

This item is the archived peer-reviewed author-version of:

Integrative transcriptomic analysis of the amyotrophic lateral sclerosis spinal cord implicates glial activation and suggests new risk genes

Reference:

Humphrey Jack, Venkatesh Sanan, Hasan Rahat, Herb Jake T., de Lopes Katia Paiva, Küçükali Fahri, Byrska-Bishop Marta, Evani Uday S., Narzisi Giuseppe, Fagegaltier Delphine,- Integrative transcriptomic analysis of the amyotrophic lateral sclerosis spinal cord implicates glial activation and suggests new risk genes

Nature neuroscience - ISSN 1546-1726 - Berlin, Nature portfolio, 26(2023), p. 150-162

Full text (Publisher's DOI): <https://doi.org/10.1038/S41593-022-01205-3>

To cite this reference: <https://hdl.handle.net/10067/1928170151162165141>

8 **Integrative transcriptomic analysis of the amyotrophic lateral sclerosis**
9 **spinal cord implicates glial activation and suggests new risk genes**

10 Jack Humphrey^{1,2,3,4*}, Sanan Venkatesh^{1,3,5}, Rahat Hasan^{1,2,3,4}, Jake T. Herb⁶, Katia de Paiva
11 Lopes^{1,2,3,4}, Fahri Küçükali^{7,8}, Marta Byrska-Bishop⁹, Uday S. Evani⁹, Giuseppe Narzisi⁹, Delphine
12 Fagegaltier^{9,10}, NYGC ALS Consortium[#], Kristel Slegers^{7,8}, Hemali Phatnani^{9,10,11}, David A. Knowles^{9,12},
13 Pietro Fratta¹³, Towfique Raj^{1,2,3,4*}

14

15 1. Nash Family Department of Neuroscience & Friedman Brain Institute, Icahn School of Medicine at Mount Sinai,
16 New York, NY, USA

17 2. Ronald M. Loeb Center for Alzheimer's disease, Icahn School of Medicine at Mount Sinai, New York, NY, USA

18 3. Department of Genetics and Genomic Sciences & Icahn Institute for Data Science and Genomic Technology,
19 Icahn School of Medicine at Mount Sinai, New York, NY, USA

20 4. Estelle and Daniel Maggin Department of Neurology, Icahn School of Medicine at Mount Sinai, New York, NY,
21 USA

22 5. Department of Psychiatry, Pamela Sklar Division of Psychiatric Genomics, Icahn School of Medicine at Mount
23 Sinai, New York, NY, 10029, USA

24 6. Graduate School of Biomedical Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA

25 7. Complex Genetics of Alzheimer's Disease Group, Center for Molecular Neurology, VIB, Antwerp, Belgium

26 8. Department of Biomedical Sciences, University of Antwerp, Antwerp, Belgium

27 9. New York Genome Center, New York, NY, USA

28 10. Center for Genomics of Neurodegenerative Disease, New York Genome Center, New York, NY, USA

29 11. Department of Neurology, Columbia University Irving Medical Center, Columbia University, New York, NY,
30 USA

31 12. Department of Computer Science, Columbia University, New York, NY, USA

32 13. Department of Neuromuscular Diseases, UCL Queen Square Institute of Neurology, London, UK

33 **#. Full list of contributors in **Supplementary Acknowledgements**.**

34 ***. Corresponding authors: Jack Humphrey (jack.humphrey@mssm.edu), Towfique Raj (towfique.raj@mssm.edu)**

35 **Abstract**

36 Amyotrophic lateral sclerosis (ALS) is a progressively fatal neurodegenerative disease affecting motor
37 neurons in the brain and spinal cord. Here we investigated gene expression changes in ALS via RNA-
38 seq in 380 post-mortem samples from cervical, thoracic, and lumbar spinal cord segments from 154
39 individuals with ALS and 49 control individuals. We observed an increase in microglia and astrocyte
40 gene expression, accompanied by a decrease in oligodendrocyte gene expression. By creating a gene
41 co-expression network in the ALS samples, we identify several activated microglia modules that
42 negatively correlate with retrospective disease duration. We map molecular quantitative trait loci and
43 find several potential ALS risk loci that may act through gene expression or splicing in the spinal cord
44 and assign putative cell-types for *FNBP1*, *ACSL5*, *SH3RF1* and *NFASC*. Finally, we outline how
45 common genetic variants associated with splicing of *C9orf72* act as proxies for the well-known repeat
46 expansion,, and use the same mechanism to suggest *ATXN3* as a putative risk gene.

47

48 Introduction

49 Amyotrophic lateral sclerosis (ALS) is a progressively fatal neurodegenerative disease characterized by
50 degeneration of upper and lower motor neurons that control voluntary movement via the corticospinal
51 tract. Most patients have a disease onset in middle age but there is a wide clinical variability in onset of
52 symptoms and the pace of disease progression before death¹. About 5-10% of ALS cases have a family
53 history of disease³, with the remaining patients deemed to be sporadic. The field has focused on rare
54 mutations of large effect size, such as large repeat expansions in the gene *C9orf72*, found in 40% of
55 familial ALS and also in 10% of sporadic ALS cases⁴. Other rare mutations, in genes such as *SOD1*,
56 *TARDBP*, *FUS*, *NEK1*, *TBK1*, and *KIF5A* make up only a small fraction of the total familial ALS
57 population⁵⁻⁸ and the majority of non-familial ALS cases have no known causative mutation. Large-scale
58 genome-wide association studies have repeatedly found common genetic variants associated with ALS
59 risk⁸⁻¹¹. Common polymorphic short-tandem repeats are a further contributor to genetic risk of ALS,
60 including *ATXN2*¹², and other *ATXN family members*¹³⁻¹⁵ where intermediate repeat lengths impart a
61 small increase in ALS risk. The interplay between rare and common genetic variants in shaping ALS
62 risk is still being explored. Crucially, there has been little progress in assigning risk genes to particular
63 cell-types. One method to achieve this is the mapping of molecular quantitative trait loci (QTLs), the
64 association between common genetic variants and a molecular phenotype such as gene expression.
65 By performing this in a relevant tissue, QTL variants can be colocalized with GWAS risk variants to
66 identify risk genes¹⁶. In Alzheimer's disease, multiple studies have applied this framework to identify
67 multiple disease risk variants as acting through gene expression and/or splicing in genes specific to
68 microglia and monocytes¹⁷⁻¹⁹.

69
70 Although motor neurons are thought to be the predominantly affected cell type within the spinal cord,
71 much research has focused on non-neuronal contributions to disease initiation and progression. Studies
72 using mouse models of *SOD1* mutations have identified a non-neuronal contribution to disease initiation
73 and length of survival^{20,21}. These studies and many others identified both astrocytes and microglia as
74 being able to modify disease duration²²⁻²⁶. As motor neurons degenerate during disease they release
75 factors which cause microglia to assume an activated pro-inflammatory state^{27,28}, which can then induce
76 an activated state in astrocytes²⁹. Both activated microglia and astrocytes are toxic to motor neurons^{30,31},
77 and blocking this microglia-astrocyte crosstalk extends survival in a *SOD1* mouse model³². Several
78 studies have profiled gene expression in human post-mortem ALS tissues, in spinal cord³³⁻³⁵, frontal
79 cortex³⁵, and motor cortex³⁶. These studies have identified a broad upregulation of inflammatory and
80 immune-related genes and a downregulation in oligodendrocyte and neuron genes. Further

81 investigation of glial activation and neuron-glia crosstalk in the context of ALS is therefore required.
82 However due to small sample sizes, these studies have been unable to identify more subtle changes in
83 gene expression, nor to compare these changes with clinically variable traits, or to leverage molecular
84 QTLs.

85 **Results**

86 *Cellular composition changes in the ALS spinal cord*

87 We aligned and processed post-mortem RNA-seq data from three spinal cord regions (cervical, lumbar,
88 and thoracic) from 154 subjects with ALS and 49 non-neurological controls from the New York Genome
89 Center ALS Consortium, contributed by 8 different medical centres. All samples underwent extensive
90 quality control (**Supplementary Fig. 1-3**). Demographic and technical information for the donors is
91 summarised in **Table 1**; full details are in **Supplementary Table 1**.

92 We performed differential gene expression between all ALS cases and controls in each spinal cord
93 section, controlling for sex, age at death, sequencing batch, submitting site, and technical factors
94 including RNA integrity number (RIN). At a false discovery rate (FDR) < 0.05 we found large numbers
95 of differentially expressed genes (DEGs) in the Cervical and Lumbar regions, with 7,349 and 4,694
96 respectively, and only 256 in the smaller Thoracic cohort (**Fig. 1a-b; Table 2; Supplementary Table**
97 **2**). Of the genes significant in both lumbar and cervical spinal cord, 238 were upregulated with LFC > 1
98 in at least one of two regions, with 109 in both, all of which were more strongly upregulated in the
99 Cervical region (**Fig. 1c**). Although highly concordant, only 12 of those 109 genes passed FDR < 0.05
100 in the thoracic spinal cord, demonstrating the added benefit of our increased sample size. A smaller
101 number of DEGs were strongly upregulated (\log_2 fold change > 2, equivalent to a 4-fold increase in
102 mean expression), including *CHIT1*, *CCL18*, *CHRNA1*, *GPNMB*, and *LYZ*, mostly genes encoding
103 proteins secreted by activated macrophages/microglia. *CHIT1*, encoding the enzyme chitotriosidase, is
104 known to be upregulated in the cerebrospinal fluid (CSF) and plasma of ALS patients³⁷. *GPNMB*,
105 encoding glycoprotein nonmetastatic melanoma B, is upregulated at the protein level in ALS patient
106 spinal cord, CSF, and sera^{38,39} and is expressed by activated microglia⁴⁰. A common genetic variant in
107 *GPNMB* is associated with Parkinson's disease^{41,42}. *CCL18* is a cytokine released by myeloid cells. *LYZ*
108 encodes human lysozyme, an antibacterial protein secreted by myeloid cells. Neither *CCL18* nor *LYZ*
109 have been previously linked to ALS. *CHRNA1*, encoding the alpha subunit of the muscle acetylcholine
110 receptor, is a known marker of denervation of muscles in *SOD1* mouse models⁴³. A marker of astrocyte
111 activation, *C3*^{29,32}, was also upregulated, albeit with a lower effect size.

112 Of the 67 genes downregulated with LFC < -1 in at least one of the two regions, only 13 were < -1
113 both, with the majority (46) more strongly downregulated in the lumbar spinal cord. The downregulated
114 genes include the small nucleolar RNA gene *SNORD3C* and *MOBP*, a marker of oligodendrocytes (**Fig.**
115 **1d**). The motor neuron marker genes *MNX1* and *ISL1* were both downregulated in the cervical and
116 lumbar spinal cord with LFC > -1 (FDR < 0.05). The top 20 strongest (by effect size) upregulated and
117 downregulated genes are presented (**Fig. 1e**). 37 of the 67 most downregulated genes (55%) were non-
118 coding, including antisense transcripts and long intergenic non-coding RNAs, compared to only 38 of
119 the 246 (15%) of the upregulated genes.

120

121 We performed Gene Set Enrichment Analysis (GSEA)⁴⁴ using both curated molecular pathways and
122 sets of cell-type marker genes. Using MSigDB curated pathways⁴⁵, we identified 21 pathways positively
123 enriched in both regions (normalised enrichment score (NES) > 1; adjusted P < 0.05), which mostly
124 reflected different immune and inflammatory pathways (**Fig. 1f; Supplementary Fig. 8**). We next
125 performed GSEA with lists of the 100 most specific human cell-type marker genes for six major brain
126 cell-types⁴⁶. We observed strong positive enrichment of microglia markers, whereas oligodendrocyte
127 markers were negatively enriched (**Fig. 1g**). Repeating the analysis with several other marker gene sets
128 resulted in concordant results and revealed positive enrichments in endothelial cells and pericytes,
129 despite there being low overlap between genes used in each set (**Supplementary Figure 9a-b;**
130 **Supplementary Table 3**).

131 We then prepared a panel of immune activation genes using four studies of microglia and/or astrocyte
132 responses to pro-inflammatory stimuli in mice. These are disease-associated microglia (DAM)⁴⁷,
133 disease-associated astrocytes (DAA)⁴⁸, reactive astrocytes (RA)⁴⁹, and plaque-associated genes
134 (PIG)⁵⁰. These gene lists only partially overlap (**Supplementary Fig. 10**), and represent signatures of
135 microglia and astrocyte responses to a range of stimuli, including amyloid plaques, neurodegeneration,
136 hypoxia (MCAO) and lipopolysaccharide (LPS). All glial activation sets were enriched in the upregulated
137 genes in both regions (**Fig. 1h**).

138

139 We then estimated cell-type proportions in the bulk RNA-seq using both single-nucleus and single-cell
140 RNA-seq from human cortical samples^{46,51}, using two different algorithms^{52,53}, producing four different
141 predictions per sample. Predictions for each cell-type were highly correlated between tools and
142 references (**Supplementary Fig. 11-15; Supplementary Table 4**). We highlight the deconvolution
143 estimates for the cervical spinal cord using single-nucleus RNA-seq reference data from human frontal
144 cortex⁴⁶ and the MuSiC algorithm⁵² (**Fig. 1i**). As a further analysis of cell-type changes we ran
145 expression-weighted cell-type enrichment (EWCE)⁵⁴ using the differentially expressed genes and the

146 same single-nucleus RNA-seq data⁴⁶, which confirmed the observations from deconvolution
147 (**Supplementary Fig. 17**).

148

149 We overlapped the 7,349 cervical spinal cord DEGs (FDR < 0.05) with a recently published mass
150 spectrometry proteomic dataset (**Supplementary Table 5**)³⁹. The study performed differential protein
151 expression in an independent cohort of post-mortem spinal cord samples from 8 ALS cases and 7
152 controls, and cerebrospinal fluid (CSF) from 24 cases and 16 controls. Of the 287 differentially
153 expressed proteins found in the spinal cord (FDR < 0.05), 153 were also DEGs in our dataset (OR =
154 2.8; P < 1e-16, Fisher's exact test), and 137 (90%) had the same direction of effect between mRNA and
155 protein (**Fig. 1j**). The top two most upregulated genes were *GPNMB* and *IQGAP2*. *PEX5L*, found to be
156 highly oligodendrocyte specific in single cell and single nucleus RNA-seq^{46,51} was downregulated at both
157 RNA and protein level. In the CSF, of the 30 genes significant at the protein level, 17 were DEGs (P =
158 0.001), with all but one upregulated (**Fig. 1k**). *GPNMB* and *CHIT1* were both upregulated in CSF,
159 validating their associations with ALS. As well as *GPNMB*, *SERPINA3* was upregulated in both RNA
160 and protein in spinal cord and CSF. Together, these results suggest that ALS spinal cord experiences
161 a robust inflammatory reaction driven by microglia and astrocytes, with dysregulation of
162 oligodendrocytes.

163 *C9orf72-ALS transcriptomes indistinguishable from sporadic ALS*

164 Analysis of frontal cortex and cerebellum has reported distinct sets of differentially expressed genes
165 between *C9orf72* repeat expansion carriers and sporadic ALS and/or FTD patients^{55,56}. We repeated
166 the differential expression analysis but split patients by *C9orf72* repeat expansion status, as assessed
167 by repeat-primed PCR or estimated through ExpansionHunter⁵⁷. Comparing each disease set to
168 controls, the directionality of expression changes in each comparison were highly concordant within
169 each spinal cord section (**Supplementary Fig. 19**). Directly comparing *C9orf72* carriers to sporadic ALS
170 cases, no differentially expressed genes were observed, with the exception of *C9orf72* itself, which was
171 downregulated in *C9orf72*-ALS (cervical spinal cord: log₂ fold change = -0.45; P = 1e-5). This has been
172 previously observed due to hypermethylation of the *C9orf72* promoter in expansion carriers⁵⁸.

173 *Gene co-expression networks associate with disease duration*

174 We then created a weighted gene co-expression network using all 303 ALS samples, adjusting for spinal
175 cord region, contributing site, and other technical factors. We identified 23 modules of co-expressed
176 genes (**Fig. 2a; Supplementary Table 6**) and labelled them in ascending order of size from M1 (50

177 genes) to M23 (3,121). For each module we created a module eigengene (ME), equivalent to the first
178 principal component of the expression of all genes within that module in each sample (**Supplementary**
179 **Table 7**). Modules are presented clustered by eigengene correlation (**Fig. 2a**). Co-expression modules
180 are known to identify cell-types⁵⁹, and 13 modules were significantly enriched with the top 100 cell-type
181 marker genes for the six major cell types of the brain⁴⁶ (**Fig. 2b**; **Supplementary Table 8**). Similarly, we
182 correlated each ME with cell-type proportion estimates in the ALS samples (created using the same
183 Mathys reference and MuSiC algorithm) and found the same modules with marker gene enrichment
184 were strongly positively correlated (Spearman's R = 0.46-0.82) with the respective cell-type proportion
185 (**Supplementary Fig. 20**). Using the same panel of glial activation gene sets as before, we found 6
186 modules enriched with genes from the different sets (**Fig. 2c**). We observed that the module enriched
187 with microglia marker genes (M17) was also enriched for disease-associated microglia and plaque-
188 induced genes, whereas of the four astrocyte marker-enriched modules, one was enriched only with
189 disease-associated astrocytes (M3). The two modules enriched with reactive astrocyte (RA) markers
190 (M9, M18) were enriched with endothelial and/or endothelial cell markers, not astrocytes.

191

192 We next performed gene ontology (GO) enrichment on each module using the GO Biological Process
193 gene sets. Overall, 22 of 23 modules had at least 1 significant GO term (**Supplementary Table 9**). We
194 manually collapsed GO terms into broad sets (**Fig. 2d**). Some sets reflect potentially cell-type specific
195 functions, such as myelination terms with oligodendrocytes, vasculature with endothelial cells/pericytes,
196 and immune response with microglia, whereas modules enriched in terms relating to gene expression
197 and translation were not enriched with cell-type specific or glial activation markers. To assess each
198 module's relevance to ALS-specific changes, we performed enrichment tests using a consensus set of
199 genes upregulated or downregulated in the ALS spinal cord versus controls (**Fig. 2e**). 3 modules were
200 enriched in downregulated genes, two of which were also enriched in oligodendrocyte markers, whereas
201 the six modules enriched with upregulated genes were also enriched with astrocyte, microglia,
202 endothelial, and glial activation markers, confirming our previous cell-type proportion analyses.

203

204 We then used the modules to find associations with clinical variables (**Supplementary Table 10**).
205 Correlating each ME with different clinical traits, we observed 1 module (M3) to correlate with age at
206 death and age of onset (**Fig. 2g**), whereas 5 modules correlated with retrospective disease duration,
207 defined as the length of time between the age at recorded disease onset and age at death. All 3
208 positively correlated modules were enriched with astrocyte marker genes (**Fig. 2h**), and of the two
209 negatively correlated modules were enriched with microglia (**Fig. 2i**) and endothelial marker genes
210 respectively. 2 modules associated with sex, an oligodendrocyte module (M16), and an astrocyte

211 module (M6), suggesting potential cell composition differences between males and females. Our
212 previous study⁶⁰ used these same samples to estimate the abundance of truncated *STMN2* (tSTMN2),
213 a novel cryptic exon transcript created by loss of nuclear TDP-43^{61,62}, which may be a biomarker of TDP-
214 43 pathology⁶⁰. 2 modules correlated with tSTMN2 abundance. One module, M20, was positively
215 correlated with tSTMN2 (**Fig. 2j**). M20 is a large module containing 2048 genes and is enriched with
216 neuronal marker genes, including the full-length *STMN2* gene, as well as motor neuron markers *MNX1*
217 and *ISL1*, though as the sole neuronal module it is likely non-specific to motor neurons. The module
218 negatively correlated with tSTMN2 (M4) is enriched with pericyte marker genes. No modules were
219 significantly associated with the site of motor onset (limb vs bulbar).

220 *Glial composition associates with disease duration*

221 To further investigate the associations with disease duration, we performed a transcriptome-wide
222 correlation analysis with disease duration as a continuous variable. 745 and 39 genes were significantly
223 associated with disease duration at FDR < 0.05 in the cervical and lumbar spinal cord, respectively
224 (**Supplementary Table 11**). Estimated fold-changes represent unit change in expression per month of
225 disease. No effect size threshold was applied. Test statistics for each gene were highly concordant
226 between the cervical and lumbar cords (Pearson R = 0.71, P < 1e-16; **Supplementary Fig 21**). Using
227 GSEA, we found that negatively correlated genes were enriched with microglia markers and microglia
228 activation genes, whereas positively correlated genes were enriched with astrocyte markers and
229 pericyte markers but not astrocyte activation gene sets (**Fig. 3b-c**). Using cell-type proportion estimates
230 from the cervical spinal cord, we observed the same negative correlation between duration and
231 microglial proportion (R = -0.31; adjusted P = 0.002), (**Fig. 3d**), but not with astrocyte proportion (R =
232 0.15; adjusted P = 0.49). One of the strongest positive correlations with disease duration was found for
233 the paraxaonase gene *PON3*, which has been previously linked to ALS through rare mutations⁶³. The
234 previously observed *CHIT1* was found to be the strongest negatively correlated gene with disease
235 duration in both cervical and lumbar spinal cord. There is a non-linear relationship between age of onset
236 and age at death in ALS, with shorter durations seen in both younger and older onset patients. We
237 confirm that the association with *CHIT1* expression is strongest with disease duration, and not with age
238 of onset or death (**Fig. 3e-f**).

239 *Mapping spinal cord QTLs*

240 We took common genetic variants (minor allele frequency > 1%) from the matched whole genome
241 sequencing for all donors of European ancestry (**Supplementary Fig. 22; Supplementary Table 12**)

242 in the cohort, including cases of non-ALS neurodegeneration. We used this to map quantitative trait loci
243 (QTLs) for gene expression and splicing, the latter using the intron-junction clustering method
244 Leafcutter⁶⁴. We identified 9,492 genes with an expression QTL (eQTL) and 5,627 with a splicing QTL
245 (sQTL) in at least one region (**Fig. 4a; Supplementary Fig. 23**). As a comparison, we downloaded
246 summary statistics for the only other available human spinal cord dataset, from GTEx (v8). We
247 discovered substantially more genes with sQTLs than the 965 found by GTEx. Using Storey's π_1 we
248 observed high sharing of QTLs between each region and with GTEx (**Fig. 4b-c**), although sharing was
249 higher in sQTLs than eQTLs, as previously observed^{65,66}. We used our previously generated cell-type
250 proportion estimates to find cell-type interaction QTLs⁶⁷ but no tissue had sufficient power to detect any
251 such associations.

252 *Putative ALS risk variants colocalise with spinal cord QTLs*

253 We then used our QTLs, in combination with GTEx, to prioritise common genetic risk loci using the latest
254 available ALS GWAS⁸ (**Fig. 4d**). Taking a relaxed approach, we extended our search from the 10
255 genome-wide significant loci ($P < 5e-8$) to 64 nominally significant subthreshold loci ($P < 1e-5$)
256 (**Supplementary Table 13**). Among genome-wide significant loci, we identified strong colocalization
257 with QTLs at a posterior probability of colocalization hypothesis 4 (PP4) > 0.8 , only in *C9orf72*. In the
258 *UNC13A* locus we observed a potentially spurious colocalization with *MVB12A* only in GTEx (PP4 =
259 0.5). Among the subthreshold loci, we observed colocalization in 16 loci, with the strongest colocalizing
260 genes (PP4 > 0.8) across our tissues and GTEx seen for *ATXN3*, *GGNBP2*, *ACSL5*, and *FNBP1*
261 (**Supplementary Table 14**).

262 We then ran transcriptome-wide association study (TWAS), an orthogonal method that uses common
263 variants, gene expression, and splicing ratios to predict cis-regulated expression and splicing. TWAS
264 then imputes those models to GWAS summary statistics to identify genes that are associated with
265 disease risk. We generated TWAS models for each spinal cord section and used summary statistics
266 from the latest available ALS GWAS⁸. In both cervical and lumbar spinal cord, splicing in *C9orf72* and
267 *ATXN3* were significantly associated with ALS (FDR < 0.05) (**Supplementary Fig. 24; Supplementary**
268 **Table 15**). The lumbar spinal cord TWAS models also identified an association with expression of
269 *MAPT-AS1* and splicing of *LINC02210* and *LINC02210-CRHR1*. These three genes are within the
270 contentious *MAPT* H1/H2 haplotype region, which has a complex linkage disequilibrium structure, and
271 so are potential false positives. As a comparison, we downloaded pre-computed expression and splicing
272 weights for the dorsolateral frontal cortex ($n = 453$;⁶⁸), which found associations with *C9orf72* in both
273 splicing and expression. In addition, the cortex TWAS models identified *SLC9A8*, *G2E3*, *SCFD1*, and
274 *GPX3* (**Supplementary Fig. 24**).

275 *Prioritised genes annotate to cell-types*

276 We took each colocalised protein-coding gene (PP4 > 0.7) in any of the three spinal cord datasets and
277 annotated them to a cell-type, and to understand how these genes might be involved in ALS. We first
278 took cell-type fidelity ratings⁶⁹, expressed as a fidelity score from 0-100, with high scores suggesting
279 greater cell-type specificity. Although most genes showed no preference towards any cell type, *FNBP1*
280 (fidelity = 92) showed high specificity to oligodendrocytes (**Fig. 5a**; **Supplementary Fig. 25**). We then
281 used the ALS co-expression network modules generated earlier to infer roles for the genes specifically
282 in ALS. Using guilt-by-association, if a gene belongs to a module enriched in a particular cell-type or
283 marker list, it may also be involved in that cell-type. Both *FNBP1* and *SH3RF1* were placed in module
284 M16, highly enriched for oligodendrocytes (**Fig. 5b**). *NFASC* was placed within M6, a module enriched
285 in both astrocyte marker genes and in disease-associated astrocytes, whereas *ACSL5* was located in
286 M14, a module enriched in disease-associated microglia genes but not microglia markers. We then
287 correlated each prioritised gene with estimated cell-type proportions for six cortical cell types⁴⁶. A
288 positive correlation with a particular cell-type proportion is suggestive evidence for specificity. *FNBP1*,
289 *SH3RF1*, and *NFASC* all positively correlated with oligodendrocyte proportions (**Fig. 5c**). *ACSL5*
290 positively correlated with microglia, endothelial and pericyte proportions, with the strongest correlation
291 seen with endothelial cells. Repeating the analysis in just the control samples replicated the correlations
292 between *FNBP1* and oligodendrocytes and *ACSL5* with endothelial cells (**Supplementary Fig. 26**).

293

294 Finally, both *FNBP1* and *SH3RF1* are downregulated in ALS cases, whereas *NFASC* expression is
295 positively associated with disease duration in the cervical spinal cord, the only colocalised gene to do
296 so (**Fig. 5d**). *GGNBP2* was upregulated in ALS patients but did not show a clear cell-type specificity.
297 Despite *C9orf72* being highly expressed in mouse microglia⁷⁰, we observed no associations between
298 *C9orf72* and any cell-type or module.

299 *Splicing QTLs implicate repeat expansions in ALS risk*

300 The *C9orf72* gene produces transcripts from two alternative promoters, exon 1a and exon 1b. The ALS-
301 associated G₄C₂ hexanucleotide repeat expansion (HRE) is located between the two exons (**Fig. 6a**),
302 with more than 30 copies of the HRE considered to be pathogenic⁷¹. The *C9orf72* GWAS locus
303 colocalizes with a splicing QTL in the *C9orf72* transcript in the NYGC lumbar spinal cord, as well as an
304 eQTL in GTEx (**Fig. 4a**). The sQTL increases the usage of the intron J1 connecting exon 1a with exon
305 2, which spans the HRE (**Fig. 6a**). The lead GWAS SNP rs8349943 and the lead sQTL SNP rs1537712
306 are in strong LD in Europeans ($R^2 = 0.75$) and we show that the GWAS SNP is also associated with J1

307 intron usage (**Fig. 6b; Supplementary Table 15**). The lead GWAS SNP rs8349943 is known to tag a
308 founder haplotype which is more susceptible to the HRE⁷². Using ExpansionHunter to estimate the
309 length of the HRE in our cohort, we replicate this finding, as carriers of rs8349943 are also enriched for
310 the HRE (**Fig. 6c**). The usage of the J1 intron is correlated with repeat length (**Fig. 6d**). Therefore, the
311 sQTL colocalization result is likely being driven by the effect of the tagged repeat expansion on the
312 splicing of intron J1.

313

314 We propose a similar mechanism for the colocalization of a subthreshold GWAS locus ($P = 3.2e-7$) with
315 the splicing of *ATXN3*, a promising potential ALS risk gene. The lead SNP rs10143310 was below
316 genome-wide significance in the European ALS GWAS⁸ but crossed the threshold in a multi-ethnic
317 meta-analysis⁷³. A CAG repeat in exon 10 of *ATXN3* is highly polymorphic, and expansions greater than
318 45 copies cause spinocerebellar ataxia type 3 (SCA3), also known as Machado-Josephs disease⁷⁴.
319 SCA3 patients have lower motor neuron loss and have detectable TDP-43 protein inclusions⁷⁵.
320 Intermediate length expansions, not sufficient to cause ataxia, have been shown to increase ALS risk in
321 several other ataxin family genes, most notably *ATXN2*¹², but also *ATXN1*¹³ and *ATXN8OS*¹⁵. Tagging
322 repeat expansions in *ATXN3* with common genetic variants has been previously explored in SCA3
323 patients⁷⁶. In both lumbar and cervical spinal cord samples, and in GTEx, the lead QTL SNP
324 rs200388434 is associated with splicing with a cluster of introns at the 3' end of the *ATXN3* gene, just
325 downstream of the site of the repeat expansion in exon 10 (**Fig. 5e**). The lead QTL SNP rs200388434
326 is in high linkage disequilibrium ($R^2 = 0.93$) with the lead GWAS SNP rs10143310 in Europeans, and
327 rs10143310 also associated with intron splicing (**Fig. 5f**). We hypothesise that the GWAS association
328 is tagging an intermediate length CAG repeat, and this may be the underlying causal genetic factor. We
329 were able to genotype the CAG repeat in 304 individuals in the cohort using ExpansionHunter, observing
330 that the lead QTL SNP associated with a narrow range of repeat lengths ≥ 16 (**Fig. 5g**). CAG repeat
331 length also correlated with splicing in the lumbar spinal cord (**Fig. 5h**).

332 Discussion

333 In this study we assembled the largest ever cohort of post-mortem ALS spinal cords. This has allowed
334 us not only to identify differentially expressed genes when compared to controls, but to identify genes
335 associated with clinical characteristics within the ALS patient cohort. By integrating common genetic
336 variants we prioritise several new candidate ALS genes that may have cell-type-specific functions. In
337 this way, we investigated both the cause (genetic risk) and likely consequence (post-mortem gene
338 expression changes) of disease.

339

340 Comparing ALS cases to controls we identified robust shifts in cell-type in the spinal cord, primarily
341 comprised of a downregulation of oligodendrocytes and motor neurons, and an upregulation in
342 astrocytes and microglia, as well as smaller upward shifts in endothelial cells and pericytes. We observe
343 this across the three spinal cord regions with multiple orthogonal techniques (GSEA, deconvolution,
344 EWCE). However, the interpretation of our results is constrained by the relatively low number of control
345 samples in the cohort, as well as the inherent limitations in the use of bulk tissue sections. The reduction
346 in oligodendrocyte gene expression may reflect genuine cell loss due to secondary demyelination
347 accompanying axonal loss⁷⁷, but this may also reflect a relative shift in proportion compared to increased
348 astrocytes and microglia. For both microglia and astrocytes, although we saw overall upregulation of
349 multiple microglia and astrocyte activation gene lists, it is currently intractable to separate changes in
350 cell-type proportion from changes in cell state in bulk tissue RNA-seq. We also cannot rule out that the
351 increased microglia and activated microglia gene expression signatures may be driven by peripheral
352 monocytes and/or T-cells, which are known to migrate into the spinal cord⁷⁸. We also observed small
353 increases in endothelial cells and pericytes. Alterations to the choroid plexus, including reductions in
354 pericytes, have been observed in ALS⁷⁹. Increases in the recently identified perivascular fibroblast cell-
355 type have been observed in ALS spinal cord RNA-seq as well as ALS mouse models⁸⁰, although we did
356 not explicitly look for this cell type. Crucially, we observed high concordance between our data and a
357 published proteomic dataset from an independent ALS spinal cord and cerebrospinal fluid cohort,
358 suggesting that the gene expression changes we identify are maintained at the protein level, increasing
359 their utility as potential biomarkers.

360 Using co-expression networks built in ALS samples only, we observed a series of associations with
361 disease duration and co-expression modules enriched in microglia and astrocyte genes, in opposing
362 directions. Increased numbers of activated microglia, as measured by CD68 staining in the spinal cord,
363 have been observed in faster progressing ALS patients⁸⁵. However, it is unclear whether microglia
364 activation accelerates neuronal death, or whether microglia activation is an attempted compensatory
365 process, with disease duration driven by some other factor. The negative correlation observed between
366 *CHIT1* expression and disease duration replicates previous findings at the protein level^{37,39,86,87}, but we
367 also find hundreds of new associations, including in the ALS-linked gene *PON3*.

368 By mapping QTLs we provide a genetic resource for the ALS and wider neuroscience community to
369 understand common genetic drivers of gene expression and splicing in the spinal cord.

370 Colocalization has allowed us to prioritise new ALS risk genes, but we must stress that the bulk of our
371 findings rely on nominally significant genetic loci. We are also mindful of the potential for false positive

372 associations due to gene co-expression and LD contamination, which affect both colocalization and
373 TWAS⁸⁹.

374 Taken together, our analyses of the ALS spinal cord point to non-neuronal cells as firmly in the heart of
375 disease in the spinal cord in responding to, and potentially driving, progression of the disease. Our
376 genetic analyses highlight potential new genes that may act on ALS through specific glial cell types.
377 Future genome-wide survival studies may highlight more glial genes in also driving ALS progression.
378 We hope our data are a useful resource for the design of future experiments.

379 **Acknowledgements**

380 We thank all members of the Raj lab for their feedback on the manuscript. This work was supported by
381 grants from NIH NIA R56-AG055824 and U01-AG068880 (J.H. and T.R.), NIH NINDS
382 U54NS123743 (J.H., T.R., P.F.), NIH NIH Medical Scientist Training Program grant T3GM007280
383 (J.T.H.), P.F. is supported by a UK Medical Research Council Senior Clinical Fellowship and Lady Edith
384 Wolfson Fellowship (MR/M008606/1 and MR/S006508/1). F.K. is supported by a BOF DOCPRO
385 fellowship of the University of Antwerp Research Fund. P.F. is supported by the UK Motor Neurone
386 Disease Association, Rosetrees Trust, and the UCLH NIHR Biomedical Research Centre. This work
387 was supported in part through the computational resources and staff expertise provided by Scientific
388 Computing at the Icahn School of Medicine at Mount Sinai. Research reported in this paper was
389 supported by the Office of Research Infrastructure of the National Institutes of Health under award
390 number S10OD018522 and S10OD026880. All NYGC ALS Consortium activities are supported by the
391 ALS Association (ALSA, 19-SI-459) and the Tow Foundation. The funders had no role in study design,
392 data collection and analysis, decision to publish or preparation of the manuscript.

393 **Author contributions**

394 JH and TR conceived and designed the project. JH led the main analysis, with SV, RH, JTH, KPL, FK,
395 KS, MBB, GN, USE contributing code and performing additional data analyses. JH and TR oversaw all
396 aspects of the study, with input from DAK, HP and PF. DF and HP designed the sample collection
397 methodology, reviewed sample and data quality, and coordinated NYGC ALS Consortium postmortem
398 RNA research activity. The NYGC ALS Consortium and the Target ALS Human Postmortem Tissue
399 Core provided human tissue samples as well as pathological, genetic, and clinical information. JH wrote
400 the manuscript with input from all co-authors.

401 **Competing interest declaration**

402 The authors have no competing interests.

403 **Tables**

404

405 **Table 1 - Clinical and technical characteristics of the differential gene expression cohort**

	Control	ALS	ALS-<i>C9orf72</i>	P-value
Donors	49	125	29	-
% Female	53.1%	43.2%	58.6%	0.29
% Bulbar onset	-	24.2%	25.1%	1
Disease duration, months	-	35 (6-156)	31 (6-90)	0.12
Age at death	66 (16-89)	66 (32-85)	64 (50-78)	0.63
Sequencing platform (NovaSeq / HiSeq 2500)	66.7%	59.3%	78.8%	0.104
Tissues				
Cervical Spinal Cord	35	111	28	-
RIN	6.4 (5.1-8.1)	7 (5-9)	6.5 (5.1-8.6)	0.0077
Lumbar Spinal Cord	32	101	21	-
RIN	5.8 (5-7.8)	6.9 (5.1-8.7)	6.3 (5.1-8)	3e-05
Thoracic Spinal Cord	10	37	5	-
RIN	6.35 (5.6-8.1)	6.5 (5-8)	7.5 (5.5-8)	0.79

406 ALS-*C9orf72*: ALS with confirmed *C9orf72* hexanucleotide expansion. RIN: RNA integrity number. Continuous
 407 variables presented as median and range. Categorical variables compared with Fisher's exact test, continuous
 408 variables with Kruskal-Wallis or Wilcoxon rank sum tests. P-values shown are uncorrected for multiple testing, P-
 409 values < 0.05 are bolded

410
411
412

Table 2 - Differentially expressed genes (DEGs) found in each spinal cord region

Region	Control	ALS	Genes tested	All DEGs (FDR < 0.05)	DEGs LFC > 1	DEGs LFC > 2
Cervical	35	139	25,389	7,349	377	29
Thoracic	10	42	19,367	256	65	9
Lumbar	32	122	25,601	4,694	233	7

413

414 **Figure Legends**

415 **Fig. 1 | Differential gene expression in the ALS spinal cord is driven by cell-type composition.** **a-b.** Volcano
416 plots comparing ALS patients to controls in each spinal cord section. P-values for each gene generated from
417 empirical Bayes moderated t-statistics (limma-voom), followed by Benjamini-Hochberg multiple testing adjustment.
418 Genes coloured by whether not differentially expressed (FDR < 0.05; grey), differentially expressed but with
419 modest effects ($|\log_2$ fold change (LFC)| < 1; orange) and with stronger effects ($|\text{LFC}| > 1$; red). Numbers of genes
420 in each category above the plot. **c-d.** Comparing LFC effect sizes between the two regions for the most upregulated
421 (**c**) or downregulated (**d**) genes. **e.** The 20 most upregulated (left) and downregulated (right) genes, ordered by
422 LFC. Asterisks represent Benjamini-Hochberg adjusted P < 0.05 across the 25,389 and 25,601 respective genes
423 from differential expression. **(f-h)** Gene Set Enrichment Analysis results. Normalised enrichment score (NES) is a
424 measure of enrichment of a gene set within a ranked list of genes compared to a permuted background. All
425 pathways are enriched in upregulated genes. Significance derived from empirical P-values from a one-sided
426 permutation test followed by Benjamini-Hochberg correction. **(f)** GSEA results for the 50 molecular signature
427 hallmark pathways genes sets. 100 tests performed. **(g)** GSEA results for the cell-type signature gene sets. 12
428 tests performed. **(h)** GSEA results for the glial activation gene sets. 10 tests performed. DAA: disease-associated
429 astrocytes; DAM: disease-associated microglia; PIG: plaque-induced genes; RA-LPS: reactive astrocytes in
430 response to lipopolysaccharide; RA-MCAO: reactive astrocytes in response to hypoxia. **i.** Estimated cell-type
431 proportions in the cervical spinal cord, between 139 ALS patients and 35 controls. n=174 biologically independent
432 samples. P-values from a two-sided Wilcoxon non-parametric test comparing residuals after regressing technical
433 covariates, followed by Bonferroni correction. 6 tests performed. Boxplots show the median, first and third quartile
434 of the distribution with whiskers extending to 1.5 times the interquartile range. **j-k.** Correlating differentially
435 expressed genes (FDR < 0.05) in the Cervical spinal cord with differentially expressed proteins (FDR < 0.05) from
436 post-mortem spinal cord (**j**) and cerebrospinal fluid (**k**). Asterisks reflect magnitude of adjusted P-values: *** q <
437 1e-4; ** q < 1e-3; * q < 0.05; . q > 0.05.

438
439 **Fig. 2 | Gene co-expression network in the ALS spinal cord.** **(a-g)** Weighted gene co-expression network
440 analysis of 303 ALS spinal cord samples identifies 23 gene modules. **a.** Modules are presented as hierarchical
441 clustering based on module eigengene (ME) correlation. **b-e:** Enrichment results between each module and **b)**
442 cell-type marker genes from Mathys et al, **c)** glial activation genes, **d)** gene ontology (biological process)
443 enrichment, manually collapsed, **e)** differentially expressed genes (FDR < 0.05, no fold change cutoff) between
444 ALS and controls, across all spinal cord regions, **f)** Spearman correlation with disease traits. **g-j).** MEs for each
445 ALS patient. M3 correlates with age of symptom onset, M8 and M17 with duration of disease in months, and M20
446 with *tSTMN2* expression. R refers to Spearman correlation. * refers to Bonferroni adjusted P < 0.05, adjusted for
447 the number of cells in each panel separately. Tests performed: 138 (**b**), 115 (**c**), 3,326 (**d**), 46 (**e**), 161 (**f**). *tSTMN2*
448 - truncated STMN2. TPM - transcripts per million. P-values for **b,c,e** from one-sided Fisher's exact test followed
449 by Bonferroni adjustment, **d** from one-sided hypergeometric test followed by **g:**SCS adjustment, **f** from two-sided
450 Spearman correlation test followed by Bonferroni adjustment.
451

452 **Fig. 3 | Gene expression correlations with duration of disease.** **a.** Volcano plots for correlation in each tissue.
453 Log₂ fold-changes represent unit change in expression per month of disease duration. P-values for each gene
454 generated from empirical Bayes moderated t-statistics (limma-voom), followed by Benjamini-Hochberg multiple
455 testing adjustment **b.** GSEA with cell-type marker genes. 8 tests performed. **c.** GSEA with glia activation gene
456 lists. 10 tests performed. GSEA P-values generated from a one-sided permutation test followed by Benjamini-
457 Hochberg correction **d.** Cell-type proportions in the cervical spinal cord estimated with deconvolution plotted
458 against disease duration. **e.** *CHIT1* is strongly upregulated in ALS in all three tissues. Sample numbers in Table
459 1. **f.** *CHIT1* expression negatively correlates with disease duration, but not with age of onset, and only weakly with
460 age at death. All correlations are Spearman rank correlations. Two-sided P-values in panels b and c are Bonferroni
461 corrected for 12 and 10 tests respectively. P-values in d. are Bonferroni-corrected for 6 tests. Asterisks reflect
462 magnitude of adjusted P-values: *** $q < 1e-4$; ** $q < 1e-3$; * $q < 0.05$; . $q > 0.05$. Boxplots show the median, first
463 and third quartile of the distribution with whiskers extending to 1.5 times the interquartile range.
464

465 **Fig. 4 | Quantitative trait loci (QTL) colocalize with putative ALS risk variants.** **a.** QTL discovery in the three
466 spinal cord tissues and compared with GTEx (v8). Numbers refer to genes with an expression QTL (eGenes) or a
467 splicing QTL (sGenes) at q value < 0.05 . **b-c.** Sharing of QTLs between tissues using Storey's π_1 metric. Values
468 are not symmetric. **d.** Colocalization of subthreshold ALS GWAS loci with spinal cord QTLs. Loci are named for
469 their nearest protein-coding gene. P-values refer to the association of the lead variant in the locus with ALS risk
470 from the GWAS (logistic regression). Numbers refer to the probability of a single shared variant in both GWAS and
471 QTL (PP4). All genes and loci shown with PP4 > 0.5 in at least one QTL dataset. Genes taken for further analysis
472 are in bold font. Circles refer to eQTLs, triangles to sQTLs. PP4: posterior probability of colocalization hypothesis
473 4.
474

475 **Fig. 5 | Annotating colocalised genes with cell-type information.** **a-d.** Each protein-coding gene with PP4 $>$
476 0.7 in at least one spinal cord QTL dataset. **a.** Cell-type fidelity scores from Kelley et al., higher scores imply higher
477 cell-type specificity. **b.** The cell-type and activation marker enrichment p-values (one-sided Fisher's exact test)
478 from Fig. 4 for the modules containing each gene. **c.** Each gene correlated with estimated cell-type proportions in
479 cervical spinal cord in the ALS samples only. Two-sided Pearson correlation test. **d.** Log₂ fold changes from
480 differential expression in ALS vs Control (upper panel) and ALS disease duration (lower panel) in cervical and
481 lumbar spinal cord. P-values in b. and c. Bonferroni adjusted for 138 (b upper panel), 115 (b lower panel), and 72
482 (c) tests. P-values in d. from limma-voom adjusted by Benjamini-Hochberg method for all genes tested in each
483 cohort. Asterisk denotes adjusted p-value < 0.05 .
484

485 **Fig. 6 | Splicing QTLs illuminate genetic associations with repeat expansions in *C9orf72* and *ATXN3*.** **a.**
486 The ALS-causing GGGGCC repeat expansion lies in between the two first exons, 1a and 1b. The intron connecting
487 the exon 1a with exon 2 (J1) has an sQTL in the lumbar spinal cord that colocalises with ALS risk (PP4 = 0.78).
488 **b.** The lead GWAS SNP rs8349943 is associated with J1 intron splicing in the lumbar spinal cord ($P = 1.3e-9$,
489 linear regression, $n=197$ independent samples). **c.** The GGGGCC expansion is only observed in carriers of
490 rs8349943. 30 copies of the repeat is considered the threshold for disease initiation. ($P=1.6e-4$, linear regression,
491 $n=139$ independent samples)⁷¹. **d.** The GGGGCC repeat expansion is associated with J1 intron splicing ($R=0.33$,
492 $P = 8.8e-5$, two-sided Pearson correlation test). **e.** The *ATXN3* gene produces multiple transcripts, including
493 several short transcripts at the 3' end of the gene. Three introns have sQTLs that colocalise with a subthreshold
494 ALS risk GWAS locus with high PP4. The introns are all immediately downstream of a CAG repeat within exon 10.
495 **f.** The lead GWAS SNP rs10143310 is associated with usage of the J1 intron, ($P=3.8e-13$, linear regression, $n=196$
496 independent samples). **g.** Carriers of rs10143310 have a CAG repeat length > 16 copies. ($P=1.1e-7$, linear
497 regression, $n=130$ independent samples). **h.** The length of the CAG repeat correlates with J1 splicing ($p = 0.05$,
498 two-sided Pearson correlation test). Boxplots show the median, first and third quartile of the distribution with
499 whiskers extending to 1.5 times the interquartile range.

500 References

- 501 1. Ravits, J. M. & La Spada, A. R. ALS motor phenotype heterogeneity, focality, and spread: deconstructing motor neuron degeneration.
502 *Neurology* **73**, 805–811 (2009).

- 503 2. Neumann, M. *et al.* Ubiquitinated TDP-43 in Frontotemporal Lobar Degeneration and Amyotrophic Lateral Sclerosis. *Science* **314**, 130–
504 133 (2006).
- 505 3. Byrne, S. *et al.* Rate of familial amyotrophic lateral sclerosis: a systematic review and meta-analysis. *J. Neurol. Neurosurg. Psychiatry*
506 **82**, 623–627 (2011).
- 507 4. Majounie, E. *et al.* Frequency of the C9orf72 hexanucleotide repeat expansion in patients with amyotrophic lateral sclerosis and
508 frontotemporal dementia: A cross-sectional study. *Lancet Neurol.* **11**, 323–330 (2012).
- 509 5. Renton, A. E., Chiò, A. & Traynor, B. J. State of play in amyotrophic lateral sclerosis genetics. *Nat. Neurosci.* **17**, 17–23 (2014).
- 510 6. Cirulli, E. T. *et al.* Exome sequencing in amyotrophic lateral sclerosis identifies risk genes and pathways. *Nat. Methods* **347**, 1436–1441
511 (2016).
- 512 7. Kenna, K. P. *et al.* NEK1 variants confer susceptibility to amyotrophic lateral sclerosis. *Nat. Genet.* (2016) doi:10.1038/ng.3626.
- 513 8. Nicolas, A. *et al.* Genome-wide Analyses Identify KIF5A as a Novel ALS Gene. *Neuron* **97**, 1268–1283.e6 (2018).
- 514 9. van Es, M. A. *et al.* Genome-wide association study identifies 19p13.3 (UNC13A) and 9p21.2 as susceptibility loci for sporadic
515 amyotrophic lateral sclerosis. *Nat. Genet.* **41**, 1083–1087 (2009).
- 516 10. Van Rheenen, W. *et al.* Genome-wide association analyses identify new risk variants and the genetic architecture of amyotrophic lateral
517 sclerosis. *Nat. Genet.* **48**, 1043–1048 (2016).
- 518 11. van Rheenen, W. *et al.* Common and rare variant association analyses in amyotrophic lateral sclerosis identify 15 risk loci with distinct
519 genetic architectures and neuron-specific biology. *Nat. Genet.* **53**, 1636–1648 (2021).
- 520 12. Elden, A. C. *et al.* Ataxin-2 intermediate-length polyglutamine expansions are associated with increased risk for ALS. *Nature* **466**, 1069–
521 1075 (2010).
- 522 13. Tazelaar, G. H. P. *et al.* repeat expansions confer risk for amyotrophic lateral sclerosis and contribute to TDP-43 mislocalization. *Brain*
523 *Commun* **2**, fcaa064 (2020).
- 524 14. Lattante, S. *et al.* ATXN1 intermediate-length polyglutamine expansions are associated with amyotrophic lateral sclerosis. *Neurobiol.*
525 *Aging* **64**, 157.e1–157.e5 (2018).
- 526 15. Hirano, M. *et al.* Noncoding repeat expansions for ALS in Japan are associated with the ATXN8OS gene. *Neurology Genetics* vol. 4
527 e252 Preprint at <https://doi.org/10.1212/nxg.0000000000000252> (2018).
- 528 16. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS*
529 *Genet.* **10**, e1004383 (2014).
- 530 17. Young, A. M. H. *et al.* A map of transcriptional heterogeneity and regulatory variation in human microglia. *Nat. Genet.* **53**, 861–868
531 (2021).
- 532 18. Novikova, G. *et al.* Integration of Alzheimer’s disease genetics and myeloid genomics identifies disease risk regulatory elements and
533 genes. *Nat. Commun.* **12**, 1610 (2021).
- 534 19. Lopes, K. de P. *et al.* Genetic analysis of the human microglial transcriptome across brain regions, aging and disease pathologies. *Nat.*
535 *Genet.* **54**, 4–17 (2022).
- 536 20. Pramatarova, A., Laganière, J., Roussel, J., Brisebois, K. & Rouleau, G. A. Neuron-specific expression of mutant superoxide dismutase
537 1 in transgenic mice does not lead to motor impairment. *J. Neurosci.* **21**, 3369–3374 (2001).

- 538 21. Jaarsma, D., Teuling, E., Haasdijk, E. D., De Zeeuw, C. I. & Hoogenraad, C. C. Neuron-specific expression of mutant superoxide
539 dismutase is sufficient to induce amyotrophic lateral sclerosis in transgenic mice. *J. Neurosci.* **28**, 2075–2088 (2008).
- 540 22. Yamanaka, K. *et al.* Astrocytes as determinants of disease progression in inherited amyotrophic lateral sclerosis. *Nat. Neurosci.* **11**,
541 251–253 (2008).
- 542 23. Lepore, A. C. *et al.* Focal transplantation-based astrocyte replacement is neuroprotective in a model of motor neuron disease. *Nat.*
543 *Neurosci.* **11**, 1294–1301 (2008).
- 544 24. Boillée, S. *et al.* Onset and progression in inherited ALS determined by motor neurons and microglia. *Science* **312**, 1389–1392 (2006).
- 545 25. Wang, L., Sharma, K., Grisotti, G. & Roos, R. P. The effect of mutant SOD1 dismutase activity on non-cell autonomous degeneration in
546 familial amyotrophic lateral sclerosis. *Neurobiol. Dis.* **35**, 234–240 (2009).
- 547 26. Phatnani, H. P. *et al.* Intricate interplay between astrocytes and motor neurons in ALS. *Proc. Natl. Acad. Sci. U. S. A.* **110**, E756–65
548 (2013).
- 549 27. Town, T., Nikolic, V. & Tan, J. The microglial ‘activation’ continuum: from innate to adaptive responses. *J. Neuroinflammation* **2**, 24
550 (2005).
- 551 28. Chiu, I. M. *et al.* A neurodegeneration-specific gene-expression signature of acutely isolated microglia from an amyotrophic lateral
552 sclerosis mouse model. *Cell Rep.* **4**, 385–401 (2013).
- 553 29. Liddelow, S. A. *et al.* Neurotoxic reactive astrocytes are induced by activated microglia. *Nature* **541**, 481–487 (2017).
- 554 30. Zhao, W. *et al.* Activated microglia initiate motor neuron injury by a nitric oxide and glutamate-mediated mechanism. *J. Neuropathol.*
555 *Exp. Neurol.* **63**, 964–977 (2004).
- 556 31. Haidet-Phillips, A. M. *et al.* Astrocytes from familial and sporadic ALS patients are toxic to motor neurons. *Nat. Biotechnol.* **29**, 824–828
557 (2011).
- 558 32. Guttenplan, K. A. *et al.* Knockout of reactive astrocyte activating factors slows disease progression in an ALS mouse model. *Nat.*
559 *Commun.* **11**, 3753 (2020).
- 560 33. D’Erchia, A. M. *et al.* Massive transcriptome sequencing of human spinal cord tissues provides new insights into motor neuron
561 degeneration in ALS. *Scientific Reports* vol. 7 Preprint at <https://doi.org/10.1038/s41598-017-10488-7> (2017).
- 562 34. Brohawn, D. G., O’Brien, L. C. & Bennett, J. P., Jr. RNAseq Analyses Identify Tumor Necrosis Factor-Mediated Inflammation as a Major
563 Abnormality in ALS Spinal Cord. *PLoS One* **11**, e0160520 (2016).
- 564 35. Andrés-Benito, P., Moreno, J., Aso, E., Povedano, M. & Ferrer, I. Amyotrophic lateral sclerosis, gene deregulation in the anterior horn of
565 the spinal cord and frontal cortex area 8: implications in frontotemporal lobar degeneration. *Aging* **9**, 823–851 (2017).
- 566 36. Dols-Icardo, O. *et al.* Motor cortex transcriptome reveals microglial key events in amyotrophic lateral sclerosis. *Neurol Neuroimmunol*
567 *Neuroinflamm* **7**, (2020).
- 568 37. Thompson, A. G. *et al.* Cerebrospinal fluid macrophage biomarkers in amyotrophic lateral sclerosis. *Ann. Neurol.* **83**, 258–268 (2018).
- 569 38. Tanaka, H. *et al.* The potential of GPNMB as novel neuroprotective factor in amyotrophic lateral sclerosis. *Sci. Rep.* **2**, 573 (2012).
- 570 39. Oeckl, P. *et al.* Proteomics in cerebrospinal fluid and spinal cord suggests UCHL1, MAP2 and GPNMB as biomarkers and underpins
571 importance of transcriptional pathways in amyotrophic lateral sclerosis. *Acta Neuropathol.* **139**, 119–134 (2020).
- 572 40. Hüttenrauch, M. *et al.* Glycoprotein NMB: a novel Alzheimer’s disease associated marker expressed in a subset of activated microglia.

- 573 *Acta Neuropathol Commun* **6**, 108 (2018).
- 574 41. Murthy, M. N. *et al.* Increased brain expression of GPNMB is associated with genome wide significant risk for Parkinson's disease on
575 chromosome 7p15.3. *Neurogenetics* **18**, 121–133 (2017).
- 576 42. Nalls, M. A. *et al.* Identification of novel risk loci, causal insights, and heritable risk for Parkinson's disease: a meta-analysis of genome-
577 wide association studies. *Lancet Neurol.* **18**, 1091–1102 (2019).
- 578 43. Halter, B. *et al.* Oxidative stress in skeletal muscle stimulates early expression of Rad in a mouse model of amyotrophic lateral sclerosis.
579 *Free Radic. Biol. Med.* **48**, 915–923 (2010).
- 580 44. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles.
581 *Proc. Natl. Acad. Sci. U. S. A.* **102**, 15545–15550 (2005).
- 582 45. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst* **1**, 417–425 (2015).
- 583 46. Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer's disease. *Nature* **570**, 332–337 (2019).
- 584 47. Keren-Shaul, H. *et al.* A Unique Microglia Type Associated with Restricting Development of Alzheimer's Disease. *Cell* **169**, 1276–
585 1290.e17 (2017).
- 586 48. Habib, N. *et al.* Disease-associated astrocytes in Alzheimer's disease and aging. *Nat. Neurosci.* **23**, 701–706 (2020).
- 587 49. Zamanian, J. L. *et al.* Genomic analysis of reactive astrogliosis. *J. Neurosci.* **32**, 6391–6410 (2012).
- 588 50. Chen, W.-T. *et al.* Spatial Transcriptomics and In Situ Sequencing to Study Alzheimer's Disease. *Cell* **182**, 976–991.e19 (2020).
- 589 51. Darmanis, S. *et al.* A survey of human brain transcriptome diversity at the single cell level. *Proceedings of the National Academy of*
590 *Sciences* vol. 112 7285–7290 Preprint at <https://doi.org/10.1073/pnas.1507125112> (2015).
- 591 52. Wang, X., Park, J., Susztak, K., Zhang, N. R. & Li, M. Bulk tissue cell type deconvolution with multi-subject single-cell expression
592 reference. *Nat. Commun.* **10**, 380 (2019).
- 593 53. Hunt, G. J., Freytag, S., Bahlo, M. & Gagnon-Bartsch, J. A. dtangle: accurate and robust cell type deconvolution. *Bioinformatics* **35**,
594 2093–2099 (2019).
- 595 54. Skene, N. G. & Grant, S. G. N. Identification of Vulnerable Cell Types in Major Brain Disorders Using Single Cell Transcriptomes and
596 Expression Weighted Cell Type Enrichment. *Front. Neurosci.* **10**, 16 (2016).
- 597 55. Prudencio, M. *et al.* Distinct brain transcriptome profiles in C9orf72-associated and sporadic ALS. *Nat. Neurosci.* **18**, 1175–1182 (2015).
- 598 56. Dickson, D. W. *et al.* Extensive transcriptomic study emphasizes importance of vesicular transport in C9orf72 expansion carriers. *Acta*
599 *Neuropathol Commun* **7**, 150 (2019).
- 600 57. Dolzhenko, E. *et al.* ExpansionHunter: A sequence-graph based tool to analyze variation in short tandem repeat regions. Preprint at
601 <https://doi.org/10.1101/572545>.
- 602 58. Jackson, J. L. *et al.* Elevated methylation levels, reduced expression levels, and frequent contractions in a clinical cohort of C9orf72
603 expansion carriers. *Molecular Neurodegeneration* vol. 15 Preprint at <https://doi.org/10.1186/s13024-020-0359-8> (2020).
- 604 59. Oldham, M. C. *et al.* Functional organization of the transcriptome in human brain. *Nat. Neurosci.* **11**, 1271–1282 (2008).
- 605 60. Prudencio, M. *et al.* Truncated stathmin-2 is a marker of TDP-43 pathology in frontotemporal dementia. *J. Clin. Invest.* (2020)
606 doi:10.1172/JCI139741.
- 607 61. Klim, J. R. *et al.* ALS-implicated protein TDP-43 sustains levels of STMN2, a mediator of motor neuron growth and repair. *Nat. Neurosci.*

- 608 22, 167–179 (2019).
- 609 62. Melamed, Z. *et al.* Premature polyadenylation-mediated loss of stathmin-2 is a hallmark of TDP-43-dependent neurodegeneration. *Nat.*
610 *Neurosci.* **22**, 180–190 (2019).
- 611 63. Ticozzi, N. *et al.* Paraoxonase gene mutations in amyotrophic lateral sclerosis. *Ann. Neurol.* **68**, 102–107 (2010).
- 612 64. Li, Y. I. *et al.* Annotation-free quantification of RNA splicing using LeafCutter. *Nat. Genet.* **50**, 151–158 (2018).
- 613 65. Lopes, K. de P. *et al.* Atlas of genetic effects in human microglia transcriptome across brain regions, aging and disease pathologies.
614 Preprint at <https://doi.org/10.1101/2020.10.27.356113>.
- 615 66. Consortium, T. G. & The GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science*
616 vol. 369 1318–1330 Preprint at <https://doi.org/10.1126/science.aaz1776> (2020).
- 617 67. Kim-Hellmuth, S. *et al.* Cell type-specific genetic regulation of gene expression across human tissues. *Science* **369**, (2020).
- 618 68. Li, Y. I., Wong, G., Humphrey, J. & Raj, T. Prioritizing Parkinson’s disease genes using population-scale transcriptomic data. *Nat.*
619 *Commun.* **10**, 994 (2019).
- 620 69. Kelley, K. W., Nakao-Inoue, H., Molofsky, A. V. & Oldham, M. C. Variation among intact tissue samples reveals the core transcriptional
621 features of human CNS cell classes. *Nat. Neurosci.* **21**, 1171–1184 (2018).
- 622 70. O’Rourke, J. G. *et al.* C9orf72 is required for proper macrophage and microglial function in mice. *Science* **351**, 1324–1329 (2016).
- 623 71. Renton, A. E. *et al.* A hexanucleotide repeat expansion in C9ORF72 is the cause of chromosome 9p21-linked ALS-FTD. *Neuron* **72**,
624 257–268 (2011).
- 625 72. DeJesus-Hernandez, M. *et al.* Expanded GGGGCC hexanucleotide repeat in noncoding region of C9ORF72 causes chromosome 9p-
626 linked FTD and ALS. *Neuron* **72**, 245–256 (2011).
- 627 73. Nakamura, R. *et al.* A multi-ethnic meta-analysis identifies novel genes, including ACSL5, associated with amyotrophic lateral sclerosis.
628 *Commun Biol* **3**, 526 (2020).
- 629 74. Paulson, H. Machado-Joseph Disease/Spinocerebellar Ataxia Type 3. *Genetic Instabilities and Neurological Diseases* 363–377 Preprint
630 at <https://doi.org/10.1016/b978-012369462-1/50025-9> (2006).
- 631 75. Seidel, K. *et al.* Axonal inclusions in spinocerebellar ataxia type 3. *Acta Neuropathol.* **120**, 449–460 (2010).
- 632 76. Prudencio, M. *et al.* Toward allele-specific targeting therapy and pharmacodynamic marker for spinocerebellar ataxia type 3. *Sci. Transl.*
633 *Med.* **12**, (2020).
- 634 77. Kang, S. H. *et al.* Degeneration and impaired regeneration of gray matter oligodendrocytes in amyotrophic lateral sclerosis. *Nat.*
635 *Neurosci.* **16**, 571–579 (2013).
- 636 78. Zondler, L. *et al.* Peripheral monocytes are functionally altered and invade the CNS in ALS patients. *Acta Neuropathol.* **132**, 391–411
637 (2016).
- 638 79. Saul, J. *et al.* Global alterations to the choroid plexus blood-CSF barrier in amyotrophic lateral sclerosis. *Acta Neuropathol Commun* **8**,
639 92 (2020).
- 640 80. Månberg, A. *et al.* Publisher Correction: Altered perivascular fibroblast activity precedes ALS disease onset. *Nat. Med.* **27**, 1308 (2021).
- 641 81. Lake, B. B. *et al.* Neuronal subtypes and diversity revealed by single-nucleus RNA sequencing of the human brain. *Science* **352**, 1586–
642 1590 (2016).

- 643 82. Maniatis, S. *et al.* Spatiotemporal dynamics of molecular pathology in amyotrophic lateral sclerosis. *Science* **364**, 89–93 (2019).
- 644 83. Blum, J. A. *et al.* Single-cell transcriptomic analysis of the adult mouse spinal cord reveals molecular diversity of autonomic and skeletal
645 motor neurons. *Nat. Neurosci.* **24**, 572–583 (2021).
- 646 84. Ho, R. *et al.* Cross-Comparison of Human iPSC Motor Neuron Models of Familial and Sporadic ALS Reveals Early and Convergent
647 Transcriptomic Disease Signatures. *Cell Syst* **12**, 159–175.e9 (2021).
- 648 85. Brettschneider, J. *et al.* Microglial activation correlates with disease progression and upper motor neuron clinical symptoms in
649 amyotrophic lateral sclerosis. *PLoS One* **7**, e39216 (2012).
- 650 86. Varghese, A. M. *et al.* Chitotriosidase, a biomarker of amyotrophic lateral sclerosis, accentuates neurodegeneration in spinal motor
651 neurons through neuroinflammation. *J. Neuroinflammation* **17**, 232 (2020).
- 652 87. Pagliardini, V. *et al.* Chitotriosidase and lysosomal enzymes as potential biomarkers of disease progression in amyotrophic lateral
653 sclerosis: a survey clinic-based study. *J. Neurol. Sci.* **348**, 245–250 (2015).
- 654 88. Zeng, B. *et al.* Trans-ethnic eQTL meta-analysis of human brain reveals regulatory architecture and candidate causal variants for brain-
655 related traits. *medRxiv* (2021).
- 656 89. Wainberg, M. *et al.* Opportunities and challenges for transcriptome-wide association studies. *Nat. Genet.* **51**, 592–599 (2019).
- 657 90. van Rheenen, W. *et al.* Author Correction: Common and rare variant association analyses in amyotrophic lateral sclerosis identify 15 risk
658 loci with distinct genetic architectures and neuron-specific biology. *Nat. Genet.* **54**, 361 (2022).
- 659 91. Aspenström, P. Formin-binding proteins: modulators of formin-dependent actin polymerization. *Biochim. Biophys. Acta* **1803**, 174–182
660 (2010).
- 661 92. Wu, C.-H. *et al.* Mutations in the profilin 1 gene cause familial amyotrophic lateral sclerosis. *Nature* **488**, 499–503 (2012).
- 662 93. Nelson, A. D. & Jenkins, P. M. Axonal Membranes and Their Domains: Assembly and Function of the Axon Initial Segment and Node of
663 Ranvier. *Front. Cell. Neurosci.* **11**, 136 (2017).
- 664 94. Efthymiou, S. *et al.* Biallelic mutations in neurofascin cause neurodevelopmental impairment and peripheral demyelination. *Brain* **142**,
665 2948–2964 (2019).
- 666 95. West, R. J. H., Ugbo, C., Gao, F.-B. & Sweeney, S. T. The pro-apoptotic JNK scaffold POSH/SH3RF1 mediates CHMP2BIntron5-
667 associated toxicity in animal models of frontotemporal dementia. *Hum. Mol. Genet.* **27**, 1382–1395 (2018).
- 668 96. Saez-Atienzar, S. *et al.* Genetic analysis of amyotrophic lateral sclerosis identifies contributing pathways and cell types. *Cold Spring*
669 *Harbor Laboratory* 2020.07.20.211276 (2020) doi:10.1101/2020.07.20.211276.
- 670 97. Klemens, J. *et al.* Neurotoxic potential of reactive astrocytes in canine distemper demyelinating leukoencephalitis. *Sci. Rep.* **9**, 11689
671 (2019).

672

673 **Methods**

674 *NYGC ALS Consortium cohort*

675 The 1,917 RNA-seq samples from the January 2020 freeze of the New York Genome Center (NYGC)
676 ALS Consortium were downloaded, comprising of samples from cortical regions, cerebellum and spinal
677 cord. This study used only the spinal cord samples. Diagnosis was determined by each contributing site.
678 Donors include non-neurological disease controls (hereafter controls), those with classical ALS
679 (hereafter ALS), frontotemporal dementia (FTD), mixed pathologies (ALS-FTD, ALS-Alzheimer's), and
680 a small number of other diseases including Primary Lateral Sclerosis, Kennedy's Disease and
681 Parkinson's Disease. *C9orf72* and *ATXN3* repeat expansion lengths were estimated by the Consortium
682 using ExpansionHunter (v2.5.5)¹ on samples that had PCR-free whole genome sequencing available.
683 Patients with greater than 30 repeats were defined as *C9orf72*-ALS. For ALS patients, age of symptom
684 onset and age at death was reported by each contributing site. Disease duration was defined as the
685 difference between age at death and symptom onset, in months. The NYGC ALS Consortium samples
686 presented in this work were acquired through various institutional review board (IRB) protocols from
687 member sites and the Target ALS postmortem tissue core and transferred to the NYGC in accordance
688 with all applicable foreign, domestic, federal, state, and local laws and regulations for processing,
689 sequencing, and analysis. The Biomedical Research Alliance of New York (BRANY) IRB serves as the
690 central ethics oversight body for NYGC ALS Consortium. Ethical approval was given and is effective
691 through 08/22/2022.

692 *RNA-seq processing and quality control*

693 The Consortium's RNA-seq sample processing has been, in part, previously described^{2,3}. In brief, RNA
694 was extracted from flash-frozen postmortem tissue using TRIzol (Thermo Fisher Scientific) chloroform,
695 followed by column purification (RNeasy Minikit, QIAGEN). RNA integrity number (RIN)⁴ was assessed
696 on a Bioanalyzer (Agilent Technologies). RNA-Seq libraries were prepared from 500 ng total RNA using
697 the KAPA Stranded RNA-Seq Kit with RiboErase (KAPA Biosystems) for rRNA depletion and Illumina-
698 compatible indexes (NEXTflex RNA-Seq Barcodes, NOVA-512915, PerkinElmer, and IDT for Illumina
699 TruSeq UD Indexes, 20022370). Pooled libraries (average insert size: 375 bp) passing the quality
700 criteria were sequenced either on an Illumina HiSeq 2500 (125 bp paired end) or an Illumina NovaSeq
701 (100 bp paired-end). Samples were subjected to extensive sequencing and RNA-Seq quality control
702 metrics at the NYGC that are described below. Notably, a set of more than 250 markers was used to
703 confirm tissue, neuroanatomical regions, and sex in the RNA-Seq data. Only samples passing these

704 metrics are available for distribution. The samples had a median sequencing depth of 42 million read
705 pairs, with a range between 16 and 167 million read pairs.

706

707 Samples were uniformly processed using RAPiD-nf, an efficient RNA-Seq processing pipeline
708 implemented in the NextFlow framework⁵. Following adapter trimming with Trimmomatic (version 0.36)
709 ⁶, all samples were aligned to the hg38 build (GRCh38.primary_assembly) of the human reference
710 genome using STAR (2.7.2a)⁷, with indexes created from GENCODE, version 30 ⁸. Gene expression
711 was quantified using RSEM (1.3.1)⁹. Quality control was performed using SAMtools (v1.9) ¹⁰ and Picard
712 (v2.22.3), and the results were collated using MultiQC (v1.8)¹¹.

713

714 Aligned RNA-seq samples were subjected to quality control modelled on the criteria of the Genotype
715 Tissue Expression Consortium¹². Any sample failing 1 of the following sequencing metric thresholds
716 was removed: a unique alignment rate of less than 90%, ribosomal bases of greater than 10%, a
717 mismatch rate of greater than 1%, a duplication rate of greater than 0.5%, intergenic bases of less than
718 10.5%, and ribosomal bases of greater than 0.1%. For tissue identity, both principal components
719 analysis and UMAP were performed on the TMM-normalised gene expression matrix, followed by k-
720 means clustering. This identified three clusters of samples, grouped by cerebellum, cortical regions, and
721 spinal cord. Samples that clustered with a non-matching tissue type were flagged and tissue identity
722 was re-confirmed using the expression of the cerebellar marker *CBNL1*, the cortical marker *NRGN* and
723 the oligodendrocyte marker *MOBP*. 19 samples were removed for having ambiguous tissue identity. For
724 duplicate samples, where samples of the same tissue from the same donor were sequenced, the sample
725 with the highest RIN was retained, this removed 15 duplicate samples. Sex was confirmed using *XIST*
726 and *UTY* expression. 11 samples with missing sex information were confirmed as males. Due to the
727 large impact of RNA integrity number (RIN) on expression, only samples with RIN ≥ 5 were included
728 in the differential expression analysis, totalling 380 spinal cord samples from 203 donors. For the QTL
729 analyses (see below), no RIN threshold was applied.

730 *Covariate selection and modelling for differential expression*

731 The following was run for each tissue separately: Clinical variables (disease status, age at death, sex,
732 contributing site) were combined with sequencing variables (RIN, sequencing preparation method,
733 sequencing platform), technical metrics of the RNA-seq libraries from Picard (% mRNA bases, 3' bias,
734 etc), and genotype principal components (see below). Using voom-normalised gene expression
735 removing lowly expressed genes, principal components analysis was performed. The top 10 principal
736 components were then associated with each potential confounding variable using a linear model,

737 estimating the variance explained (r^2) of the confounder on each principal component (**Supplementary**
738 **Fig. 3a**). Using an orthogonal approach, variancePartition (v1.21.6)¹³ was run on a reduced set of
739 confounding variables, taking only the nominally independent sequencing metrics (**Supplementary Fig.**
740 **3b**).

741

742 For performing differential gene expression between ALS and control samples, multiple model designs
743 were fitted to account for differences in sequencing batch and contributing site, both of which are
744 correlated with disease status. To account for potentially non-linear dependence of RIN and age at
745 death, squared terms were included. To account for potential confounding differences due to genetic
746 background, the first 5 genotype principal components (gPCs) from smartpca (v6.0.1)¹⁴ were included.
747 For filtering lowly expressed genes, a permission threshold of median TPM > 0 was applied, resulting
748 in 24-25,000 genes being kept for each analysis. For Cervical and Lumbar spinal cord, the following
749 model was fitted: $expression \sim disease + sex + library\ preparation\ method + contributing\ site + age +$
750 $age^2 + RIN + RIN^2 + \% mRNA\ bases + gPC1 + gPC2 + gPC3 + gPC4 + gPC5$. For the smaller set of
751 Thoracic spinal cord samples, a reduced model was fitted as it maximised the gene-gene correlation of
752 differential expression effect sizes with the other two regions: $expression \sim disease + sex + RIN + RIN^2$
753 $+ age + age^2 + library\ preparation\ method + gPC1 + gPC2 + gPC3 + gPC4 + gPC5$. Differential gene
754 expression was fitted using limma voom (v3.46.0)¹⁵ on TMM-normalized¹⁶ read counts. P-values were
755 adjusted for multiple testing using FDR correction, with genes were considered differentially expressed
756 at FDR < 0.05. A gene was considered to have a moderate effect size at $|\log_2\ \text{fold change}| > 1$.

757 For transcriptome-wide correlations with disease duration, the same models as before were used in the
758 ALS samples only, with disease duration (years) used as continuous variables. Downsampling was
759 performed by taking random subsets of either the Cervical or Lumbar samples, without replacement.

760 *Gene set enrichment analysis*

761 Sets of genes were collected from multiple sources and compared to the full differential expression
762 results for each tissue using Gene Set Enrichment Analysis (GSEA)¹⁷, as implemented in the
763 ClusterProfiler R package (v3.18.1)¹⁸. As input we included all tested genes from the differential
764 expression or disease duration analysis for each tissue at nominal (unadjusted) P-value < 0.05, ranked
765 by \log_2 fold change. For each gene set, a running cumulative tally is made of whether genes in a set are
766 present or absent during a walk down the list. The maximal score during the walk is the enrichment
767 score (ES), which reflects the degree of which a gene set is enriched at either the top or bottom of a list.
768 Labels are then randomly permuted to generate an empirical null ES distribution and a P-value is
769 calculated. To aid comparison between sets, each ES is then divided by the mean null ES to create a

770 normalised enrichment score (NES). Hallmark pathway gene sets (h.all.v7.2.symbols.gmt) were
771 downloaded from the molecular signatures database¹⁹. Cell-type marker genes were created using
772 single cell RNA-seq²⁰ and single nucleus RNA-seq²¹ from human cortex. For each dataset, the top 100
773 cell-type specific genes were calculated by comparing gene expression of each cell-type group against
774 the mean of all cells in Limma Voom. Marker genes for astrocytes, microglia, neurons, oligodendrocytes,
775 and pericytes were downloaded from the Kelley et al²², PanglaoDB²³ and Neuroexpresso²⁴ websites
776 (**see URLs**). Disease-associated Microglia (DAM) signature genes²⁵, Disease-associated astrocytes²⁶,
777 Plaque-associated genes²⁷, and LPS and MCAO-activated astrocyte genes²⁸ were downloaded from
778 their respective supplementary materials. Mouse genes were lifted over to their human homologues
779 using Homologene²⁹. Any duplicate gene name, or gene name without a matching Ensembl ID in
780 GENCODE v30 was removed.

781 *Re-analysis of proteomics data*

782 Summary statistics from a published study³⁰ applying isobaric tags for relative and absolute
783 quantification (iTRAQ) proteomics for cerebrospinal fluid (26 ALS, 16 Control) and label-free proteomics
784 to human spinal cord (8 ALS, 7 control) were downloaded from the study's supplementary data files. A
785 total of 1,929 proteins were tested in the cerebrospinal fluid, of which 32 were called significant at FDR
786 < 0.05. 5,115 peptides were tested in spinal cord samples, of which 292 were called significant at FDR
787 < 0.05. Peptides assigned to multiple genes were discarded, resulting in 287 genes in the spinal cord
788 and 30 in CSF.

789 *Cell-type deconvolution*

790 Filtered counts and cell-type labels for single nucleus RNA-seq from 80,660 cells from 48 human
791 dorsolateral prefrontal cortex samples²¹ were downloaded from Synapse (syn18681734). Only cells
792 from the 14 donors without dementia were kept. Single cell RNA-seq data of 466 cells from 12 donors²⁰
793 was downloaded from Gene Expression Omnibus (GSE67835) using the count matrices and cell-type
794 labels provided. Bulk spinal cord RNA-seq data was voom-normalized before deconvolution was
795 estimated using MuSiC (v0.1.1)³¹, a method which incorporates the variance between multiple donors
796 from single cell/nucleus RNA-seq. In addition, we ran dtangle (v2.0.9)³², which requires marker genes
797 to be generated for each cell-type. Markers were created using Voom to compare each gene in purified
798 cell-type to the mean of all cell-types. The top 100 genes ranked by effect size were used as cell-type
799 markers. Estimated proportions of each cell-type were compared between ALS and control using non-
800 parametric Wilcoxon tests after regressing the same technical covariates above. P-values were

801 corrected for multiple testing using Bonferroni correction. For comparing duration of onset, estimated
802 cell-type proportions were correlated using a Spearman correlation.

803 *Expression-weighted Cell-type Enrichment*

804 Expression-weighted cell-type enrichment analysis was performed using the EWCE package³³. Cell-
805 type specificity scores for each gene were created using human frontal cortex single-nucleus RNA-
806 seq²¹. Cell-type enrichment results were generated using the top 250 upregulated and downregulated
807 genes, ordered by t-statistic, for the differential expression results for each segment. Specificity scores
808 for each set were then compared to the mean of the empirical null distribution from 10,000 random gene
809 sets. Enrichment was expressed as the number of standard deviations from the mean. P-values were
810 Bonferroni corrected for multiple testing. Significance was set at adjusted $P < 0.05$.

811 *Gene co-expression Networks*

812 Gene expression from all 303 ALS samples from the three spinal cord regions was combined into a
813 single matrix. Genes annotated as protein-coding by Ensembl were kept, and only then if each gene
814 had at least 1 read count per million in at least 50% of samples, resulting in 16,992 genes. Gene counts
815 were then transformed using Voom and TMM normalization. The following covariates were then
816 regressed out using `removeBatchEffect()`: library preparation, contributing site, spinal cord section, RIN,
817 % mRNA bases, and genomic PCs 1-5. Co-expression network analysis was performed using Weighted
818 Gene Correlation Network Analysis (WGCNA; v1.70-3) following a standard pipeline. Scale-free
819 topology ($R^2 > 0.8$) was achieved by applying a soft threshold power of 8 into a signed network model.
820 The adjacency matrices were constructed using the average linkage hierarchical clustering of the
821 topological overlap dissimilarity matrix (1-TOM). Co-expression modules were defined using a dynamic
822 tree cut method with minimum module size of 50 genes and deep split parameter of 4. Modules highly
823 correlated with each other, corresponding to a module eigengene (ME) correlation > 0.75 , were merged,
824 resulting in a total of 23 modules. Modules were labelled according to their size.

825 We calculated the Spearman correlation between each module eigengene and the following clinical
826 variables: age of disease onset, age at death, disease duration (years), site of disease onset (bulbar or
827 limb), *C9orf72* status, sex, and *tSTMN2* abundance. *tSTMN2* abundance in TPM for the matching
828 samples was extracted from the supplementary data from³.

829 Cell-type and glial activation genes were tested for enrichment within each module using Fisher's exact
830 test using a background set of 16,922 genes, followed by Bonferroni correction for the number of tests
831 performed. Gene ontology biological process terms were tested for enrichment using the gProfiler2

832 package (v0.2.0)³⁴. Terms with less than 10 genes were removed before correction for multiple testing.
833 Enriched terms were then manually grouped into sets for presentation. Full module assignments,
834 eigengenes, and enrichment results are shared as **Supplementary Tables 6-10**.

835 *Quantitative Trait Loci mapping*

836 To perform expression QTL (eQTL) mapping, we created a pipeline based on the one created by the
837 GTEX consortium. We completed a separate normalization and filtering method to previous analyses.
838 Gene expression matrices were created from the RSEM output using tximport⁴⁷. Matrices were then
839 converted to GCT format, TMM normalized, filtered for lowly expressed genes, removing any gene with
840 less than 0.1 TPM in 20% of samples and at least 6 counts in 20% of samples. Each gene was then
841 inverse-normal transformed across samples. PEER⁴⁸ factors were calculated to estimate hidden
842 confounders within our expression data. We created a combined covariate matrix that included the
843 PEER factors and the first 5 genotyping principal component values as input to the analysis. We tested
844 numbers of PEER factors from 0 to 30 and found that between 10 and 30 factors produced the largest
845 number of eGenes in each region (**Supplementary Fig. 23**).

846 To test for cis-eQTLs, linear regression was performed using the tensorQTL (v1.0.5)⁴⁹ *cis_nominal* mode
847 for each SNP-gene pair using a 1 megabase window within the transcription start site (TSS) of a gene.
848 To test for association between gene expression and the top variant in cis we used tensorQTL cis
849 permutation pass per gene with 1000 permutations. To identify eGenes, we performed q-value
850 correction of the permutation P-values for the top association per gene at a threshold of 0.05.

851 We performed splicing quantitative trait loci (sQTL) analysis using the splice junction read counts
852 generated by regtools (v0.5.1)⁵⁰. Junctions were clustered using Leafcutter (psi_2019 branch)⁵¹,
853 specifying for each junction in a cluster a maximum length of 100kb. Following the GTEx pipeline, introns
854 without read counts in at least 50% of samples or with fewer than 10 read counts in at least 10% of
855 samples were removed. Introns with insufficient variability across samples were removed. Filtered
856 counts were then normalized using *prepare_phenotype_table.py* from Leafcutter, merged, and
857 converted to BED format, using the coordinates from the middle of the intron cluster. We created a
858 combined covariate matrix that included the PEER factors and the first 5 genotype principal components
859 as input to the analysis. We mapped sQTLs with between 0 and 30 PEER factors as covariates in our
860 QTL model and determined 5 and 15 factors produce the largest number of sGenes (**Supplementary**
861 **Fig. 23**).

862
863 To test for cis sQTLs, linear regression was performed using the tensorQTL nominal pass for each SNP-
864 junction pair using a 100kb window from the center of each intron cluster. To test for association between

865 intronic ratio and the top variant in cis we used tensorQTL permutation pass, grouping junctions by their
866 cluster using --grp option. To identify significant clusters, we performed q-value⁵² correction using a
867 threshold of 0.05.

868 We estimated pairwise replication (π_1) of eQTLs and sQTLs using the q-value R package. This involves
869 taking the SNP-gene pairs that are significant at q-value < 0.05 in the discovery dataset and extracting
870 the unadjusted P-values for the matched SNP-gene pairs in the replication dataset.

871 *GTEx Spinal Cord QTL summary statistics*

872 Full summary statistics for the cervical spinal cord expression QTLs (v8) were downloaded from the
873 eQTL catalogue (**see URLs**). The splicing QTLs were downloaded from the Google Cloud portal. Top
874 associations for each gene were downloaded from the GTEx portal.

875 *Genome-wide association study summary statistics*

876 Full summary statistics for the 2018 ALS GWAS⁵³ were downloaded from the EBI GWAS Catalogue,
877 which have lifted over the variants to the hg38 build. Genome-wide significant loci were taken to be the
878 most significant variants within 1 megabase at a threshold of $P < 5e-8$. Subthreshold loci were defined
879 at a relaxed threshold of $P < 1e-5$. Loci were named by their nearest protein-coding gene using
880 SNPnexus (v4)⁵⁴.

881 *Colocalization analysis*

882 We used coloc (v3.2-1)⁵⁵ to test whether SNPs from different loci in the ALS GWAS colocalized with
883 expression and splicing QTLs from the spinal cord. For each genome-wide and subthreshold locus in the
884 ALS GWAS we extracted the nominal summary statistics of association for all SNPs within 1 megabase
885 either upstream/downstream of the top lead SNP (2Mb-wide region total). In each QTL dataset we then
886 extracted all nominal associations for all SNP-gene pairs within that range and tested for colocalization
887 between the GWAS locus and each gene. To avoid spurious colocalization caused by long range linkage
888 disequilibrium, we restricted our colocalizations to GWAS SNP - eQTL SNP pairs where the distance
889 between their respective top SNPs was $\leq 500\text{kb}$ or the two lead SNPs were in moderate linkage
890 disequilibrium ($r^2 > 0.1$), taken from the 1000 Genomes (Phase 3) European populations using the

891 LDLinkR package (v1.1.2)⁵⁶. For splicing QTLs we followed the same approach but collapsed junctions
892 to return only the highest PP4 value for each gene in each locus. Due to the smaller window of
893 association (100kb from the center of the intron excision cluster) we restricted reported colocalizations
894 to cases where the GWAS SNP and the top sQTL SNP were either within 100kb of each other or in
895 moderate linkage disequilibrium ($r^2 > 0.1$).

896

897 All plots were created using ggplot2 (v3.3.3)⁶¹ in R (version 4.0.4), with ggrepel (v0.9.1)⁶², ggfortify
898 (v0.4.11)⁶³, patchwork (v1.1.1)⁶⁴, ggbreak (v0.0.9)⁶⁵, and ggbio (v1.38.0)⁶⁶ for additional layers of
899 visualization.

900 **Data availability**

901 All raw RNA-seq data can be accessed via the NCBI's GEO database (GEO GSE137810, GSE124439,
902 GSE116622, and GSE153960). Processed gene expression count matrices with de-identified metadata
903 have been deposited on Zenodo (10.5281/zenodo.6385747) and we provide an RMarkdown vignette
904 on downloading them and performing differential expression (**see URLs**). In addition, we provide an
905 interactive R Shiny app to visualise the gene expression and other clinical variable associations (**see**
906 **URLs**). Full summary statistics for expression and splicing QTLs have been deposited on Zenodo
907 (10.5281/zenodo.5248758). All TWAS weight files have been deposited on Zenodo
908 (10.5281/zenodo.5256613). All RNA-seq and whole genome sequencing data generated by the NYGC
909 ALS Consortium are made immediately available to all members of the Consortium and with other
910 consortia with whom we have a reciprocal sharing arrangement. To request immediate access to new
911 and ongoing data generated by the NYGC ALS Consortium and for samples provided through the Target
912 ALS Postmortem Core, complete a genetic data request form at CGND_help@nygenome.org. All whole
913 genome sequencing data will be deposited on dbGaP at the conclusion of the project in late 2023.

914 **Code availability**

915 All analysis code written in R is available in Rmarkdown workbooks in a Github repository, and specific
916 data processing pipelines are in separate repositories (**see URLs**).

917 **URLs**

918 Website associated with this manuscript, including all code notebooks written for this project:
919 https://jackhump.github.io/ALS_SpinalCord_QTLs/

920 Gene expression counts and TPMs with de-identified metadata:
921 <https://zenodo.org/record/6385747>
922 Code vignette demonstrating how to download data and perform differential expression with R:
923 https://jackhump.github.io/ALS_SpinalCord_QTLs/html/DE_Vignette.html
924 R Shiny app for visualisation:
925 https://jackhumphrey.shinyapps.io/als_spinal_cord_browser/
926 Full QTL summary statistics:
927 <https://zenodo.org/record/5248758>
928 Full TWAS weights:
929 <https://doi.org/10.5281/zenodo.5256613>
930 Molecular Signatures Database (MSigDb):
931 <http://www.gsea-msigdb.org/gsea/msigdb/index.jsp>
932 Kelley et. al. gene fidelity marker genes:
933 <http://oldhamlab.ctec.ucsf.edu/data-download/>
934 Neuroexpresso marker genes:
935 <http://neuroexpresso.org/>
936 PanglaoDB marker genes:
937 <https://panglaodb.se/>
938 ENCODE Blacklist:
939 <https://github.com/Boyle-Lab/Blacklist/blob/master/lists/hg38-blacklist.v2.bed.gz>
940 WGS QC pipeline:
941 <https://github.com/jackhump/WGS-QC-Pipeline>
942 QTL mapping pipeline:
943 <https://github.com/RajLabMSSM/QTL-mapping-pipeline>
944 DLPFC TWAS weights:
945 <http://gusevlab.org/projects/fusion/#reference-functional-data>
946 ExpansionHunter:
947 <https://github.com/Illumina/ExpansionHunter>
948 [SNPNexus:](https://www.snp-nexus.org/v4/)
949 <https://www.snp-nexus.org/v4/>
950 VCFs of 1000 Genomes samples:
951 ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/release/20190312_bi
952 [allelic SNV and INDEL/](https://www.1000genomes.org/data/1000Gomes/1000Gomes_b37/VCFs/1000Gomes_b37_VCFs.html)

953 **Methods-only references**

- 954 1. Dolzhenko, E. *et al.* ExpansionHunter: A sequence-graph based tool to analyze variation in short tandem repeat regions. Preprint at
955 <https://doi.org/10.1101/572545>.
- 956 2. Tam, O. H. *et al.* Postmortem Cortex Samples Identify Distinct Molecular Subtypes of ALS: Retrotransposon Activation, Oxidative
957 Stress, and Activated Glia. *Cell Rep.* **29**, 1164–1177.e5 (2019).
- 958 3. Prudencio, M. *et al.* Truncated stathmin-2 is a marker of TDP-43 pathology in frontotemporal dementia. *J. Clin. Invest.* (2020)
959 [doi:10.1172/JCI139741](https://doi.org/10.1172/JCI139741).
- 960 4. Schroeder, A. *et al.* The RIN: an RNA integrity number for assigning integrity values to RNA measurements. *BMC Mol. Biol.* **7**, 3 (2006).
- 961 5. Di Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nat. Biotechnol.* **35**, 316–319 (2017).
- 962 6. Bolduc, B. Quality Control of Reads Using Trimmomatic (Cyverse) v1 (protocols.io.ewbbfan). *protocols.io* Preprint at
963 <https://doi.org/10.17504/protocols.io.ewbbfan>.
- 964 7. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 965 8. Harrow, J., Frankish, A., Gonzalez, J. M. & Frazer, K. A. GENCODE : The reference human genome annotation for The ENCODE
966 Project. *Genome Res.* **22**, 1760–1774 (2012).
- 967 9. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC*
968 *Bioinformatics* **12**, 323 (2011).
- 969 10. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
- 970 11. Ewels, P., Magnusson, M., Lundin, S. & Källér, M. MultiQC: summarize analysis results for multiple tools and samples in a single report.
971 *Bioinformatics* **32**, 3047–3048 (2016).
- 972 12. Consortium, T. G. & The GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science*
973 vol. 369 1318–1330 Preprint at <https://doi.org/10.1126/science.aaz1776> (2020).
- 974 13. Hoffman, G. E. & Schadt, E. E. variancePartition: interpreting drivers of variation in complex gene expression studies. *BMC*
975 *Bioinformatics* **17**, 483 (2016).
- 976 14. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38**, 904–909
977 (2006).
- 978 15. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: Precision weights unlock linear model analysis tools for RNA-seq read counts.
979 *Genome Biol.* **15**, R29 (2014).
- 980 16. Risso, D., Schwartz, K., Sherlock, G. & Dudoit, S. GC-content normalization for RNA-Seq data. *BMC Bioinformatics* **12**, 480 (2011).
- 981 17. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles.
982 *Proc. Natl. Acad. Sci. U. S. A.* **102**, 15545–15550 (2005).
- 983 18. Yu, G., Wang, L.-G., Han, Y. & He, Q.-Y. clusterProfiler: an R Package for Comparing Biological Themes Among Gene Clusters.
984 *OMICS: A Journal of Integrative Biology* vol. 16 284–287 Preprint at <https://doi.org/10.1089/omi.2011.0118> (2012).
- 985 19. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst* **1**, 417–425 (2015).
- 986 20. Darmanis, S. *et al.* A survey of human brain transcriptome diversity at the single cell level. *Proceedings of the National Academy of*

- 987 *Sciences* vol. 112 7285–7290 Preprint at <https://doi.org/10.1073/pnas.1507125112> (2015).
- 988 21. Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer’s disease. *Nature* **570**, 332–337 (2019).
- 989 22. Kelley, K. W., Nakao-Inoue, H., Molofsky, A. V. & Oldham, M. C. Variation among intact tissue samples reveals the core transcriptional
990 features of human CNS cell classes. *Nat. Neurosci.* **21**, 1171–1184 (2018).
- 991 23. Franzén, O., Gan, L.-M. & Björkegren, J. L. M. PanglaoDB: a web server for exploration of mouse and human single-cell RNA
992 sequencing data. *Database* **2019**, (2019).
- 993 24. Mancarci, B. O. *et al.* Cross-Laboratory Analysis of Brain Cell Type Transcriptomes with Applications to Interpretation of Bulk Tissue
994 Data. *eNeuro* **4**, (2017).
- 995 25. Keren-Shaul, H. *et al.* A Unique Microglia Type Associated with Restricting Development of Alzheimer’s Disease. *Cell* **169**, 1276–
996 1290.e17 (2017).
- 997 26. Habib, N. *et al.* Disease-associated astrocytes in Alzheimer’s disease and aging. *Nat. Neurosci.* **23**, 701–706 (2020).
- 998 27. Chen, W.-T. *et al.* Spatial Transcriptomics and In Situ Sequencing to Study Alzheimer’s Disease. *Cell* **182**, 976–991.e19 (2020).
- 999 28. Zamanian, J. L. *et al.* Genomic analysis of reactive astrogliosis. *J. Neurosci.* **32**, 6391–6410 (2012).
- 1000 29. Mancarci, O. & French, L. Homologene: quick access to homologene and gene annotation updates. *R package version 1*, 68 (2019).
- 1001 30. Oeckl, P. *et al.* Proteomics in cerebrospinal fluid and spinal cord suggests UCHL1, MAP2 and GPNMB as biomarkers and underpins
1002 importance of transcriptional pathways in amyotrophic lateral sclerosis. *Acta Neuropathol.* **139**, 119–134 (2020).
- 1003 31. Wang, X., Park, J., Susztak, K., Zhang, N. R. & Li, M. Bulk tissue cell type deconvolution with multi-subject single-cell expression
1004 reference. *Nat. Commun.* **10**, 380 (2019).
- 1005 32. Hunt, G. J., Freytag, S., Bahlo, M. & Gagnon-Bartsch, J. A. dtangle: accurate and robust cell type deconvolution. *Bioinformatics* **35**,
1006 2093–2099 (2019).
- 1007 33. Skene, N. G. & Grant, S. G. N. Identification of Vulnerable Cell Types in Major Brain Disorders Using Single Cell Transcriptomes and
1008 Expression Weighted Cell Type Enrichment. *Front. Neurosci.* **10**, 16 (2016).
- 1009 34. Reimand, J., Kull, M., Peterson, H., Hansen, J. & Vilo, J. g:Profiler—a web-based toolset for functional profiling of gene lists from large-
1010 scale experiments. *Nucleic Acids Res.* **35**, W193–W200 (2007).
- 1011 35. Regier, A. A. *et al.* Functional equivalence of genome sequencing analysis pipelines enables harmonized variant calling across human
1012 genetics projects. *Nat. Commun.* **9**, 4038 (2018).
- 1013 36. Poplin, R. *et al.* Scaling accurate genetic variant discovery to tens of thousands of samples. Preprint at <https://doi.org/10.1101/201178>.
- 1014 37. Amemiya, H. M., Kundaje, A. & Boyle, A. P. The ENCODE Blacklist: Identification of Problematic Regions of the Genome. *Sci. Rep.* **9**,
1015 9354 (2019).
- 1016 38. Aguet, F. *et al.* The GTEx Consortium atlas of genetic regulatory effects across human tissues. *bioRxiv* 787903 (2019)
1017 doi:10.1101/787903.
- 1018 39. Adelson, R. P. *et al.* Empirical design of a variant quality control pipeline for whole genome sequencing data using replicate
1019 discordance. *Sci. Rep.* **9**, 16156 (2019).
- 1020 40. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, (2021).
- 1021 41. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).

1022 42. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience* vol. 4 Preprint at
1023 <https://doi.org/10.1186/s13742-015-0047-8> (2015).

1024 43. Köster, J. & Rahmann, S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics* **28**, 2520–2522 (2012).

1025 44. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).

1026 45. Fort, A. *et al.* MBV: a method to solve sample mislabeling and detect technical bias in large combined genotype and sequencing assay
1027 datasets. *Bioinformatics* **33**, 1895–1897 (2017).

1028 46. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867–2873 (2010).

1029 47. Love, M. I., Soneson, C. & Robinson, M. D. Importing transcript abundance datasets with tximport. *dim (txi. inf. rep \$ infReps \$ sample1)*
1030 **1**, 5 (2017).

1031 48. Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of expression residuals (PEER) to obtain increased
1032 power and interpretability of gene expression analyses. *Nat. Protoc.* **7**, 500 (2012).

1033 49. Taylor-Weiner, A. *et al.* Scaling computational genomics to millions of individuals with GPUs. *Genome Biol.* **20**, 228 (2019).

1034 50. Feng, Y.-Y. *et al.* RegTools: Integrated analysis of genomic and transcriptomic data for discovery of splicing variants in cancer. *bioRxiv*
1035 436634 (2018) doi:10.1101/436634.

1036 51. Li, Y. I. *et al.* Annotation-free quantification of RNA splicing using LeafCutter. *Nat. Genet.* **50**, 151–158 (2018).

1037 52. Storey, J. D. The positive false discovery rate: a Bayesian interpretation and the q-value. *Ann. Stat.* **31**, 2013–2035 (2003).

1038 53. Nicolas, A. *et al.* Genome-wide Analyses Identify KIF5A as a Novel ALS Gene. *Neuron* **97**, 1268–1283.e6 (2018).

1039 54. Oscanoa, J. *et al.* SNPnexus: a web server for functional annotation of human genome sequence variation (2020 update). *Nucleic Acids*
1040 *Res.* **48**, W185–W192 (2020).

1041 55. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS*
1042 *Genet.* **10**, e1004383 (2014).

1043 56. Myers, T. A., Chanock, S. J. & Machiela, M. J. LDlinkR: An R Package for Rapidly Calculating Linkage Disequilibrium Statistics in
1044 Diverse Populations. *Front. Genet.* **11**, 157 (2020).

1045 57. Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association studies. *Nat. Genet.* **48**, 245–252 (2016).

1046 58. Lowy-Gallego, E. *et al.* Variant calling on the GRCh38 assembly with the data from phase three of the 1000 Genomes Project.
1047 *Wellcome Open Res* **4**, 50 (2019).

1048 59. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* **42**, 565–569 (2010).

1049 60. Li, Y. I., Wong, G., Humphrey, J. & Raj, T. Prioritizing Parkinson’s disease genes using population-scale transcriptomic data. *Nat.*
1050 *Commun.* **10**, 994 (2019).

1051 61. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New York, 2009).

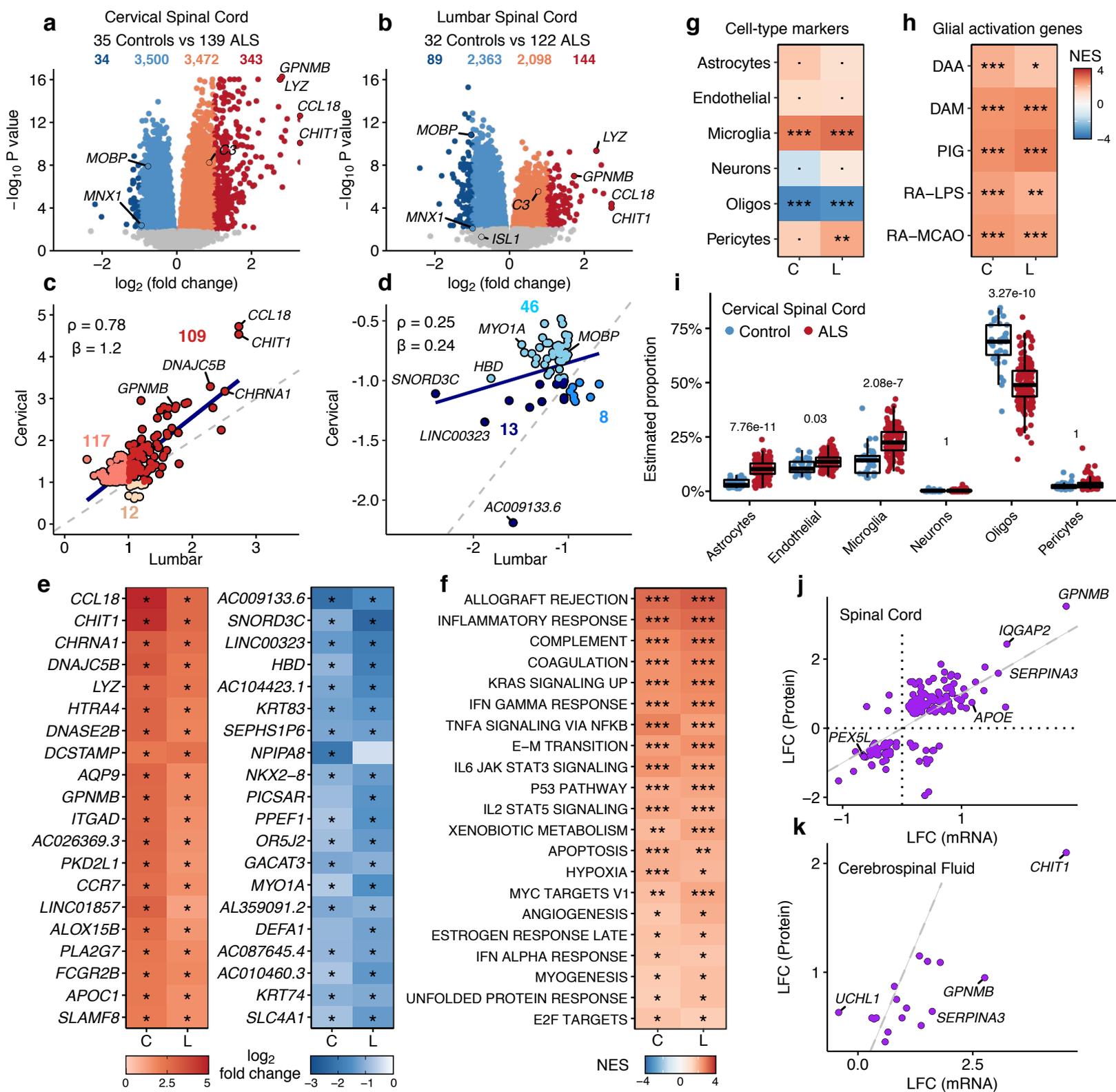
1052 62. Slowikowski, K. ggrepel: Repulsive Text and Label Geoms for ‘ggplot2’, 2016. *R package version 0.5*.

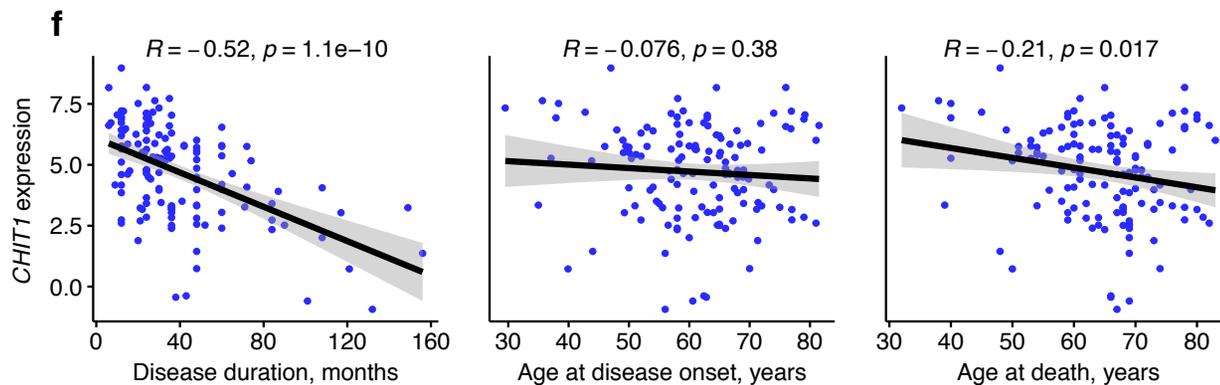
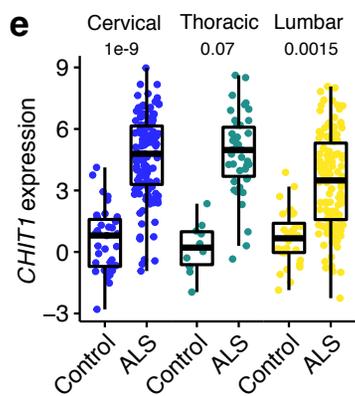
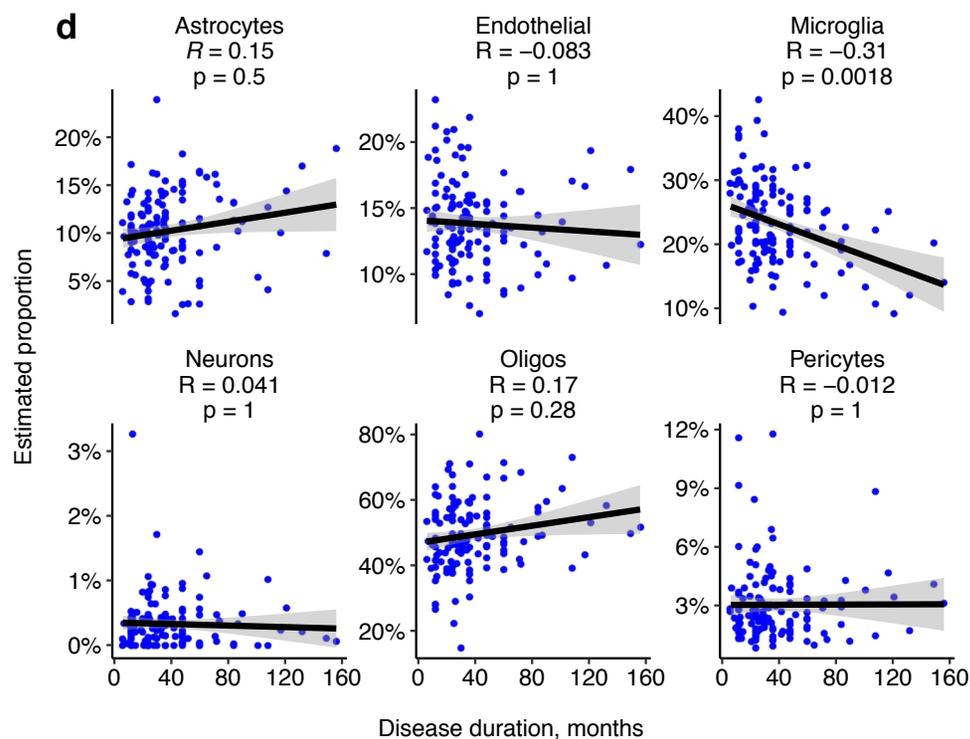
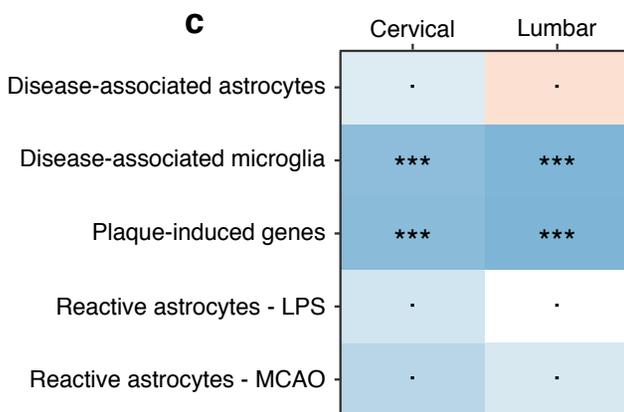
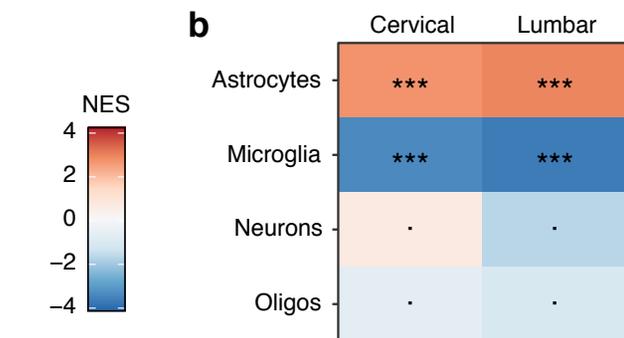
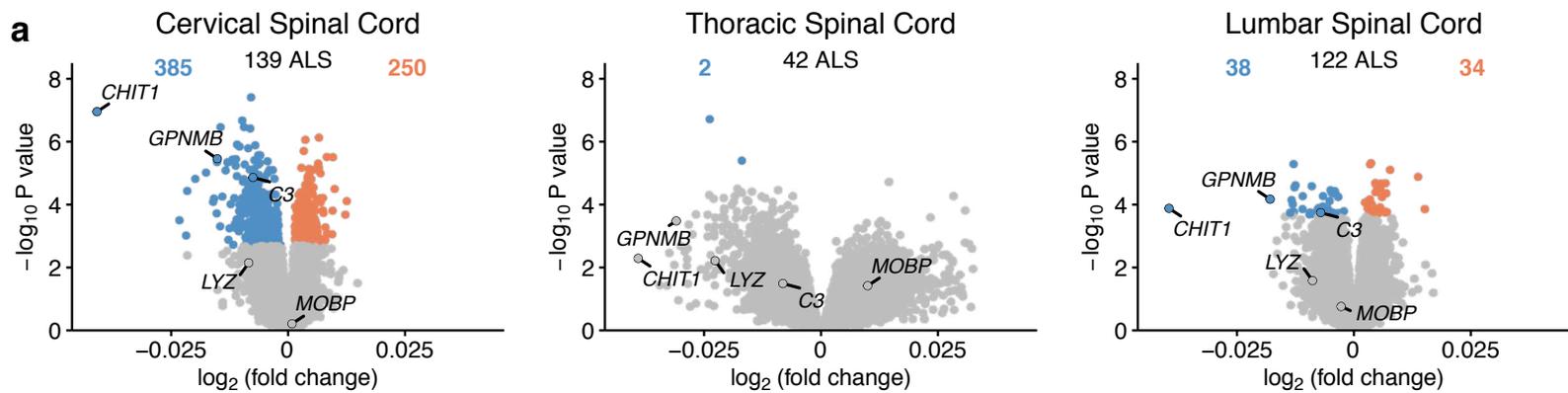
1053 63. Tang, Y., Horikoshi, M. & Li, W. ggfortify: Unified interface to visualize statistical results of popular R packages. *R J.* **8**, 474 (2016).

1054 64. Pedersen, T. L. patchwork: The Composer of Plots. *R package version 1*, 410 (2019).

1055 65. Xu, S. *et al.* Use ggbreak to Effectively Utilize Plotting Space to Deal With Large Datasets and Outliers. *Front. Genet.* **12**, 774846
1056 (2021).

- 1057 66. Yin, T., Cook, D. & Lawrence, M. ggbio: an R package for extending the grammar of graphics for genomic data. *Genome Biol.* **13**, R77
1058 (2012).
1059
1060





a

Tissue	N	eGenes	sGenes
Cervical	216	7890	4730
Lumbar	197	6751	4302
Thoracic	68	962	1387
<i>GTEx Cervical</i>	126	3414	965

