocyte<sup>36</sup>.  
287 We reasoned that genes expressed by human but not mouse astrocytes might contribute to this pheno-  
288 type. We were particularly intrigued by peripheral myelin protein 2 (*PMP2*; **Fig. 8B**), which encodes a  
289 fatty-acid binding protein made by Schwann cells that is important for maintaining membrane lipid  
290 composition<sup>37</sup>. In the human CNS, expression of *PMP2* was extremely high (mean percentile: 96.2) and  
291 largely explained by variation in astrocyte abundance, while in the mouse CNS expression of *PMP2* was  
292 effectively absent (mean percentile: 11.2) and unrelated to variation in astrocyte abundance (**Fig. 8B-D**).  
293 Furthermore, independent RNA-seq data from human, chimpanzee, macaque, and mouse prefrontal cor-  
294 tex<sup>38</sup> revealed a monotonic increase in *PMP2* expression from mouse to human (**Fig. 8E**).

295 Immunostaining showed widespread *PMP2* in human neocortical astrocytes (**Fig. 8F**). In con-  
296 trast, *PMP2* was undetectable in mouse neocortex (**Fig. 8F**), despite robust expression by Schwann cells  
297 (**Fig. S8E**). To test whether *PMP2* could increase mouse astrocyte size *in vivo*, we delivered a viral con-  
298 struct expressing *PMP2* under an astrocyte-specific promoter to neonatal mouse brains and analyzed the  
299 morphology of transduced astrocytes after 42d (**Fig. 8G**). Forced expression of *PMP2* in mouse astro-  
300 cytes significantly increased their maximum diameter and number of primary processes (**Fig. 8H-I**). The

301 increase in maximum diameter corresponded to an increase in mouse astrocyte volume of ~50% (assum-  
302 ing sphericity). To further validate this finding, we repeated the experiment with a different viral con-  
303 struct and obtained nearly identical results (**Fig. S8F**). To our knowledge, these data provide the first  
304 molecular explanation for morphological differences between human and mouse astrocytes. More gen-  
305 erally, our findings illustrate how variation among intact tissue samples can predict cell-class-specific  
306 transcriptional features with important functional implications for human neurobiology.

## 307 **DISCUSSION**

308 We have described a novel, ‘top-down’ approach to reveal the core transcriptional features of cellular  
309 identity via integrative gene coexpression analysis of intact tissue samples. Compared to ‘bottom-up’  
310 methods such as FACS, IP, and SC/SN, the main advantages of our approach are as follows: i) elimina-  
311 tion of the need for fresh tissue; ii) applicability to huge amounts of existing data; iii) elimination of  
312 technical variability caused by tissue dissociation and cDNA preamplification; iv) elimination of sam-  
313 pling bias associated with cell/nucleus capture; and v) ability to derive highly robust inferences about  
314 the core transcriptional features of cellular identity based on aggregate analysis of billions of cells.

315 Our approach also has important limitations. False-positive associations can result from technical  
316 factors such as batch effects or biological factors such as cellular collinearity. For example, we consist-  
317 ently observed that genes with high expression fidelity for oligodendrocytes had higher expression fidel-  
318 ity for microglia (and vice versa) than they did for astrocytes or neurons. Because oligodendrocytes and  
319 microglia are more abundant in white matter than gray matter<sup>39</sup>, variation in the ratio of white matter to  
320 gray matter in CNS samples drives covariation in the abundance of these cell classes and the genes that  
321 they express. False-negative associations can result from technical factors such as limitations in dynamic  
322 range/transcriptome coverage or probe failures, as well as biological factors such as alternative splicing.  
323 Notwithstanding these limitations, the genes with the highest expression fidelity for major CNS cell  
324 classes are already remarkably stable.

325 It is interesting to consider the ability of our approach to detect transcriptional signatures of less  
326 abundant cell classes (e.g. **Figs. 4, S4**). The ability to discern a gene coexpression signature of a cell  
327 class in transcriptomes from intact tissue samples depends on many factors, including its representation,  
328 the uniqueness and abundance of its transcripts, its stoichiometry with other cell types, the technology  
329 platform, the algorithmic approach, and the sampling strategy<sup>8</sup>. Some of these factors can be optimized  
330 to improve sensitivity. Ultimately, however, we envision future studies that combine the benefits of top-  
331 down and bottom-up strategies to fully deconstruct the transcriptional architecture of biological systems.



332 Our estimates of gene expression fidelity for major cell classes were highly robust to the choice  
333 of gene set used for enrichment analysis, but more so for glia than neurons. This result indicates that  
334 neuronal diversity may require additional strategies to optimize estimates of neuronal expression fidelity,  
335 particularly on a regional basis. For example, the neuronal gene sets used in this study do not capture  
336 the transcriptional profile of cerebellar granule neurons, which is highly distinct<sup>14, 22, 26</sup>. To better account  
337 for neuronal diversity, future studies may utilize additional neuron subtype-specific or composite  
338 gene sets for enrichment analyses.

339 Our results suggest that the functional identity of a cell class can be conceived as a vector of  
340 genes ranked by the fidelity with which they are expressed in that cell class relative to all other cells in  
341 the biological system of interest. An advantage of this framing is that it is inherently context-dependent.  
342 Beyond revealing novel biomarkers and cellular phenotypes, such definitions can provide ‘molecular  
343 rulers’ for measuring the validity of human cells derived *in vitro* for disease modeling and cell replacement  
344 therapies. In addition, these definitions can be tested in *de novo* CNS transcriptomes for their ability  
345 to predict gene expression levels through mathematical modeling.

346 Multivariate analyses of CNS transcriptomes often use module detection/clustering methods or  
347 projection methods such as principal component analysis. Although these methods have produced many  
348 important insights, they are inherently descriptive and do not lend themselves easily to comparisons  
349 among independent datasets. Because the building block of any biological system is the cell, and cells  
350 are distinguished by the genes that they express, an alternative approach is to model expression levels of  
351 individual genes as a function of variation in cellular composition. We have shown how expression patterns  
352 of high-fidelity genes can be used as covariates in multiple linear regression models for this purpose.  
353 The resulting models are grounded in biology, easily compared among independent datasets, and  
354 capable of extracting cell-class-specific insights from intact tissue samples. Using this approach, we explored  
355 how predictive models of gene expression in transcriptomes from intact CNS samples can inform  
356 studies of aging, disease genes, pathological samples, regional heterogeneity, and species differences.  
357 We elaborate upon our findings in the Supplementary Discussion.

358 The analyses presented in this study are based on a simple idea: variation in cellular composition  
359 among intact tissue samples will drive covariation of transcripts that are uniquely or predominantly expressed  
360 in specific kinds of cells. Although we have focused here on gene expression, our approach can  
361 also be applied to other types of molecular data, thereby offering a generalizable strategy for determining  
362 the core molecular features of cellular identity in intact biological systems.

## **REFERENCES**

- 363 1. Darmanis, S., *et al.* A survey of human brain transcriptome diversity at the single cell level.  
364 *PNAS* **112**, 7285-7290 (2015).
- 365 2. Paul, G., *et al.* The adult human brain harbors multipotent perivascular mesenchymal stem  
366 cells. *PLoS One* **7**, e35577 (2012).
- 367 3. Zhang, Y., *et al.* Purification and Characterization of Progenitor and Mature Human  
368 Astrocytes Reveals Transcriptional and Functional Differences with Mouse. *Neuron* **89**, 37-53  
369 (2016).
- 370 4. Habib, N., *et al.* Massively parallel single-nucleus RNA-seq with DroNc-seq. *Nat Methods* **14**,  
371 955-958 (2017).
- 372 5. Lake, B.B., *et al.* Neuronal subtypes and diversity revealed by single-nucleus RNA  
373 sequencing of the human brain. *Science* **352**, 1586-1590 (2016).
- 374 6. Liu, S. & Trapnell, C. Single-cell transcriptome sequencing: recent advances and remaining  
375 challenges. *F1000Res* **5** (2016).
- 376 7. Azevedo, F.A., *et al.* Equal numbers of neuronal and nonneuronal cells make the human  
377 brain an isometrically scaled-up primate brain. *J. Comp. Neurol.* **513**, 532-541 (2009).
- 378 8. Oldham, M.C. Transcriptomics: from differential expression to coexpression. in *The OMICs:*  
379 *applications in neurosciences.* (ed. G. Coppola) 85-113 (Oxford, 2014).
- 380 9. Oldham, M.C., *et al.* Functional organization of the transcriptome in human brain. *Nat.*  
381 *Neurosci.* **11**, 1271-1282 (2008).
- 382 10. Fertuzinhos, S., *et al.* Laminar and temporal expression dynamics of coding and noncoding  
383 RNAs in the mouse neocortex. *Cell Rep* **6**, 938-950 (2014).
- 384 11. Ponomarev, I., Rau, V., Eger, E.I., Harris, R.A. & Fanselow, M.S. Amygdala transcriptome and  
385 cellular mechanisms underlying stress-enhanced fear learning in a rat model of posttraumatic  
386 stress disorder. *Neuropsychopharmacology* **35**, 1402-1411 (2010).
- 387 12. Hilliard, A.T., Miller, J.E., Fraley, E.R., Horvath, S. & White, S.A. Molecular microcircuitry  
388 underlies functional specification in a basal ganglia circuit dedicated to vocal learning. *Neuron* **73**,  
389 537-552 (2012).
- 390 13. Bakken, T.E., *et al.* A comprehensive transcriptional map of primate brain development.  
391 *Nature* **535**, 367-375 (2016).
- 392 14. Hawrylycz, M., *et al.* Canonical genetic signatures of the adult human brain. *Nat. Neurosci.*  
393 **18**, 1832-1844 (2015).
- 394 15. Cahoy, J.D., *et al.* A transcriptome database for astrocytes, neurons, and oligodendrocytes: a  
395 new resource for understanding brain development and function. *J. Neurosci.* **28**, 264-278 (2008).

- 396 16. Hickman, S.E., *et al.* The microglial sensome revealed by direct RNA sequencing. *Nat.*  
397 *Neurosci.* **16**, 1896-1905 (2013).
- 398 17. Horvath, S. & Dong, J. Geometric interpretation of gene coexpression network analysis.  
399 *PLoS Comput. Biol.* **4**, e1000117 (2008).
- 400 18. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network  
401 analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 402 19. Tasic, B., *et al.* Adult mouse cortical cell taxonomy revealed by single cell transcriptomics.  
403 *Nat. Neurosci.* **19**, 335-346 (2016).
- 404 20. Hodges, A., *et al.* Regional and cellular gene expression changes in human Huntington's  
405 disease brain. *Hum. Mol. Genet.* **15**, 965-977 (2006).
- 406 21. Berchtold, N.C., *et al.* Gene expression changes in the course of normal brain aging are  
407 sexually dimorphic. *PNAS* **105**, 15605-15610 (2008).
- 408 22. Kang, H.J., *et al.* Spatio-temporal transcriptome of the human brain. *Nature* **478**, 483-489  
409 (2011).
- 410 23. Hernandez, D.G., *et al.* Integration of GWAS SNPs and tissue specific expression profiling  
411 reveal discrete eQTLs for human traits in blood and brain. *Neurobiol. Dis.* **47**, 20-28 (2012).
- 412 24. Li, J.Z., *et al.* Circadian patterns of gene expression in the human brain and disruption in  
413 major depressive disorder. *PNAS* **110**, 9950-9955 (2013).
- 414 25. Ramasamy, A., *et al.* Genetic variability in the regulation of gene expression in ten regions of  
415 the human brain. *Nat. Neurosci.* **17**, 1418-1428 (2014).
- 416 26. GTExConsortium. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis:  
417 multitissue gene regulation in humans. *Science* **348**, 648-660 (2015).
- 418 27. Lek, M., *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**,  
419 285-291 (2016).
- 420 28. Binder, J.X., *et al.* COMPARTMENTS: unification and visualization of protein subcellular  
421 localization evidence. *Database* **Feb 25**, bau012 (2014).
- 422 29. Szklarczyk, D., *et al.* STRING v10: protein-protein interaction networks, integrated over the  
423 tree of life. *Nucleic Acids Res.* **43**, D447-452 (2015).
- 424 30. Lein, E.S., *et al.* Genome-wide atlas of gene expression in the adult mouse brain. *Nature* **445**,  
425 168-176 (2007).
- 426 31. Yu, W., Clyne, M., Khoury, M.J. & Gwinn, M. Phenopedia and Genopedia: disease-centered  
427 and gene-centered views of the evolving knowledge of human genetic associations. *Bioinformatics*  
428 **26**, 145-146 (2010).

- 429 32. Cribbs, D.H., *et al.* Extensive innate immune gene activation accompanies brain aging,  
430 increasing vulnerability to cognitive decline and neurodegeneration: a microarray study. *J.*  
431 *Neuroinflammation* **9**, 179 (2012).
- 432 33. Zhang, B., *et al.* Integrated systems approach identifies genetic nodes and networks in late-  
433 onset Alzheimer's disease. *Cell* **153**, 707-720 (2013).
- 434 34. Hokama, M., *et al.* Altered expression of diabetes-related genes in Alzheimer's disease  
435 brains: the Hisayama study. *Cereb. Cortex* **24**, 2476-2488 (2014).
- 436 35. Kuhn, A., Thu, D., Waldvogel, H.J., Faull, R.L. & Luthi-Carter, R. Population-specific  
437 expression analysis (PSEA) reveals molecular changes in diseased brain. *Nat Methods* **8**, 945-947  
438 (2011).
- 439 36. Oberheim, N.A., Goldman, S.A. & Nedergaard, M. Heterogeneity of astrocytic form and  
440 function. *Methods Mol. Biol.* **814**, 23-45 (2012).
- 441 37. Zenker, J., *et al.* A role of peripheral myelin protein 2 in lipid homeostasis of myelinating  
442 Schwann cells. *Glia* **62**, 1502-1512 (2014).
- 443 38. Bozek, K., *et al.* Exceptional evolutionary divergence of human muscle and brain  
444 metabolomes parallels human cognitive and physical uniqueness. *PLoS Biol.* **12**, e1001871 (2014).
- 445 39. Mittelbronn, M., Dietz, K., Schluesener, H.J. & Meyermann, R. Local distribution of microglia  
446 in the normal adult human central nervous system differs by up to one order of magnitude. *Acta*  
447 *Neuropathol.* **101**, 249-255 (2001).
- 448 40. Storey, J.D. & Tibshirani, R. Statistical significance for genomewide studies. *PNAS* **100**,  
449 9440-9445 (2003).
- 450 41. Zhang, Y., *et al.* An RNA-sequencing transcriptome and splicing database of glia, neurons,  
451 and vascular cells of the cerebral cortex. *J. Neurosci.* **34**, 11929-11947 (2014).

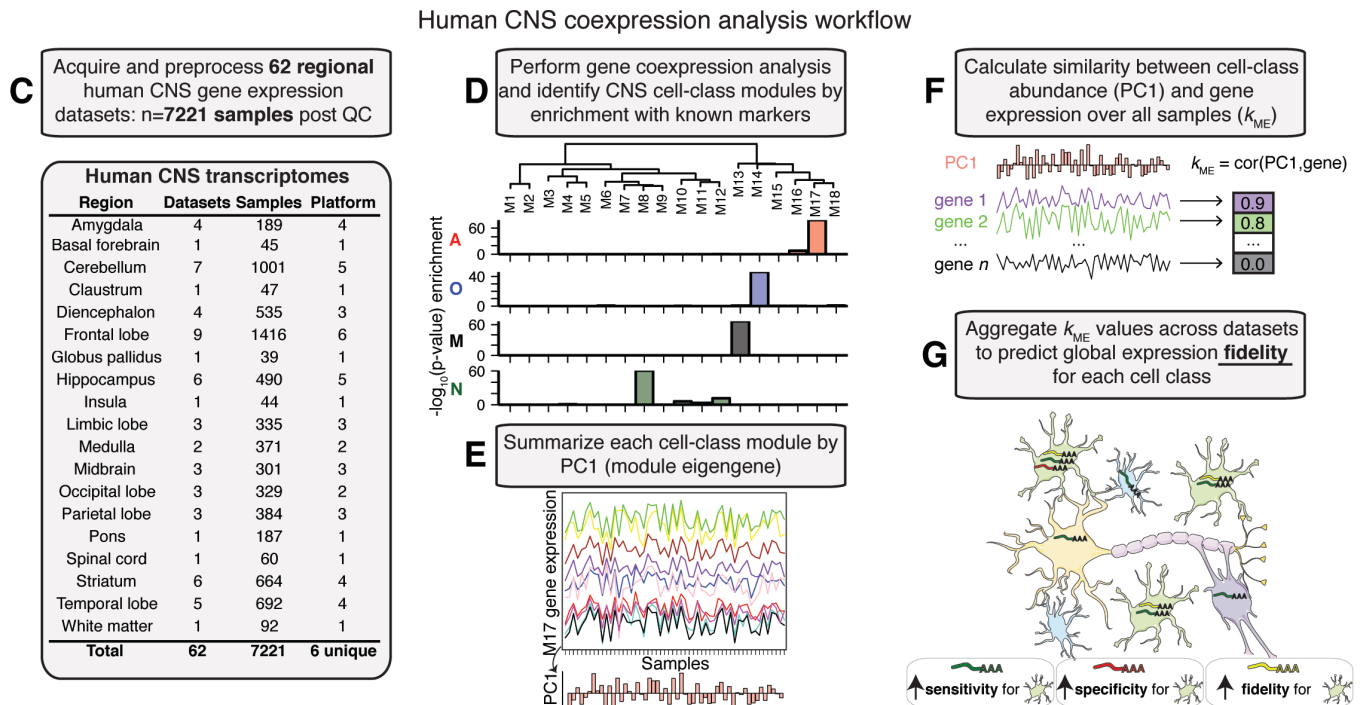
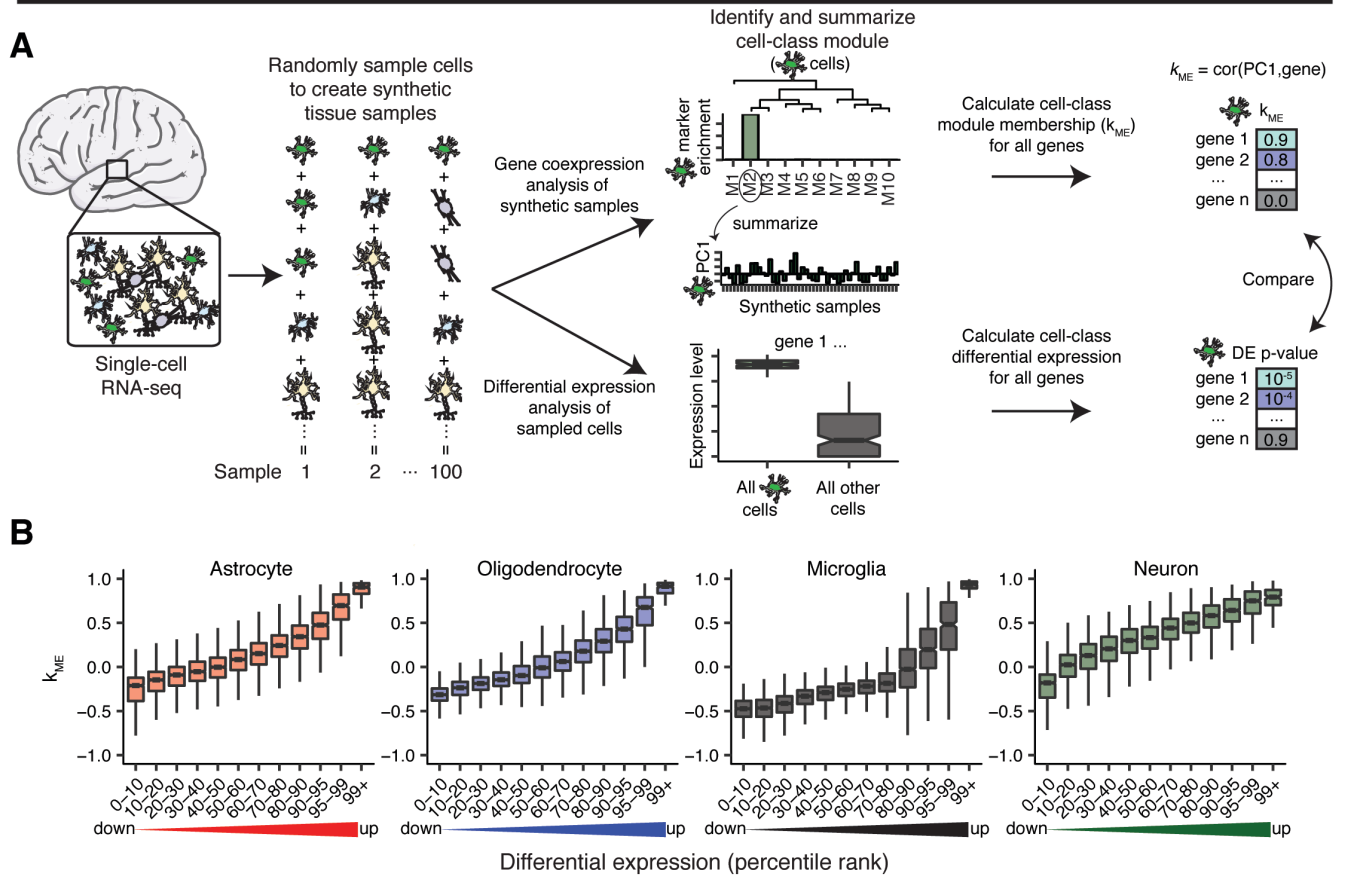
452 **ACKNOWLEDGMENTS**

453 We are grateful to Brad Dispensa, Joe Hesse, Dirk Kleinhesselink, and Jason Jed for technical support.  
454 We thank Annette Molinaro for statistical consultations, David Rowitch for astrocyte discussions, and  
455 Eric Huang and Mercedes Paredes for human brain samples. Due to space limitations, we apologize that  
456 many relevant and important publications are not cited. This work was supported by the UCSF Program  
457 for Breakthrough Biomedical Research (MCO), which is funded in part by the Sandler Foundation, a  
458 Scholar Award from the UCSF Weill Institute for Neurosciences (MCO), and NIMH R01MH113896  
459 (MCO).

460 **AUTHOR CONTRIBUTIONS**

461 KWK and MCO conceived and designed the analytical strategies and wrote the manuscript. KWK per-  
462 formed most data analyses and histological experiments. KWK and HO performed *PMP2* expression  
463 experiments under supervision from AVM.

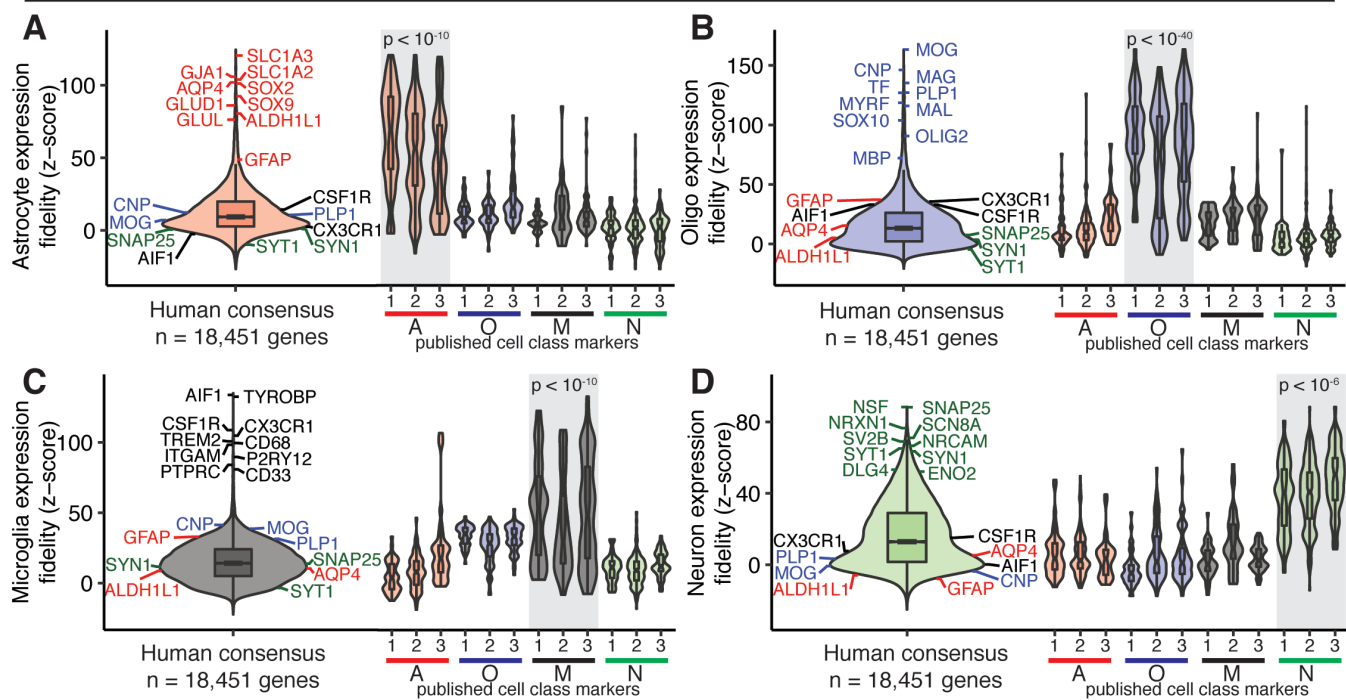
# Figure 1



464 **Fig. 1 | Rationale and workflow.** **A)** Left: Single-cell RNA-seq data from adult human brain samples<sup>1</sup>  
465 were randomly aggregated to create 100 synthetic tissue samples. Right (top): Unsupervised gene coex-  
466 pression analysis of synthetic samples revealed CNS cell-class modules that were highly enriched with  
467 markers of astrocytes, oligodendrocytes, microglia, or neurons. Cell-class module membership strength  
468 ( $k_{ME}$ ) was calculated for all genes. Right (bottom): Using the same cells that were selected to create syn-  
469 thetic samples, single-cell differential expression analysis was performed for all genes with respect to  
470 each cell class. **B)**  $k_{ME}$  values for synthetic cell-class modules accurately predicted the results of differ-  
471 ential expression analysis for each cell class (n=10 synthetic datasets; ‘up’ / ‘down’ denote up- and  
472 down-regulated genes for each cell class). **C)** 62 datasets representing diverse CNS regions and technol-  
473 ogy platforms were acquired and preprocessed (**Table S1**). **D)** Unsupervised gene coexpression analysis  
474 was performed for each dataset to identify modules of genes with similar expression patterns. Each  
475 module was summarized by PC1 (module eigengene). **E)** Published markers of each cell class were  
476 cross-referenced with all modules (Fisher’s exact test; **Table S2**). **F)** Cell-class module eigengenes were  
477 used to calculate the similarity between cellular abundance and genome-wide expression patterns ( $k_{ME}$ )  
478 over all samples. **G)** Genome-wide  $k_{ME}$  values for significant cell-class modules were combined to yield  
479 a global measure of expression ‘fidelity’ for each gene with respect to each cell class. Schematic: A gene  
480 has high fidelity for a cell class if its expression is sensitive (it is consistently expressed by members of  
481 that cell class) and specific (it is not expressed by members of other cell classes).

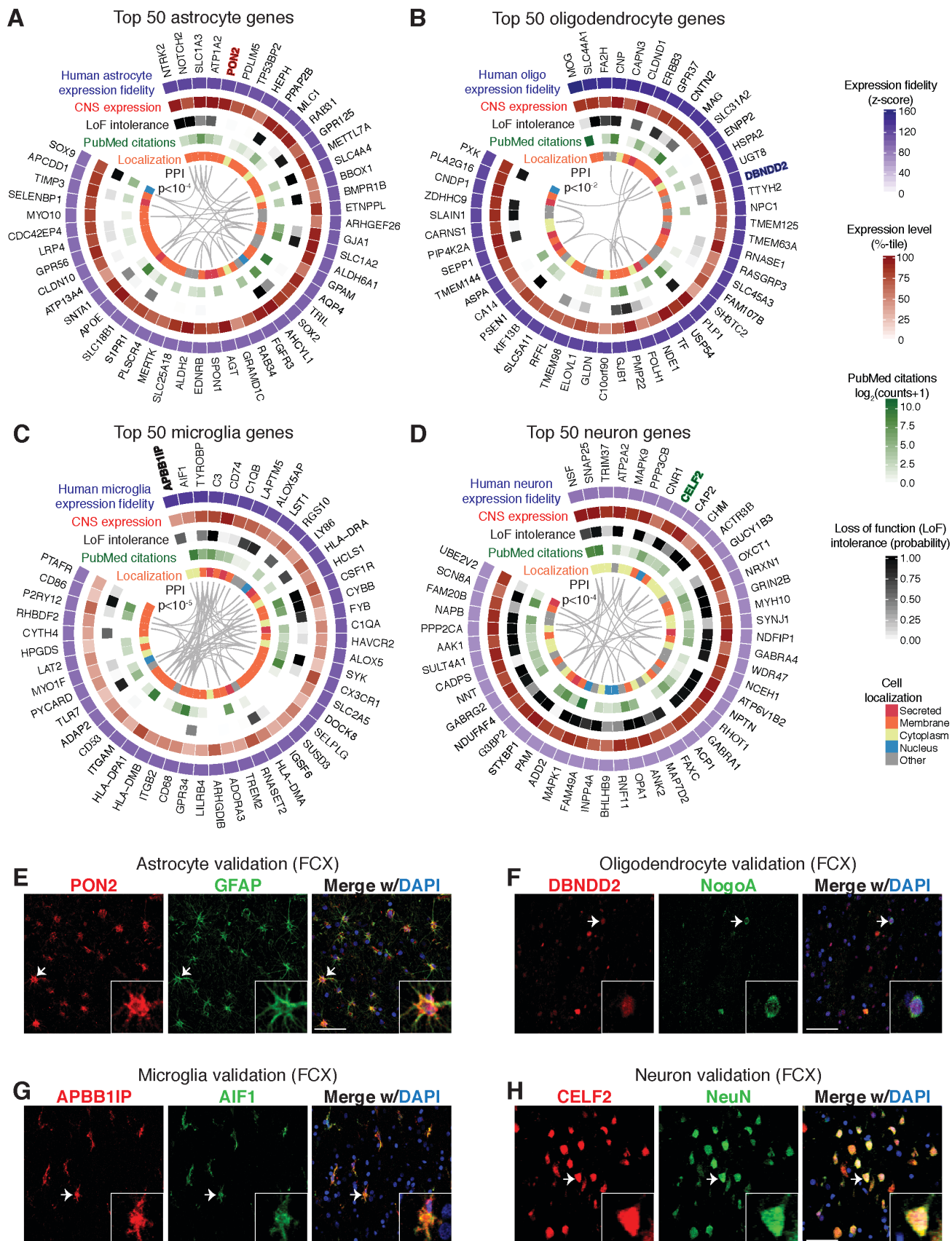


## Figure 2



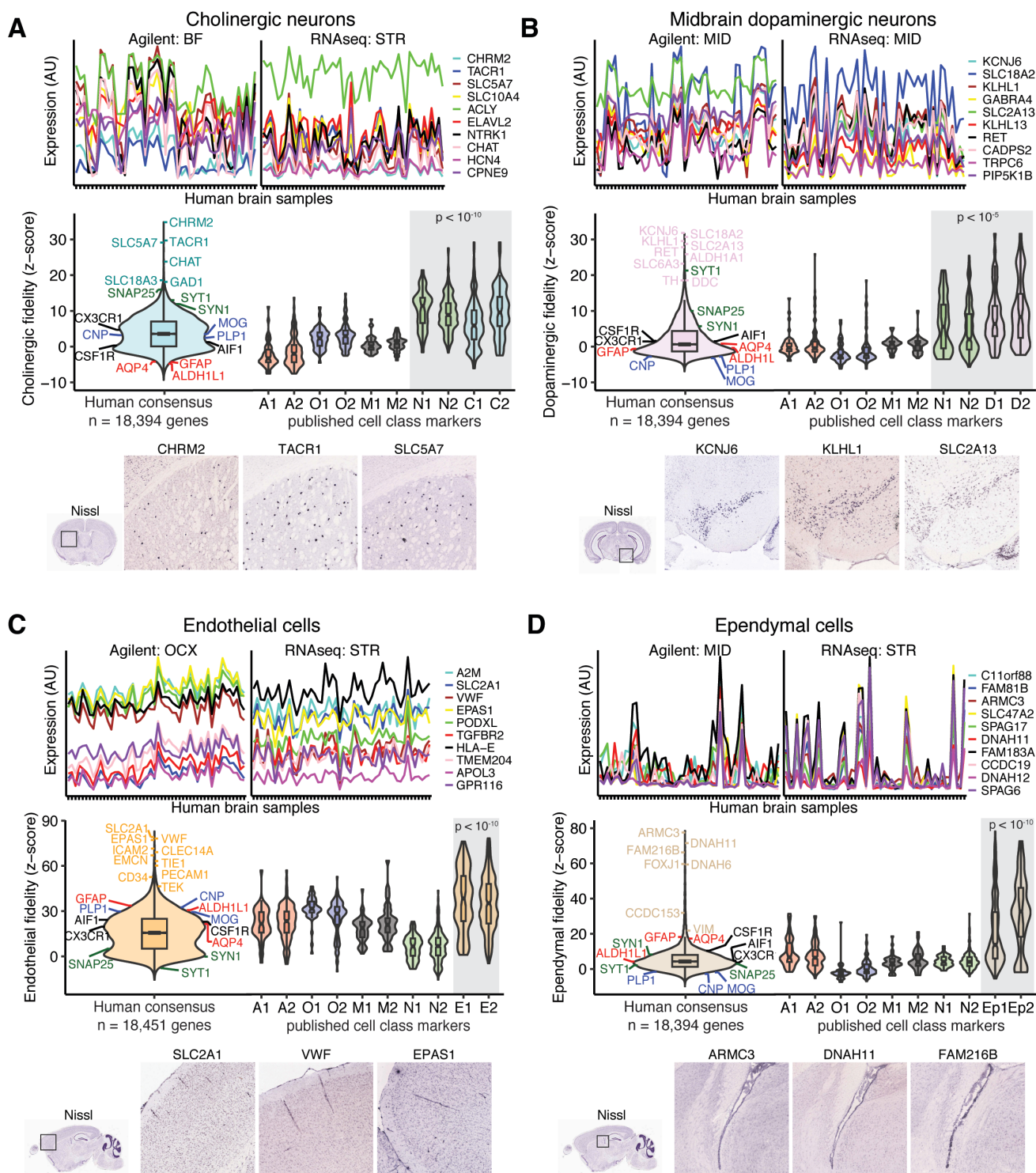
482 **Fig. 2 | Integrative gene coexpression analysis of intact CNS transcriptomes reveals consensus**  
483 **transcriptional profiles of human astrocytes, oligodendrocytes, microglia, and neurons. A-D) Left:**  
484 consensus gene expression fidelity distributions for human astrocytes (A), oligodendrocytes (O), micro-  
485 glia (M), and neurons (N). Canonical markers are labeled in red (A), blue (O), black (M), and green (N).  
486 Right: gene expression fidelity distributions for published sets of markers (A1-3, O1-3, M1-3, N1-3;  
487 Methods) were cross-referenced with high-fidelity genes (z-score >50). Gray shading: significant en-  
488 richment (Fisher's exact test).

## Figure 3



489 **Fig. 3 | The core transcriptional identities of human astrocytes, oligodendrocytes, microglia, and**  
490 **neurons include known and novel biomarkers. A-D)** The top 50 genes ranked by consensus expres-  
491 sion fidelity for astrocytes, oligodendrocytes, microglia, or neurons. Expression levels: averages of mean  
492 percentile ranks for all datasets. Mutation intolerance: ExAC<sup>27</sup>. PubMed citations: queries for gene + cell  
493 class (e.g. gene symbol + 'astrocyte'). Cellular localization: COMPARTMENTS<sup>28</sup>. Predicted protein-  
494 protein interactions (PPI): STRING<sup>29</sup>. Link=combined score >350. P-values for observed # links based  
495 on 100K random samples of 50 genes. **E-H)** Novel markers of human astrocytes (PON2), oligodendro-  
496 cytes (DBNDD2), microglia (APBB1IP), and neurons (CELF2) in adult human dorsolateral prefrontal  
497 cortex (DLPFC; **E**: L5/6; **F,G**: white matter; **H**: L2/3). Arrowhead: cell in inset. Scale bar: 50µm.

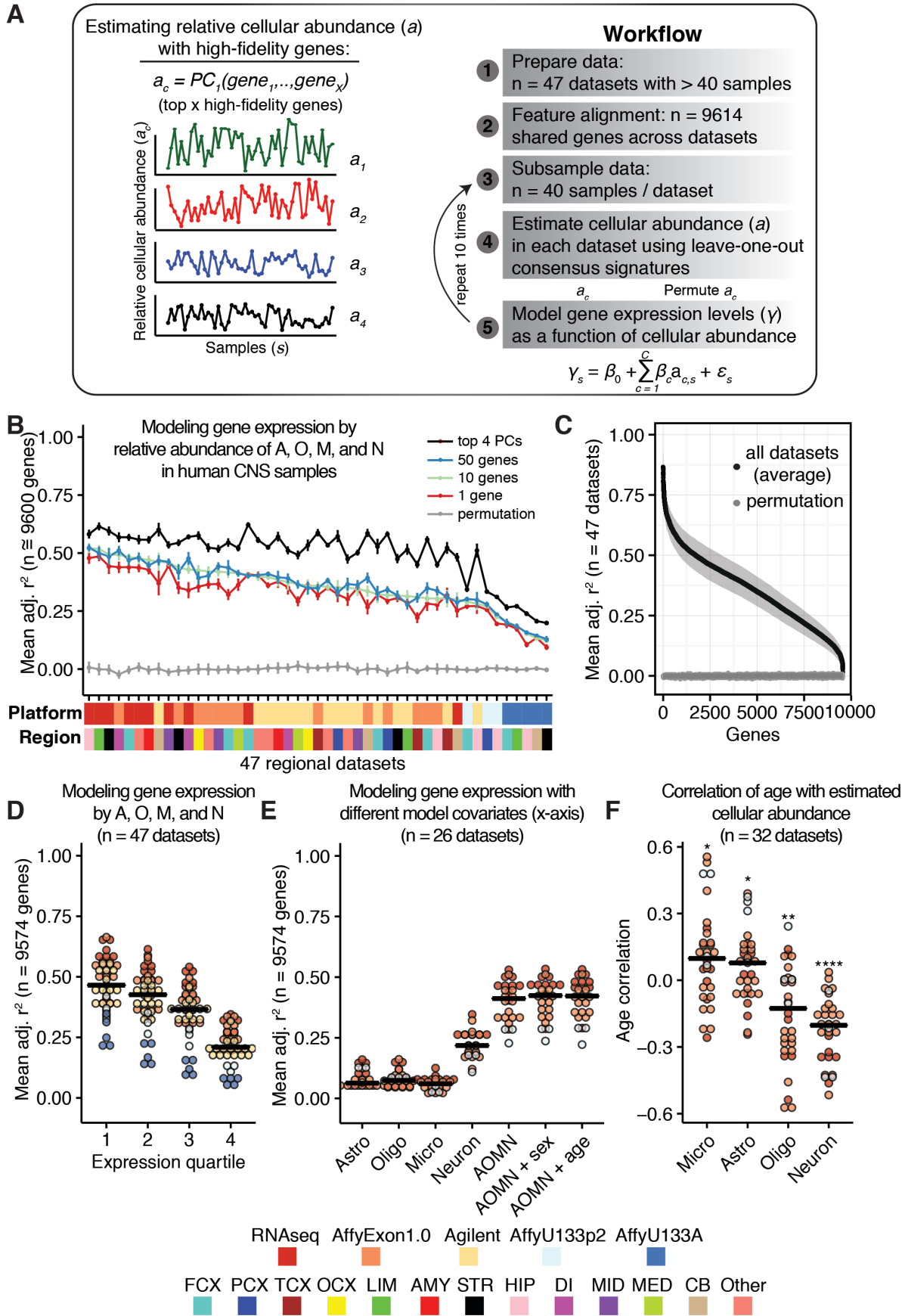
## Figure 4



498 **Fig. 4 | Variation among intact tissue samples reveals transcriptional signatures of human cholin-**  
499 **ergic neurons, midbrain dopaminergic neurons, endothelial cells, and ependymal cells. A-D)** Top:  
500 high-fidelity genes for each cell class (top 10 are shown) are consistently coexpressed in independent  
501 datasets. Middle: consensus gene expression fidelity distributions for each cell class with canonical  
502 markers of major cell classes labeled in green (neurons), red (astrocytes), blue (oligodendrocytes), and  
503 black (microglia). Gene expression fidelity distributions for published sets of markers (A1, A2, O1, O2,  
504 M1, M2, N1, N2, C1, C2, D1, D2, E1, E2, Ep1, Ep2; Methods) were cross-referenced with high-fidelity  
505 genes (top 3 percentile). Gray shading: significant enrichment (Fisher's exact test). Bottom: mouse *in*  
506 *situ* hybridization data<sup>30</sup> for high-fidelity genes in dorsal striatum (**A**), ventral midbrain (**B**), cortex (**C**),  
507 and lateral ventricle (**D**).



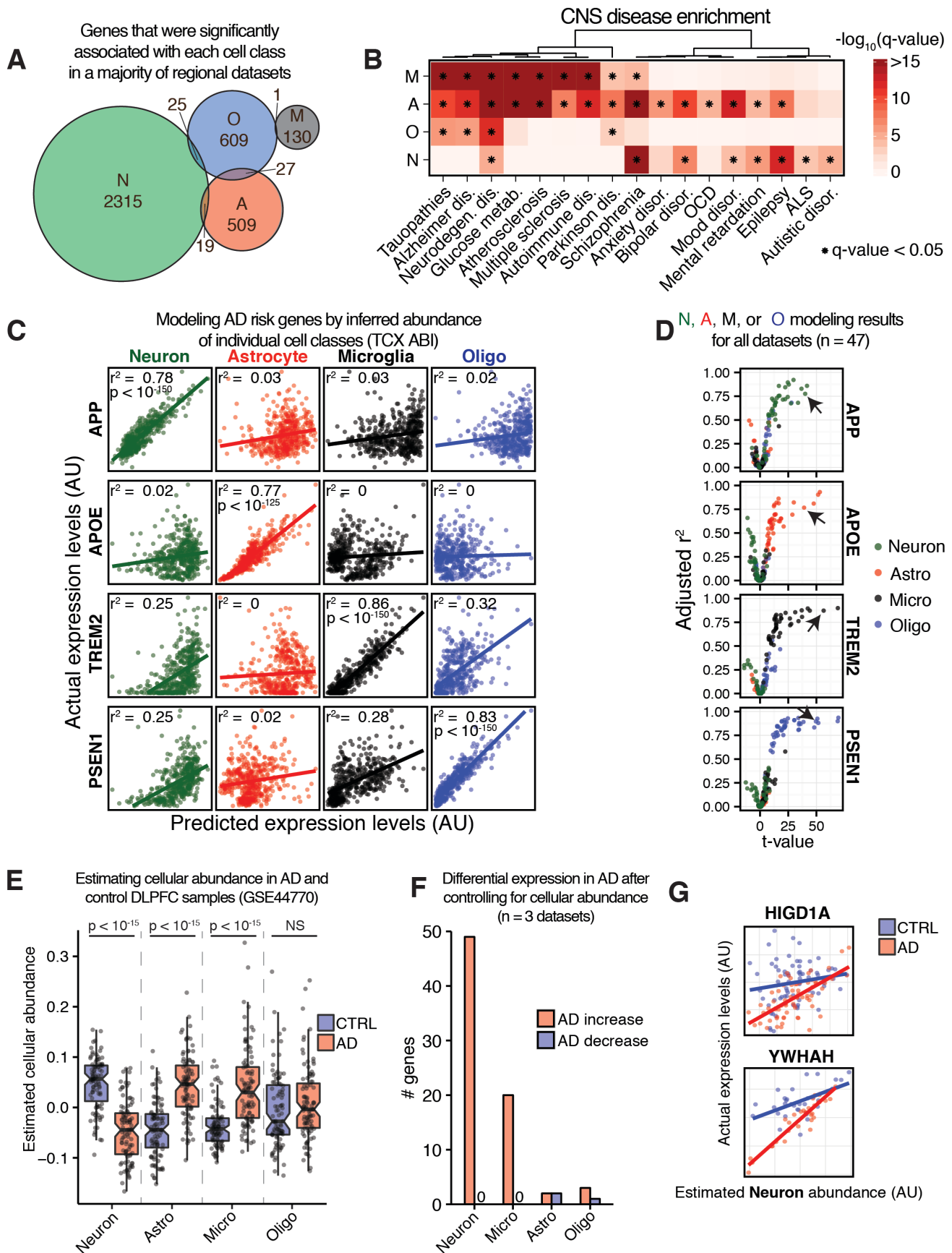
## Figure 5



508 **Fig. 5 | Variation in cellular abundance predicts gene expression in transcriptomes from intact**  
509 **CNS samples. A)** Strategy for modeling gene expression in intact human CNS samples as a function of  
510 inferred cellular abundance. **B)** Total % variance explained (mean adj.  $r^2$ ) for ~9600 genes whose ex-  
511 pression levels were modeled as a function of inferred astrocyte, oligodendrocyte, microglia, and neuron  
512 abundance in 47 datasets (subset to 40 samples; values are mean  $\pm$  2 s.e.m., 10 iterations). **C)** Mean  
513 adj.  $r^2$  values for individual genes from **(B)** over 47 datasets. Grey envelope: loess smoothed C.I. ( $\pm$  2  
514 s.e.m., 10 iterations). **D)** Mean adj.  $r^2$  values for genes from **(B)** grouped by mean expression quartiles  
515 (each point is one dataset). **E)** Mean adj.  $r^2$  values for 7 different models (restricted to datasets w/ sex  
516 and age: GSE46706, GTE<sub>x</sub>, GSE11882, GSE25219). **F)** Pearson correlation of inferred cellular abun-  
517 dance with age (\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*\*  $p < 0.0001$ , one-sample Wilcoxon signed-rank test). Horizontal  
518 bars **(D-F)**: median; points colored by technology platform.

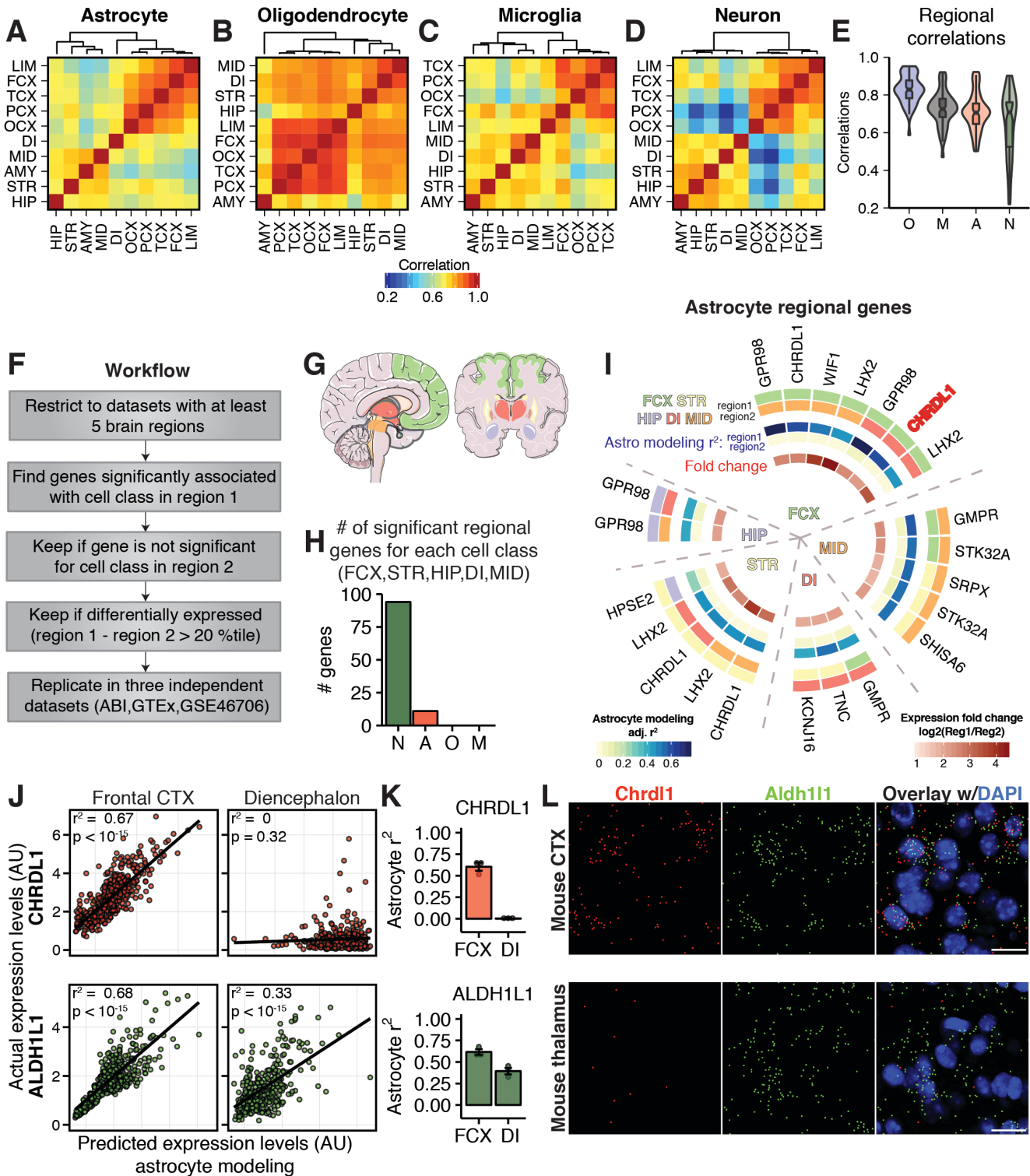


## Figure 6



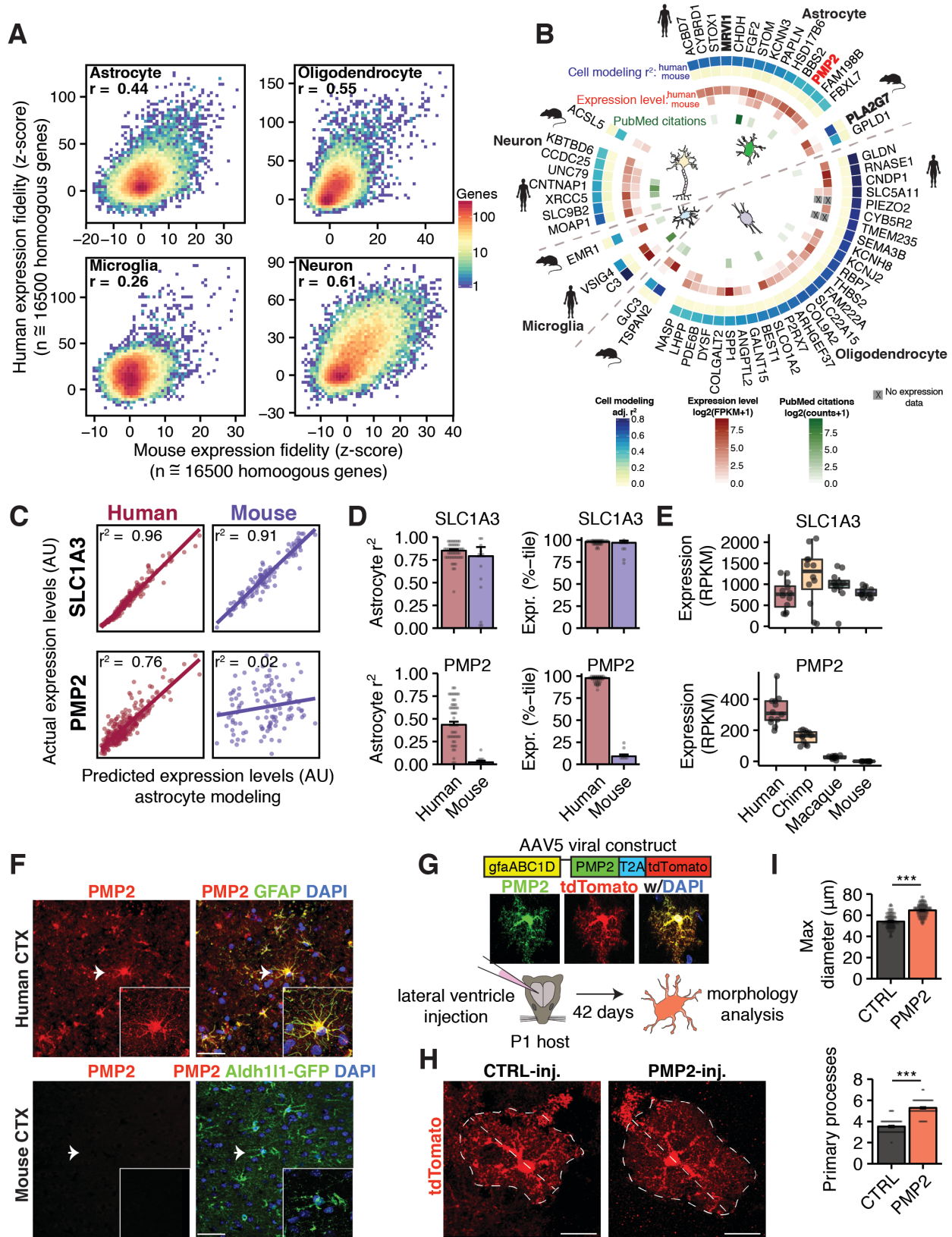
519 **Fig. 6 | Gene expression modeling offers new avenues for studying human CNS diseases.** **A)** # genes  
520 associated with each human cell class ( $p < 8.37 \times 10^{-9}$ , Bonferroni correction for # gene models). **B)** En-  
521 richment analysis (Fisher's exact test) of genes from **(A)** with human CNS disease genes from Pheno-  
522 pedia<sup>31</sup>. FDR-adjusted p-values (q-values) are shown<sup>40</sup>. **C)** Modeling results in human temporal cortex  
523 (TCX ABI; i.e. 1 dataset) for 4 AD risk genes. **D)** Modeling results for genes from **(C)** in 47 datasets  
524 ( $\geq 40$  samples). **E)** Top 10 high-fidelity genes were used to estimate the relative abundance of neurons,  
525 astrocytes, microglia, and oligodendrocytes in DLPFC from control (CTRL) and AD<sup>33</sup> (**Fig. 4A**). P-  
526 values: Wilcoxon rank-sum test. **F)** Gene expression modeling in 3 datasets<sup>32-34</sup> reveals consistent cell-  
527 class-specific expression changes in AD after controlling for differences in cellular abundance ( $p < 0.05$   
528 based on 1000 permutations of sample labels). **G)** Examples of two genes that are up-regulated in AD  
529 neurons (top<sup>33</sup>; bottom<sup>34</sup>).

# Figure 7



530 **Fig. 7 | Regional expression fidelity and predictive modeling reveal astrocyte heterogeneity in the**  
531 **human brain. (A-D)** Hierarchical clustering of human brain regions (excluding cerebellum) based on  
532 Pearson correlations among regional expression fidelity for each cell class (n=18451 genes,  $\geq 3$  da-  
533 taset/region). **E)** Distributions of correlations in **(A-D)**. **F)** Workflow to predict regional expression dif-  
534 ferences in specific cell classes. Significance threshold:  $p < 2.67 \times 10^{-8}$  (Bonferroni correction for total # of  
535 gene models). **G)** Analyzed brain regions: frontal cortex (FCX), striatum (STR), hippocampus (HIP),  
536 diencephalon (DI), and midbrain (MID). **H)** Total # of region-specific genes conservatively predicted  
537 for each cell class. **I)** Genes predicted to be expressed by human astrocytes in restricted brain regions. **J)**  
538 Modeling of *CHRDLL1* and *ALDH1L1* (+ control) as a function of inferred astrocyte abundance in exam-  
539 ple datasets (FCX/DI from ABI). **K)** Modeling results for same genes in 3 datasets (ABI, GTEx, and  
540 GSE46706). **L)** Single-molecule FISH of *Chrdll1* and *Aldh1l1* in mature mouse brain (P30). Scale bar:  
541 20 $\mu$ m.

## Figure 8



542 **Fig. 8 | Gene expression modeling identifies cell-class-specific transcriptional differences between**  
543 **humans and mice. A)** Comparison of gene expression fidelity in humans and mice for each cell class.  
544 **B)** Predicted cell-class-specific transcriptional differences between humans and mice. Expression levels  
545 are from independent datasets<sup>3, 41</sup> that were not used to predict species differences. PubMed citations  
546 obtained as in **Fig. 3. C)** Example modeling results in humans (Hs.PCX.ABI) and mice (Ms.GSE64398)  
547 (**Table S1**). *SLCIA3* is expressed by astrocytes in both species and *PMP2* by astrocytes in humans but  
548 not mice. **D)** Astrocyte modeling results and mean expression percentiles for genes in **(C)** from all da-  
549 taset. Error bars: s.e.m. **E)** *SLCIA3* and *PMP2* expression in human, chimpanzee, macaque, and mouse  
550 prefrontal cortex<sup>38</sup>. **F)** Immunostaining for PMP2 in adult human DLPFC and P42 mouse neocortex. Ar-  
551 rowheads: cells in insets. Scale bar: 40 $\mu$ m. **G)** Experimental strategy for studying *PMP2* effects on  
552 mouse astrocytes. **H)** Representative examples of CTRL and *PMP2*-infected astrocytes in mouse neo-  
553 cortex. Dashed lines outline cell and max diameter through nucleus. Scale bar: 20 $\mu$ m. **I)** Quantification  
554 of max diameter and # of primary processes in CTRL and *PMP2*-infected astrocytes. n=4 animals/group,  
555 n>15 astrocytes/animal, mean  $\pm$  s.e.m., Welch's *t*-test, \*\*\* p<0.001.