



This is a repository copy of *Exome sequencing in amyotrophic lateral sclerosis implicates a novel gene, DNAJC7, encoding a heat-shock protein.*

White Rose Research Online URL for this paper:

<http://eprints.whiterose.ac.uk/155061/>

Version: Accepted Version

Article:

Farhan, S.M.K., Howrigan, D.P., Abbott, L.E. et al. (32 more authors) (2019) Exome sequencing in amyotrophic lateral sclerosis implicates a novel gene, DNAJC7, encoding a heat-shock protein. *Nature Neuroscience*, 22 (12). pp. 1966-1974. ISSN 1097-6256

<https://doi.org/10.1038/s41593-019-0530-0>

© 2019 The Authors. This is an author-produced version of a paper subsequently published in *Nature Neuroscience*. Uploaded in accordance with the publisher's self-archiving policy.

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

1 **Exome sequencing in amyotrophic lateral sclerosis implicates a novel gene, DNAJC7**
2 **encoding a heat-shock protein**

3 Sali M.K. Farhan^{1-3*}, Daniel P. Howrigan¹⁻³, Liam E. Abbott¹⁻³, Joseph R. Klim⁴, Simon D.
4 Topp⁵, Andrea E. Byrnes¹⁻³, Claire Churchhouse¹⁻³, Hemali Phatnani⁶, Bradley N. Smith⁵,
5 Evadnie Rampersaud⁷, Gang Wu⁷, Joanne Wu⁸, Aleksey Shatunov⁹, Alfredo Iacoangeli^{9,10},
6 Ahmad Al Khleifat⁹, Daniel A. Mordes⁴, Sulagna Ghosh^{3,4}, ALSGENS Consortium, FALS
7 Consortium, Project MinE Consortium, CReATe Consortium, Kevin Eggan^{3,4}, Rosa
8 Rademakers¹¹, Jacob L. McCauley^{12,13}, Rebecca Schüle¹⁴, Stephan Züchner^{12,13}, Michael
9 Benatar⁸, J. Paul Taylor^{15,16}, Michael Nalls^{17,18}, Marc Gotkine¹⁹, Pamela J. Shaw²⁰, Karen E.
10 Morrison²¹, Ammar Al-Chalabi^{9,22}, Bryan Traynor^{17,23}, Christopher E. Shaw^{5,24}, David B.
11 Goldstein²⁵, Matthew B. Harms²⁶, Mark J. Daly¹⁻³, and Benjamin M. Neale^{1-3*}

12 Affiliations:

13 ¹Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General
14 Hospital and Harvard Medical School, Boston, MA, USA.

15 ²Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge,
16 MA, USA.

17 ³Stanley Center for Psychiatric Research, Broad Institute of Harvard and MIT, Cambridge, MA,
18 USA.

19 ⁴Department of Stem Cell and Regenerative Biology, Harvard Stem Cell Institute, Harvard
20 University, Cambridge, MA, USA.

21 ⁵United Kingdom Dementia Research Institute Centre, Maurice Wohl Clinical Neuroscience
22 Institute, Institute of Psychiatry, Psychology and Neuroscience, King's College London, London,
23 U.K.

24 ⁶Center for Genomics of Neurodegenerative Disease, New York Genome Center, New York,
25 NY, USA.

26 ⁷Department of Computational Biology, St. Jude Children’s Research Hospital, Memphis, TN,
27 USA.

28 ⁸Department of Neurology, University of Miami, Miami, FL, USA.

29 ⁹Maurice Wohl Clinical Neuroscience Institute, King’s College London, Department of Basic
30 and Clinical Neuroscience, London, UK.

31 ¹⁰Department of Biostatistics and Health Informatics, Institute of Psychiatry, Psychology and
32 Neuroscience, King’s College London, London, UK.

33 ¹¹Department of Neuroscience, Mayo Clinic, Jacksonville, FL, USA.

34 ¹²John P. Hussman Institute for Human Genomics, University of Miami, Miller School of
35 Medicine, Miami, FL, USA.

36 ¹³The Dr. John T. Macdonald Foundation Department of Human Genetics, University of Miami,
37 Miller School of Medicine, Miami, FL, USA.

38 ¹⁴Center for Neurology and Hertie Institute für Clinical Brain Research, University of Tübingen,
39 German Center for Neurodegenerative Diseases, Tübingen, Germany.

40 ¹⁵Howard Hughes Medical Institute, Chevy Chase, MD, USA.

41 ¹⁶Department of Cell and Molecular Biology, St. Jude Children’s Research Hospital, Memphis,
42 TN, USA.

43 ¹⁷Molecular Genetics Section, Laboratory of Neurogenetics, National Institute on Aging,
44 Bethesda, MD, USA.

45 ¹⁸Data Tecnica International, Glen Echo, MD, USA.

46 ¹⁹Department of Neurology, The Agnes Ginges Center for Human Neurogenetics, Hadassah-
47 Hebrew University Medical Center, Jerusalem, Israel.

48 ²⁰Sheffield Institute for Translational Neuroscience, Sheffield University, UK.

49 ²¹Faculty of Medicine, University of Southampton and Department of Neurology, University
50 Hospital Southampton, UK.

51 ²²Department of Neurology, King's College Hospital, London, UK.

52 ²³Department of Neurology, Johns Hopkins University, Baltimore, MD, USA.

53 ²⁴Centre for Brain Research, University of Auckland, Auckland, New Zealand.

54 ²⁵Institute for Genomic Medicine, Columbia University, New York, NY, USA.

55 ²⁶Department of Neurology, Columbia University, New York, NY, USA.

56 Equal contributions statement: not applicable.

57 *Corresponding authors: sfarhan@broadinstitute.org, bneale@broadinstitute.org

58

59 **ABSTRACT**

60 To discover novel genes underlying amyotrophic lateral sclerosis (ALS), we aggregated exomes
61 from 3,864 cases and 7,839 ancestry matched controls. We observed a significant excess of rare
62 protein-truncating variants among ALS cases, which was concentrated in constrained genes.
63 Through gene level analyses, we replicated known ALS genes including SOD1, NEK1, and FUS.
64 We also observed multiple distinct protein-truncating variants in a highly constrained gene,
65 DNAJC7. The signal in DNAJC7 exceeded genome-wide significance and immunoblotting
66 assays showed depletion of DNAJC7 protein in fibroblasts in an ALS patient carrying the
67 p.Arg156Ter variant. DNAJC7 encodes a member of the heat shock protein family (HSP40),
68 which along with HSP70 proteins, facilitate protein homeostasis including folding of newly

69 synthesized polypeptides and clearance of degraded proteins. When these processes are not
70 regulated, misfolding and accumulation of aberrant proteins can occur leading to protein
71 aggregation, a pathological hallmark of neurodegeneration. Our results highlight DNAJC7 as a
72 novel gene for ALS.

73

74 **KEYWORDS**

75 Amyotrophic lateral sclerosis; protein truncating variants; neurodegeneration; rare variants;
76 DNAJC7.

77

78 **INTRODUCTION**

79 Amyotrophic lateral sclerosis (ALS) is a late-onset neurodegenerative disease characterized
80 primarily by degeneration of motor neurons leading to progressive weakness of limb, bulbar, and
81 respiratory muscles^{1,2}. Genetic variation is an important risk factor for ALS. Given that 5-10% of
82 patients report a positive family history² and ~10% of sporadic patients carry known familial
83 ALS gene mutations, the distinction between familial and sporadic disease is increasingly
84 blurred³. Until recently, ALS gene discoveries were made through large multigenerational
85 pedigrees in which the gene and the causal variant segregated in an autosomal dominant
86 inheritance pattern with very few cases of autosomal recessive inheritance reported. Collecting
87 sporadic case samples has been valuable for gene discovery in more common disorders such as
88 schizophrenia⁴, inflammatory bowel disease⁵, and type 2 diabetes⁶, and can have profound
89 effects on the success of targeted therapeutic approaches^{2,7,8}. The most recent ALS genetic
90 discoveries using large massively parallel sequencing data yielded several gene discoveries

91 including TBK1, TUBA4A, ANXA11 and NEK1 and KIF5A⁹⁻¹³; in addition to other risk loci in
92 C21orf2, MOBP, and SCFD1¹⁴.

93 Herein, we have assembled the largest ALS exome case-control study to date, consisting
94 of 11,703 individuals (3,864 cases and 7,839 controls). We complemented our analysis by
95 leveraging allele frequencies from large external exome sequencing databases in DiscovEHR
96 (>50,000 samples) and a subset of ExAC (>45,000 samples). In our analysis, we observed an
97 excess of rare protein truncating variants in ALS cases, which primarily resided in genes under
98 strong purifying selection and therefore, are less likely to tolerate deleterious mutations
99 (constrained genes). Furthermore, through gene burden testing in which multiple independent
100 variants are harbored in the same gene therefore, implicating that gene in a disease, we
101 confirmed the known association of SOD1, NEK1, and FUS in ALS. Interestingly, we observed
102 multiple, distinct protein-truncating variants in DNAJC7 in our cohort and in an independent,
103 replication cohort. In our analysis, the signal in DNAJC7 exceeded genome-wide significance
104 and immunoblotting showed depletion of DNAJC7 in fibroblasts from an ALS patient carrying
105 the p.Arg156Ter protein truncating variant. DNAJC7 is a highly constrained gene, and encodes a
106 DNAJ molecular chaperone, which facilitates protein maintenance and quality control, such as
107 folding of newly synthesized polypeptides, and clearance of degraded proteins¹⁵. Dysregulation
108 of these processes can lead to aberrant protein aggregation, one of the pathological hallmarks of
109 neurodegenerative diseases.

110

111 **RESULTS**

112 **Patient demographics and dataset overview**

113 We processed our initial dataset of 15,722 samples through a rigorous quality control pipeline
114 using Hail, an open-source, scalable framework for exploring and analyzing genomic data
115 <https://hail.is/>. All samples were screened for the C9orf72 hexanucleotide expansion (G4C2) and
116 positive samples were excluded from our study. We removed samples with poor sequencing
117 quality, high levels of sequence contamination, closely related with one another, ambiguous sex
118 status, or population outliers per PCA (Supplementary Table 1; Supplementary Fig. 1-2). Our
119 final data set consisted of 3,864 cases and 7,839 controls for a total of 11,703 samples.
120 Individuals were of European descent with 7,355 (62.8%) and 4,348 (37.2%) of samples
121 classified as males and females, respectively. Of 3,864 cases, 2,274 (58.9%) and 1,590 (41.1%)
122 samples were classified as males and females, respectively; where 5,081 (64.8%) and 2,758
123 (35.2%) were classified as males in controls.

124

125 **Excess of exome-wide rare protein truncating variants**

126 We assessed four models that incorporated different covariates and assessed their stringency and
127 performance by controlling for benign or synonymous variation. Specifically, each model uses
128 firth based logistic regression and incorporates some or all the covariates: 1) sample sex, 2) PC1-
129 PC10, and either 3) the total exome count (summation of synonymous variants, benign missense
130 variants, damaging missense variants, and protein-truncating variants) or 4) benign variation
131 (summation of synonymous and benign missense variants). We show the results from the most
132 conservative model (model 3), which used all the covariates and the total exome count. Under
133 these models, we evaluated four classes of allele frequency thresholds: (1) singletons, which are
134 variants present in a single individual in our dataset (allele count, AC = 1); (2) doubletons, which
135 are present in two individuals in our dataset (AC = 2); (3) ultra-rare singletons, which are

136 singletons in our dataset and are absent in DiscovEHR, a large, independent exome dataset (AC
137 = 1, 0 in DiscovEHR); and finally, (4) rare variants, which have an allele frequency of of <0.01%
138 in our dataset (11,703 samples), in ExAC (non-psychiatric studies, >45,000 samples) and in
139 DiscovEHR (>50,000 samples). For a full explanation of these models and allele frequency
140 thresholds, please see the Methods section.

141 Using model 3, we observed a significant enrichment of singleton protein-truncating
142 variants in ALS cases relative to controls (OR: 1.07, P: 5.00×10^{-7}); ultra-rare singleton PTVs
143 (OR: 1.08, P: 1.97×10^{-6}); and rare PTVs (OR: 1.04, P: 1.77×10^{-7}) (Fig. 1). These values all
144 passed multiple test correction ($P < 0.0125$). The number of doubletons (AC=2) was too low to
145 detect any significant enrichment.

146 When using model 4 where we restrict to ‘benign variation’ as the final covariate, the
147 protein-truncating variants signal is further enriched among singletons (OR: 1.12, P: $< 2 \times 10^{-16}$);
148 ultra-rare singletons (OR: 1.10, P: 1.53×10^{-10}); and rare variants (OR: 1.04, P: 1.47×10^{-7}).
149 Interestingly, in this analysis, there is a consistent and a significant enrichment of damaging
150 missense variants not observed in the previous analysis: singletons (OR: 1.06, P: $< 2 \times 10^{-16}$);
151 ultra-rare singletons (OR: 1.03, P: 6.33×10^{-5}); and rare variants (OR: 1.01, P: 3.24×10^{-3}).

152 In our analyses, we use a standard definition of protein-truncating variants as frameshift
153 variants, splice acceptor variants, splice donor variants, or stop gained variants, which are due to
154 insertions or deletions (indels), or single nucleotide variants (SNVs). Given the known elevated
155 error rate in indels we divided all protein-truncating variants as either SNVs or indels and
156 repeated the exome-wide analysis to eliminate any false positive signals. The significant signal is
157 present in both SNVs and indels: SNV singletons (OR: 1.05, P: 2.99×10^{-3}); indel singletons
158 (OR: 1.10, P: 5.75×10^{-6}); SNV ultra-rare singletons (OR: 1.06, P: 4.34×10^{-3}); indel ultra-rare

159 singletons (OR: 1.12, P: 1.96×10^{-5}); and SNV rare variants (OR: 1.03, P: 6.48×10^{-4}); indel rare
160 variants (OR: 1.05, P: 3.30×10^{-5}) (Supplementary Fig. 4). This additional quality control test
161 ensures that the protein-truncating variants signal is driven by both indels and SNVs and is
162 unlikely to be false.

163

164 **Gene set testing: enrichment of rare variants in constrained genes**

165 To determine whether we could identify the source of the protein-truncating variants enrichment,
166 we assessed multiple different gene sets. We evaluated: (1) constrained genes, which are a set of
167 genes under strong purifying selection; (2) genes known to confer risk to ALS; (3) genes
168 associated with clinically overlapping diseases such as other motor neuron diseases (primary
169 lateral sclerosis, progressive muscular atrophy, progressive bulbar palsy, and spinal muscular
170 atrophy) as well as genes associated with frontotemporal dementia, Parkinson's disease, Pick's
171 disease, and Alzheimer's disease; and finally, (4) genes in which their expression is specific to
172 the brain.

173 Among constrained genes we observed a significant enrichment of singleton protein-
174 truncating variants (OR: 1.23, P: 7.74×10^{-7}); ultra-rare singletons (OR: 1.27, P: 5.76×10^{-8}), and
175 rare variants (OR: 1.33, P: $< 2 \times 10^{-16}$) (Fig. 2A, Supplementary Fig. 5A). We obtained similar
176 results using model 4 (Supplementary Fig. 5A). To determine whether the entire signal can be
177 explained by constrained genes, we removed them genes and reconducted the analysis. The
178 significant enrichment signal persists however, the effect sizes are attenuated: singleton protein-
179 truncating variants (OR: 1.05, P: 3.30×10^{-4}); ultra-rare singleton protein-truncating variants
180 (OR: 1.05, P: 1.96×10^{-3}); and rare protein-truncating variants (OR: 1.02, P: 2.93×10^{-3}) (Fig. 2B,

181 Supplementary Fig. 5B). This enrichment was also observed in model 4 (Supplementary Fig.
182 5B).

183 Next, we evaluated the potential effects of known ALS genes. We did not include the
184 ALS genes TBK1, NEK1, KIF5A, C21orf2, MOBP, or SCFD1 as these genes were discovered
185 using datasets that contained a large subset of the same samples and can generate an amplified
186 signal. The known ALS genes had negligible, insignificant effects (Fig. 3A, Supplementary Fig.
187 6). When including variants from TBK1, NEK1, KIF5A, C21orf2, MOBP, or SCFD1, the
188 negligible signals persist therefore, the initial observation of the exome-wide protein-truncating
189 variant enrichment is not driven by known effects of ALS genes and is likely due to other
190 genomic loci.

191 Although ALS is traditionally considered to be a disease of upper and lower motor
192 neurons, more than 50% of ALS patients exhibit neuropsychological and cognitive deficits, with
193 up to 30% of ALS patients meeting some diagnostic criteria for frontotemporal dementia, and
194 some patients may also exhibit Parkinsonism or Parkinsonism-dementia^{1,16-20}. We tabulated a list
195 of genes associated with other motor neuron diseases such as primary lateral sclerosis,
196 progressive muscular atrophy, progressive bulbar palsy, and spinal muscular atrophy. We also
197 included genes associated with frontotemporal dementia, Parkinson's disease, Pick's disease, and
198 Alzheimer's disease (Supplementary Table 5). We did not observe a significant enrichment of
199 variants in any class of variation, suggesting that the initial observation of protein-truncating
200 variant enrichment is unlikely to be explained by only these genes (Fig. 3B, Supplementary Fig.
201 7).

202 Finally, we tested whether there is a signal in brain specific genes as ALS is a
203 neurodegenerative disease with the predominant symptoms affecting the central nervous system.

204 We extracted a list of genes with specific brain expression generated using GTEx and performed
205 the same burden analysis across classes of variation. We did not observe any significant
206 differences in protein-truncating variants or damaging missense variation in any allele frequency
207 threshold (Fig. 3C, Supplementary Fig. 8).

208

209 **Single gene burden analysis replicates previous ALS associations**

210 To determine whether a single gene is enriched for variation in ALS cases (ALS-associated) or
211 depleted in ALS cases (ALS-protective), we evaluated ultra-rare (AC=1, absent in DiscovEHR)
212 and rare (MAF <0.001% in our dataset, DiscovEHR, and ExAC) protein-truncating variants and
213 damaging missense variants. Within the ultra-rare variant category, no individual gene passed
214 exome-wide significance. However, the top genes were known ALS genes: (1) NEK1 (PTVs,
215 OR: 12.21, P: 7.32×10^{-5}); (2) OPTN (PTVs, OR: 20.33, P: 1.2×10^{-4}); and (3) SOD1 (dmis, OR:
216 46.91, P: 5.03×10^{-6}) (Supplementary Fig. 9). Within rare protein-truncating variants, only NEK1
217 (OR: 12.8, P: 4.59×10^{-9}), passed exome-wide significance; the next top 9 most significant genes,
218 which include FUS, a known ALS gene (OR: 26.4, P: 1.29×10^{-3}), are displayed in Table 1, Fig.
219 4A. Similarly, within damaging missense variants, SOD1 (OR: 87.7, P: 7.5×10^{-11}) was the only
220 gene to pass exome-wide significance; the top 9 most significant genes are displayed in Table 1,
221 Fig. 4B. In Supplementary Tables 2 and 3, we tabulate the results of the single gene burden
222 analysis for the proposed ALS genes based on the literature, as well as their odds ratio and P-
223 values.

224 To determine if we can reproduce the initial signals observed, we included an additional
225 21,071 controls from ExAC that are of European descent (non-Finnish) and were not a part of
226 any psychiatric or brain related studies, to eliminate any sample overlap. We performed the same

227 burden analyses using 3,864 cases and 28,910 controls (7,839 controls within our dataset and
228 21,071 additional controls). In Tables 1 and 2, we display the most significant genes that were
229 identified in the initial discovery and tabulate their OR and P-values for both the initial discovery
230 cohort (3,864 cases and 7,839 controls) and the secondary analysis (3,864 cases and 28,910
231 controls). Within protein-truncating variants, NEK1 is still the only gene that exceeds exome-
232 wide significance (OR: 6.5, P: 3.03×10^{-10}) (Fig. 4C). Of the next 9 most significant genes in the
233 initial analysis, the only signal that was strengthened was in FUS (OR: 97.4, P: 2.68×10^{-6}). This
234 finding suggests that the other genes may not be true positives or will need further evidence to
235 support their association with ALS. Interestingly, the signal in OPTN, a proposed ALS
236 associated gene, decreased (OR: 6.6, P: 3.0×10^{-3} to OR: 2.6, P: 6.9×10^{-3}) however, this may be
237 explained in part by the observation that OPTN protein-truncating variants tend to manifest as a
238 recessive form of ALS, which may not be detected in our burden model. With the additional
239 controls, multiple genes had similar ORs as the discovery analysis, with their respective P-values
240 approaching significance (P-values ranging from 7.7×10^{-5} - 1.4×10^{-3}). Most notably, the signal
241 in TBK1, a proposed ALS gene based on Cirulli et al. strengthened: (initial analysis; OR: 22.3, P:
242 3.9×10^{-3} ; secondary analysis: OR: 12.5, P: 9.35×10^{-4}). Within damaging missense variants,
243 SOD1 is still the only gene that exceeds exome-wide significance (OR: 79.0, P: 6.0×10^{-18});
244 however, the next 9 most significant genes no longer approach statistical significance. Similarly,
245 when integrating additional controls, multiple genes approach significance (P-values ranging
246 from 1.2×10^{-4} - 6.2×10^{-4}) (Fig. 4D).

247

248 **Loss of function variants in DNAJC7 in ALS patients**

249 DNAJC7, which is a highly constrained gene (pLI = 0.99) had 4 protein-truncating variants
250 carriers in cases (3,864) and 0 in controls (7,839) in the discovery analysis (OR: 18.3, P: 0.01);
251 and 0 protein-truncating variants in total controls (28,910) (OR: 96.1, P: 1.9×10^{-4}). While
252 DNAJC7 did not initially exceed genome-wide significance, its high constraint score and role in
253 neurodegeneration as a member of the heat shock protein 40 (HSP40) family, encouraged us to
254 evaluate additional datasets to determine its loss of function mutation frequency.

255 We surveyed data from the UK Motor Neurone Disease Association (n=1,135) and The
256 Agnes Ginges Center for Human Neurogenetics at the Hadassah-Hebrew University Medical
257 Center in Israel (n=96). We observed an additional 4 carriers for a total of 6 distinct protein-
258 truncating variants in 8 individuals with ALS (cases: 5,095; controls: 28,910; OR: 96.6, P:
259 2.5×10^{-7}) (Table 2). These DNAJC7 variants are extremely rare or completely absent from large
260 population datasets such as gnomAD (Table 2). The DNAJC7 p.Phe163fs variant was observed
261 in the Israeli cohort. As gnomAD does not currently provide variant frequency on individuals of
262 Middle Eastern ethnicity, we screened an additional 3,244 controls from a mixture of Middle
263 Eastern ethnicities for the p.Phe163fs variant and did not observe any carriers further
264 demonstrating its rarity in the general population and an ancestry matched population. In
265 addition, we also observed 15 rare missense variants in DNAJC7, of which 4 are predicted to
266 exert a damaging effect in 5 ALS cases and 1 in control (Table 2).

267 We next proceeded to ask if any of the protein-truncating variants in DNAJC7 can affect
268 its mRNA or protein levels. Accordingly, we collected total RNA from human fibroblasts
269 derived from healthy controls and a patient with a DNAJC7 protein-truncating variant
270 p.Arg156Ter and performed qRT-PCR with two different sets of primer pairs to investigate
271 DNAJC7 transcript levels (Supplementary Fig. 10A and B). These data indicate that DNAJC7

272 mRNA abundance is not significantly altered in fibroblasts harboring a DNAJC7 protein-
273 truncating variant (Fig. 5A). We next carried out immunoblot assays on protein lysates from
274 fibroblasts and determined that DNAJC7 protein levels were significantly reduced in the ALS
275 patient fibroblasts (Fig. 5B). Although this protein-truncating variant could potentially yield a
276 17.5 kDa protein, no evidence for such a product was detected (Supplementary Fig. 10C).
277 Together, our findings indicate the protein-truncating variants we identified in DNAJC7 leads to
278 decreased protein levels of this heat shock protein co-chaperone.

279

280 **DISCUSSION**

281 Herein, we have assembled the exomes of 3,864 ALS cases and 7,839 controls and observed an
282 exome-wide enrichment of protein-truncating variants, which typically result in protein loss-of-
283 function. The abundance of protein-truncating variants in ALS cases seems to be primarily
284 driven by constrained genes, which are under strong purifying selection. When removing
285 constrained genes, the initial exome-wide enrichment of protein-truncating variants remains;
286 however, the effect sizes are much smaller, suggesting that while constraint genes may explain
287 much of protein-truncating variant enrichment, there may be minor residual effects elsewhere in
288 the genome. Accordingly, we examined the effects of ALS associated genes and did not observe
289 any significant enrichment. Importantly, a subset of cases was pre-screened for known
290 pathogenic variants in a select number of known ALS genes and positive cases were eliminated
291 prior to assembling the dataset, which attenuated the effect size estimates and significance for
292 genes in this gene set.

293 Acknowledging the phenotypic variability of ALS, we also evaluated the effects of genes
294 implicated in other motor neuron diseases such as primary lateral sclerosis, progressive muscular

295 atrophy, progressive bulbar palsy, and spinal muscular atrophy; as well as genes associated with
296 frontotemporal dementia, Parkinson's disease, Pick's disease, and Alzheimer's disease. We did
297 not observe a significant enrichment in any class of variation, suggesting that the initial
298 observation of excess protein-truncating variants do not reside in these genes. Lastly, the genes
299 implicated in the development of ALS are not specifically expressed in motor neurons, nor are
300 they brain specific, despite the specific degree of degeneration of upper and lower motor
301 neurons. Nevertheless, we tested whether the signal in protein-truncating variants is concentrated
302 in brain specific genes, a much larger gene set than ALS genes only. We did not observe any
303 significant enrichment within brain specific genes.

304 The single gene burden analysis identified the most significant genes as SOD1, NEK1,
305 and FUS, which are known ALS genes. No other individual gene passed exome-wide
306 significance within our dataset (3,864 cases and 7,839 controls) and the additional controls in the
307 secondary analysis (3,864 cases and 28,910 controls). Notably, in the secondary analysis,
308 multiple genes with consistent OR and lower P-values than the initial analysis, surfaced. Within
309 protein-truncating variants, these include: GRIN3B, HRCT1, IL3, and DNAJC7. Interestingly,
310 protein-truncating variants in GRIN3B and HRCT1 may offer protection against ALS: OR: 0.05,
311 P: 7.7×10^{-5} ; OR: 0.05, P: 1.2×10^{-4} , respectively; while protein-truncating variants in IL3 and
312 DNAJC7 may confer risk: OR: 10.5, P: 1.8×10^{-4} ; OR: 67.4, P: 1.9×10^{-4}).

313 In this analysis, DNAJC7 had 4 protein-truncating variant carriers in 3,864 cases and 0 in
314 7,839 and 28,910 controls additionally, when integrating data from the UK Motor Neurone
315 Disease Association, we observed an additional 4 protein-truncating variant carriers for a total of
316 6 distinct protein-truncating variants in 8 individuals (initial analysis P: 0.01; secondary analysis
317 P: 1.9×10^{-4} ; replication analysis P: 2.5×10^{-7}). According to the HPA RNA-seq normal tissues

318 project²¹ and the Genotype-Tissue Expression (GTEx) project²², DNAJC7 is ubiquitously
319 expressed with elevated expression in the brain. DNAJC7 encodes a molecular chaperone, DnaJ
320 heat shock protein family (HSP40) member C7, and like all DNAJ proteins, contains an
321 approximately 70 amino acid J-domain, which is critical for binding to HSP70 proteins²³. There
322 are approximately 50 DNAJ proteins, which are also classified as HSP40 proteins, that facilitate
323 protein maintenance and quality control, such as folding of newly synthesized polypeptides, and
324 clearance of degraded proteins^{15,24,25}. Specifically, DNAJs act as co-chaperones for HSP70
325 proteins by regulating ATPase activity, aid in polypeptide binding, and prevention of premature
326 polypeptide folding^{25,26}.

327 Aberrant protein aggregation due to accumulation of misfolded proteins, is one of the
328 pathological hallmarks of neurodegenerative diseases like Alzheimer's disease, Parkinson's
329 disease, Huntington's disease, prion disease, and ALS²⁷⁻³². HSP proteins have a conserved and
330 central role in protein function by aiding in their folding and stabilization, and the clearance of
331 misfolded proteins, ultimately diminishing protein aggregates and the associated pathologies.
332 However, genetic aberrations or cellular stress such as exposure to environmental toxins,
333 fluctuations in temperature, chemical stress, cell injury, or aging, can influence the dynamics of
334 the protein quality control network allowing misfolded proteins to go undetected thereby
335 triggering neurotoxicity^{33,34}. Furthermore, abnormal expression of HSP70 and DNAJ genes leads
336 to the formation of protein aggregates in models of Alzheimer's disease³⁵, Parkinson's
337 disease^{36,37}, Huntington's disease^{35,38}, prion disease^{39,40}, and ALS⁴¹⁻⁴³. In light of these studies,
338 elevated HSP expression is thought to be beneficial in preventing or in halting neurodegenerative
339 disease progression⁴⁴. For example, overexpression of DNAJB6b and DNAJB8 suppressed toxic
340 protein aggregation⁴⁵; while overexpression of HSP70 in neuroglioma cells decreased the

341 formation of alpha-synuclein fibrils⁴⁶. Within ALS models, overexpression of HSPB8 promoted
342 clearance of mutant SOD1⁴⁷; double transgenic mice overexpressing HSP27 and mutated SOD1
343 exhibited increased survival of spinal motor neurons than mice overexpressing a SOD1 mutation
344 only, however, the neuroprotective effects were not sustained in later stages of the disease⁴⁸.
345 Finally, DNAJB2, which when mutated can cause autosomal recessive spinal muscular atrophy,
346 was overexpressed in mice motor neurons also expressing a SOD1 mutation (p.Gly93Ala), and
347 led to reduced mutant SOD1 aggregation and improved motor neuron survival⁴⁹. In
348 Supplementary Table 4, we tabulated additional HSP genes that have been reported to harbor
349 pathogenic or likely pathogenic mutations in patients with neurodegenerative diseases.

350 In summary, we observed a significant exome-wide enrichment of protein-truncating
351 variants, which seem to primarily reside in constrained genes. Through gene burden tests, we
352 confirmed the known association of ALS genes SOD1, NEK1, and FUS, and also observed
353 multiple protein-truncating variants in ALS cases in a highly constrained, HSP40 gene, DNAJC7.
354 Our replication of protein-truncating variants in DNAJC7 in an independent ALS cohort as well
355 as functional validation highlights loss of DNAJC7 as a novel genetic risk factor for ALS.

356

357 **ACKNOWLEDGEMENTS**

358 We thank and acknowledge the consent and cooperation of all study participants. Many thanks to
359 F. Cerrato for helping us assemble the dataset and providing general project management; and to
360 T. Poterba, J. Bloom, D. King, and C. Seed for their assistance in Hail. Data used in this research
361 were in part obtained from the UK MND Collections for MND Research, funded by the MND
362 Association and the Wellcome Trust. We would like to thank people with MND and their
363 families for their participation in this project. The project is supported through the following

364 funding organisations under the aegis of JPND - www.jpnd.eu (United Kingdom, Medical
365 Research Council (MR/L501529/1; MR/R024804/1) and Economic and Social Research Council
366 (ES/L008238/1)) and through the Motor Neurone Disease Association. This study represents
367 independent research part funded by the National Institute for Health Research (NIHR)
368 Biomedical Research Centre at South London and Maudsley NHS Foundation Trust and King's
369 College London. Samples used in this research were in part obtained from the UK National DNA
370 Bank for MND Research, funded by the MND Association and the Wellcome Trust. We
371 acknowledge sample management undertaken by Biobanking Solutions funded by the Medical
372 Research Council at the Centre for Integrated Genomic Medical Research, University of
373 Manchester. The CReATe consortium (U54NS092091) is part of Rare Diseases Clinical
374 Research Network (RDCRN), an initiative of the Office of Rare Diseases Research (ORDR),
375 NCATS. This consortium is funded through collaboration between NCATS, and the NINDS.
376 Additional support is provided by the ALS Association (17-LGCA-331). S.M.K. Farhan is
377 supported by the ALS Canada Tim E. Noël Postdoctoral Fellowship. J.R. Klim was supported by
378 the Project ALS Tom Kirchhoff Family Postdoctoral Fellowship and acknowledges K. Mamia
379 and L.T. Kane for their work banking fibroblasts.

380

381 **AUTHOR CONTRIBUTIONS**

382 S.M.K.F., M.J.D., and B.M.N. conceived and designed the experiments. S.M.K.F., S.D.T., H.P.,
383 B.N.S., E.R., G.W., J.W., A.S., A.I., A.A.K., D.A.M., S.G., A.G., K.E., R.R., J.L.M., R.S., S.Z.,
384 M.B., J.P.T., M.N., M.G., P.J.S., K.E.M., A.A.C., B.T., C.E.S., D.B.G., M.B.H., and B.M.N.
385 collected samples, prepared samples for analysis, or were involved in clinical evaluation. M.B.
386 and J.P.T. were the lead contacts for the CReATe Consortium. S.D.T. and C.E.S. were the lead

387 contacts for the FALS Consortium. D.B.G. and M.B.H. were the lead contacts for the ALSGENS
388 Consortium. S.M.K.F. performed all experiments and executed data analyses. D.P.H., L.E.A.,
389 A.E.B., and S.D.T. provided analysis suggestions. J.R.K. completed the cell culture, RNA, and
390 protein analyses. S.M.K.F. performed the primary writing of the manuscript with input from
391 D.P.H., C.C., M.J.D., and B.M.N. All authors approved the final manuscript. M.J.D. and B.M.N.
392 supervised the research.

393

394 **COMPETING INTERESTS**

395 MN participation is supported by a consulting contract between Data Tecnica International and
396 the National Institute on Aging, NIH, Bethesda, MD, USA, as a possible conflict of interest. MN
397 also consults for Lysosomal Therapeutics Inc, the Michael J. Fox Foundation and Vivid
398 Genomics among others. The other authors declare no competing interests.

399

400 **REFERENCES**

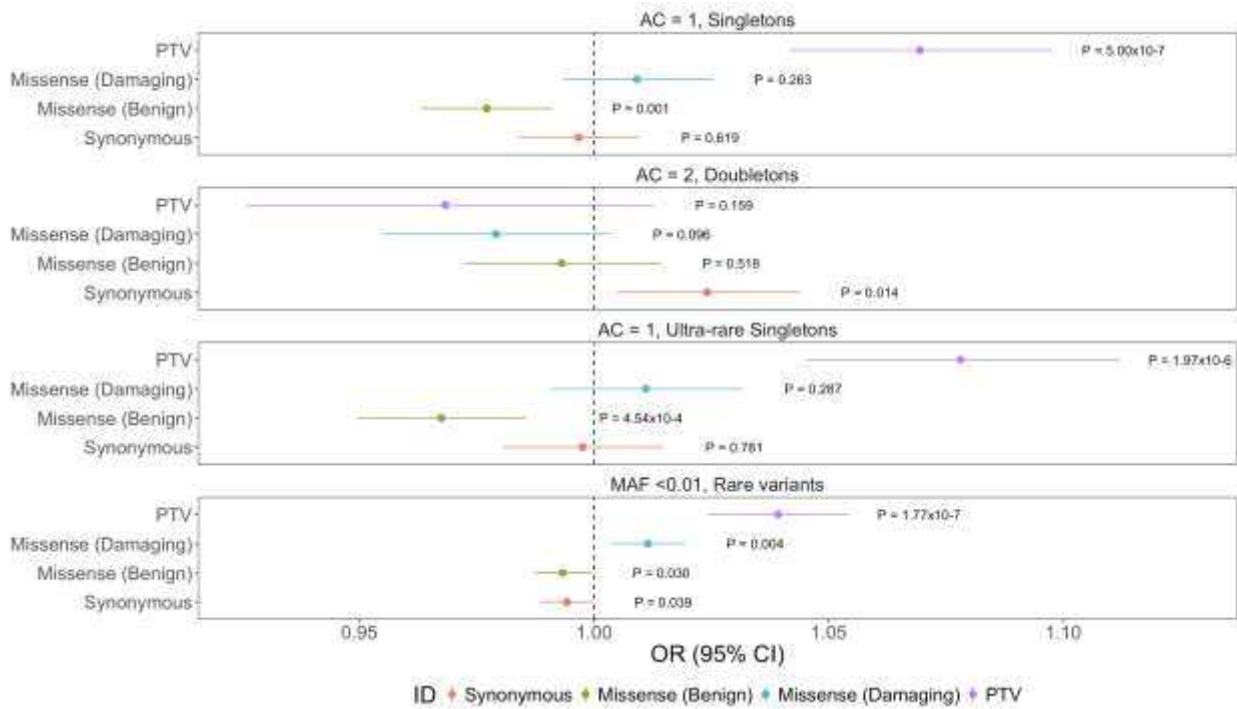
- 401 1. Strong, M.J. et al. Amyotrophic lateral sclerosis - frontotemporal spectrum disorder
402 (ALS-FTSD): Revised diagnostic criteria. *Amyotroph Lateral Scler Frontotemporal*
403 *Degener* **18**, 153-174 (2017).
- 404 2. Al-Chalabi, A., van den Berg, L.H. & Veldink, J. Gene discovery in amyotrophic lateral
405 sclerosis: implications for clinical management. *Nat Rev Neurol* **13**, 96-104 (2017).
- 406 3. Al-Chalabi, A. Perspective: Don't keep it in the family. *Nature* **550**, S112 (2017).
- 407 4. Singh, T. et al. Rare loss-of-function variants in SETD1A are associated with
408 schizophrenia and developmental disorders. *Nat Neurosci* **19**, 571-7 (2016).
- 409 5. Mohanan, V. et al. C1orf106 is a colitis risk gene that regulates stability of epithelial
410 adherens junctions. *Science* **359**, 1161-1166 (2018).
- 411 6. Manning, A. et al. A Low-Frequency Inactivating AKT2 Variant Enriched in the Finnish
412 Population Is Associated With Fasting Insulin Levels and Type 2 Diabetes Risk. *Diabetes*
413 **66**, 2019-2032 (2017).

- 414 7. Hamburg, M.A. & Collins, F.S. The path to personalized medicine. *N Engl J Med* **363**,
415 301-4 (2010).
- 416 8. Nelson, M.R. et al. The support of human genetic evidence for approved drug
417 indications. *Nat Genet* **47**, 856-60 (2015).
- 418 9. Cirulli, E.T. et al. Exome sequencing in amyotrophic lateral sclerosis identifies risk genes
419 and pathways. *Science* **347**, 1436-41 (2015).
- 420 10. Smith, B.N. et al. Exome-wide rare variant analysis identifies TUBA4A mutations
421 associated with familial ALS. *Neuron* **84**, 324-31 (2014).
- 422 11. Smith, B.N. et al. Mutations in the vesicular trafficking protein annexin A11 are
423 associated with amyotrophic lateral sclerosis. *Sci Transl Med* **9**(2017).
- 424 12. Kenna, K.P. et al. NEK1 variants confer susceptibility to amyotrophic lateral sclerosis.
425 *Nat Genet* **48**, 1037-42 (2016).
- 426 13. Nicolas, A. et al. Genome-wide Analyses Identify KIF5A as a Novel ALS Gene. *Neuron*
427 **97**, 1268-1283 e6 (2018).
- 428 14. van Rheenen, W. et al. Genome-wide association analyses identify new risk variants and
429 the genetic architecture of amyotrophic lateral sclerosis. *Nat Genet* **48**, 1043-8 (2016).
- 430 15. Lackie, R.E. et al. The Hsp70/Hsp90 Chaperone Machinery in Neurodegenerative
431 Diseases. *Front Neurosci* **11**, 254 (2017).
- 432 16. Swinnen, B. & Robberecht, W. The phenotypic variability of amyotrophic lateral
433 sclerosis. *Nat Rev Neurol* **10**, 661-70 (2014).
- 434 17. Farhan, S.M. et al. The Ontario Neurodegenerative Disease Research Initiative (ONDRI).
435 *Can J Neurol Sci* **44**, 196-202 (2017).
- 436 18. Farhan, S.M.K., Gendron, T.F., Petrucelli, L., Hegele, R.A. & Strong, M.J. OPTN
437 p.Met468Arg and ATXN2 intermediate length polyQ extension in families with C9orf72
438 mediated amyotrophic lateral sclerosis and frontotemporal dementia. *Am J Med Genet B*
439 *Neuropsychiatr Genet* **177**, 75-85 (2018).
- 440 19. Aarsland, D., Zaccai, J. & Brayne, C. A systematic review of prevalence studies of
441 dementia in Parkinson's disease. *Mov Disord* **20**, 1255-63 (2005).
- 442 20. Hely, M.A., Reid, W.G., Adena, M.A., Halliday, G.M. & Morris, J.G. The Sydney
443 multicenter study of Parkinson's disease: the inevitability of dementia at 20 years. *Mov*
444 *Disord* **23**, 837-44 (2008).
- 445 21. Fagerberg, L. et al. Analysis of the human tissue-specific expression by genome-wide
446 integration of transcriptomics and antibody-based proteomics. *Mol Cell Proteomics* **13**,
447 397-406 (2014).

- 448 22. Consortium, G.T. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**, 580-5
449 (2013).
- 450 23. Jiang, J. et al. Structural basis of J cochaperone binding and regulation of Hsp70. *Mol*
451 *Cell* **28**, 422-33 (2007).
- 452 24. Kampinga, H.H. & Craig, E.A. The HSP70 chaperone machinery: J proteins as drivers of
453 functional specificity. *Nat Rev Mol Cell Biol* **11**, 579-92 (2010).
- 454 25. Mayer, M.P. & Bukau, B. Hsp70 chaperones: cellular functions and molecular
455 mechanism. *Cell Mol Life Sci* **62**, 670-84 (2005).
- 456 26. Clerico, E.M., Tilitsky, J.M., Meng, W. & Gierasch, L.M. How hsp70 molecular
457 machines interact with their substrates to mediate diverse physiological functions. *J Mol*
458 *Biol* **427**, 1575-88 (2015).
- 459 27. Uddin, M.S. et al. Autophagy and Alzheimer's Disease: From Molecular Mechanisms to
460 Therapeutic Implications. *Front Aging Neurosci* **10**, 04 (2018).
- 461 28. Irwin, D.J., Lee, V.M. & Trojanowski, J.Q. Parkinson's disease dementia: convergence of
462 alpha-synuclein, tau and amyloid-beta pathologies. *Nat Rev Neurosci* **14**, 626-36 (2013).
- 463 29. Imarisio, S. et al. Huntington's disease: from pathology and genetics to potential
464 therapies. *Biochem J* **412**, 191-209 (2008).
- 465 30. Brundin, P., Melki, R. & Kopito, R. Prion-like transmission of protein aggregates in
466 neurodegenerative diseases. *Nat Rev Mol Cell Biol* **11**, 301-7 (2010).
- 467 31. Ross, C.A. & Poirier, M.A. Protein aggregation and neurodegenerative disease. *Nat Med*
468 **10 Suppl**, S10-7 (2004).
- 469 32. Winklhofer, K.F., Tatzelt, J. & Haass, C. The two faces of protein misfolding: gain- and
470 loss-of-function in neurodegenerative diseases. *EMBO J* **27**, 336-49 (2008).
- 471 33. Gidalevitz, T., Ben-Zvi, A., Ho, K.H., Brignull, H.R. & Morimoto, R.I. Progressive
472 disruption of cellular protein folding in models of polyglutamine diseases. *Science* **311**,
473 1471-4 (2006).
- 474 34. Voisine, C., Pedersen, J.S. & Morimoto, R.I. Chaperone networks: tipping the balance in
475 protein folding diseases. *Neurobiol Dis* **40**, 12-20 (2010).
- 476 35. Brehme, M. et al. A chaperome subnetwork safeguards proteostasis in aging and
477 neurodegenerative disease. *Cell Rep* **9**, 1135-50 (2014).
- 478 36. Roodveldt, C. et al. Chaperone proteostasis in Parkinson's disease: stabilization of the
479 Hsp70/alpha-synuclein complex by Hip. *EMBO J* **28**, 3758-70 (2009).

- 480 37. Auluck, P.K., Chan, H.Y., Trojanowski, J.Q., Lee, V.M. & Bonini, N.M. Chaperone
481 suppression of alpha-synuclein toxicity in a Drosophila model for Parkinson's disease.
482 *Science* **295**, 865-8 (2002).
- 483 38. Wacker, J.L. et al. Loss of Hsp70 exacerbates pathogenesis but not levels of fibrillar
484 aggregates in a mouse model of Huntington's disease. *J Neurosci* **29**, 9104-14 (2009).
- 485 39. Kovacs, G.G. et al. Prominent stress response of Purkinje cells in Creutzfeldt-Jakob
486 disease. *Neurobiol Dis* **8**, 881-9 (2001).
- 487 40. Jones, G., Song, Y., Chung, S. & Masison, D.C. Propagation of *Saccharomyces*
488 *cerevisiae* [PSI⁺] prion is impaired by factors that regulate Hsp70 substrate binding. *Mol*
489 *Cell Biol* **24**, 3928-37 (2004).
- 490 41. Chen, H.J. et al. The heat shock response plays an important role in TDP-43 clearance:
491 evidence for dysfunction in amyotrophic lateral sclerosis. *Brain* **139**, 1417-32 (2016).
- 492 42. Udan-Johns, M. et al. Prion-like nuclear aggregation of TDP-43 during heat shock is
493 regulated by HSP40/70 chaperones. *Hum Mol Genet* **23**, 157-70 (2014).
- 494 43. Zhang, Y.J. et al. Phosphorylation regulates proteasomal-mediated degradation and
495 solubility of TAR DNA binding protein-43 C-terminal fragments. *Mol Neurodegener* **5**,
496 33 (2010).
- 497 44. Benatar, M. et al. Randomized, double-blind, placebo-controlled trial of arimoclomol in
498 rapidly progressive SOD1 ALS. *Neurology* **90**, e565-e574 (2018).
- 499 45. Hageman, J. et al. A DNAJB chaperone subfamily with HDAC-dependent activities
500 suppresses toxic protein aggregation. *Mol Cell* **37**, 355-69 (2010).
- 501 46. Outeiro, T.F. et al. Formation of toxic oligomeric alpha-synuclein species in living cells.
502 *PLoS One* **3**, e1867 (2008).
- 503 47. Crippa, V. et al. The small heat shock protein B8 (HspB8) promotes autophagic removal
504 of misfolded proteins involved in amyotrophic lateral sclerosis (ALS). *Hum Mol Genet*
505 **19**, 3440-56 (2010).
- 506 48. Sharp, P.S. et al. Protective effects of heat shock protein 27 in a model of ALS occur in
507 the early stages of disease progression. *Neurobiol Dis* **30**, 42-55 (2008).
- 508 49. Novoselov, S.S. et al. Molecular chaperone mediated late-stage neuroprotection in the
509 SOD1(G93A) mouse model of amyotrophic lateral sclerosis. *PLoS One* **8**, e73944 (2013).
- 510
- 511
- 512

513 **FIGURES**

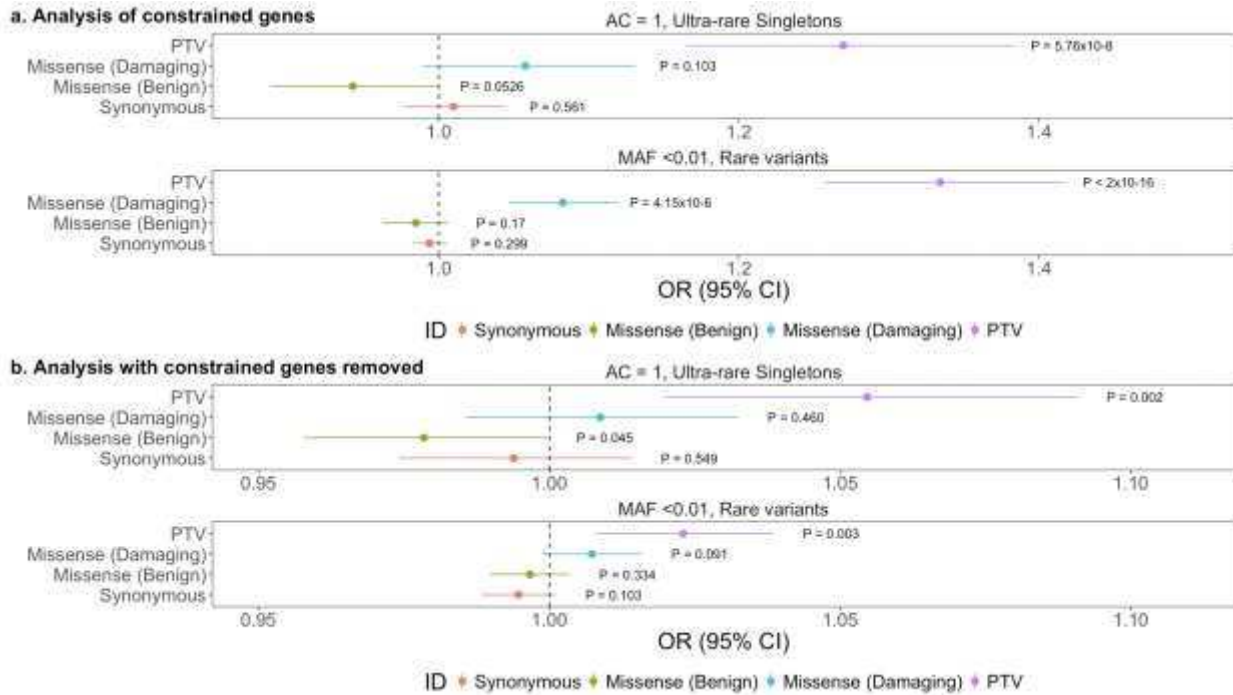


514

515 **Fig. 1. Exome wide enrichment of protein-truncating variants in ALS cases**

516 Exome wide analysis of synonymous variants, benign missense variants, damaging missense
 517 variants, and protein-truncating variants within singletons, doubletons, ultra-rare singletons, and
 518 rare variants. Odds ratios and 95% confidence intervals for each class of variation are depicted
 519 by different colors. P-values from firth logistic regression test are also displayed. Multiple test
 520 correction P-value: 0.0125. N=3,864 ALS cases; N=7,839 controls. The graph display the mean
 521 and standard deviation.

522



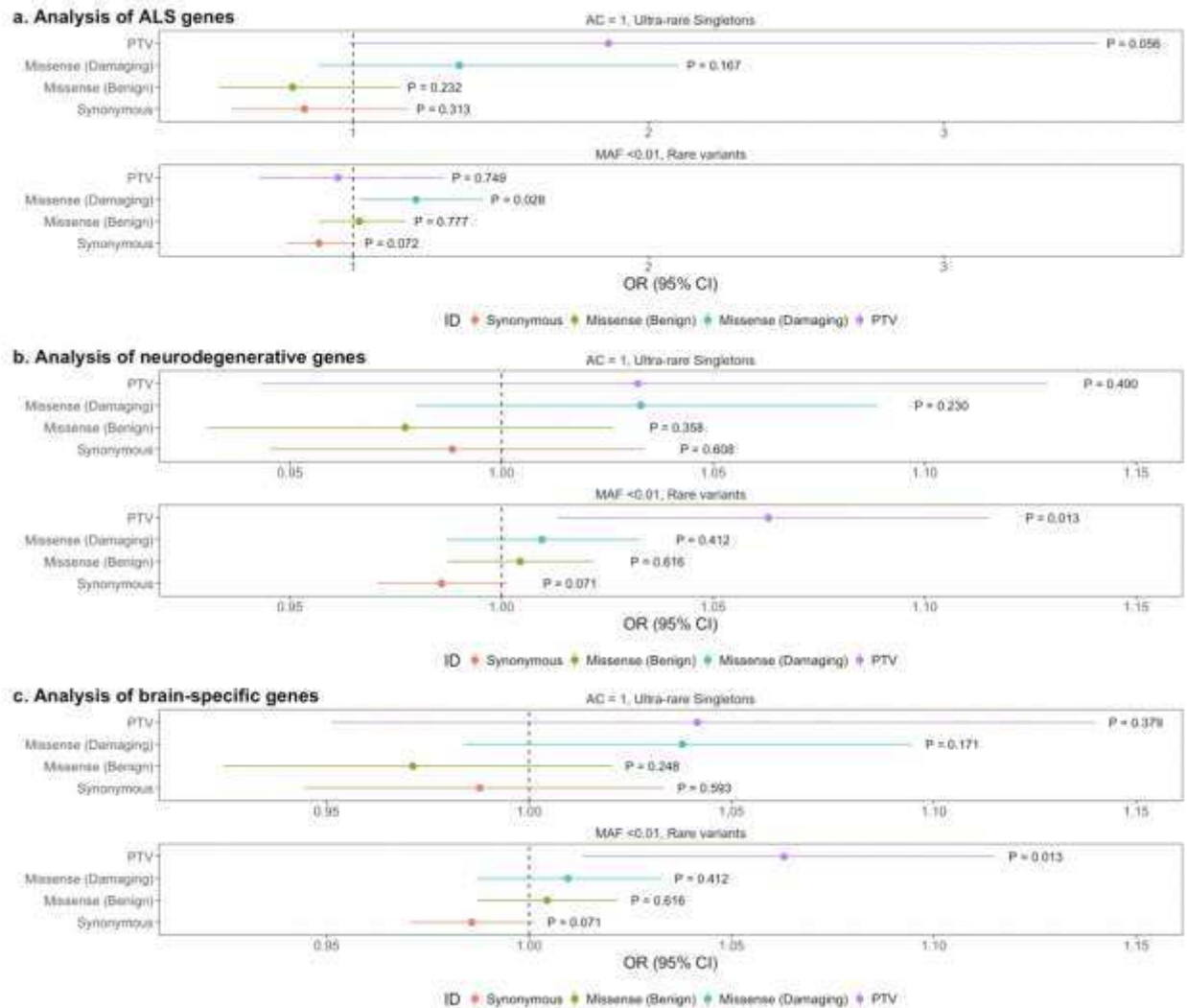
523

524 **Fig. 2. Enrichment of protein-truncating variants in constrained genes in ALS cases**

525 **a**, Analysis of constrained genes only in synonymous variants, benign missense variants,
 526 damaging missense variants, and protein-truncating variants within ultra-rare singletons and rare
 527 variants. Odds ratios and 95% confidence intervals for each class of variation are depicted by
 528 different colors. P-values from firth logistic regression test are also displayed. Multiple test
 529 correction P-value: 0.0125. N=3,864 ALS cases; N=7,839 controls. The graphs display the mean
 530 and standard deviation.

531 **b**, Exome-wide analysis with constrained genes removed.

532

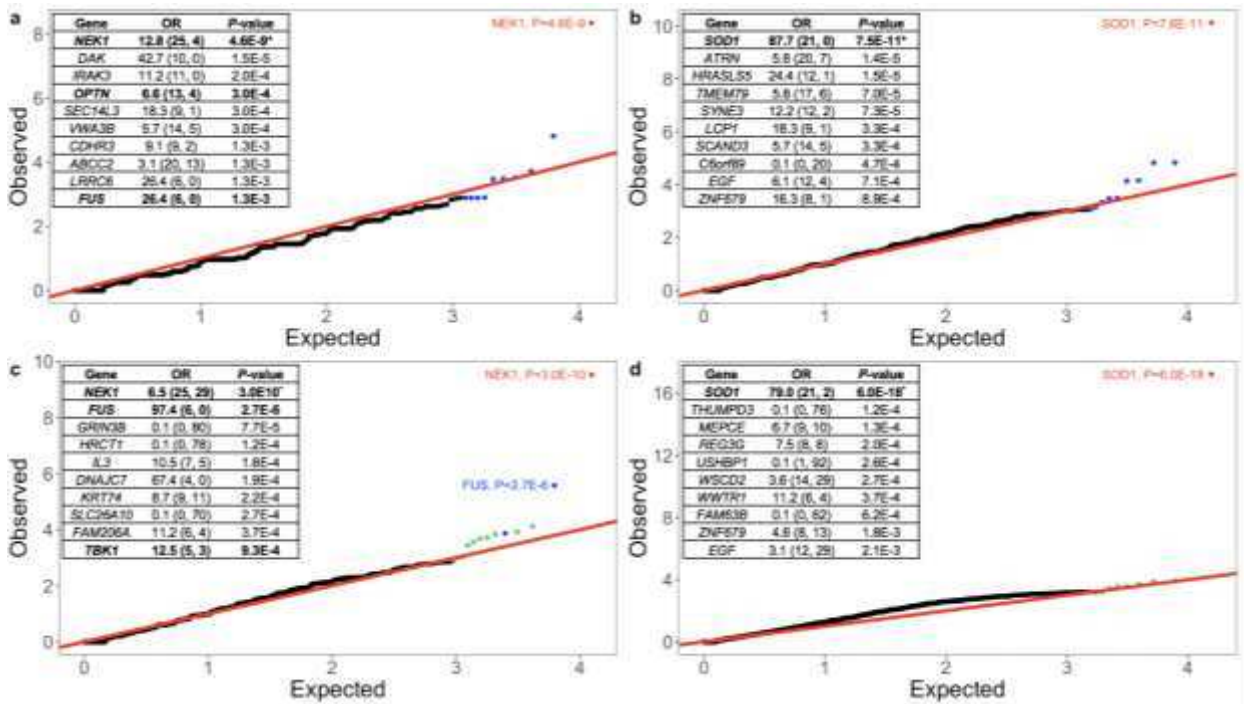


533
 534 **Fig. 3. No enrichment of variants in known ALS genes, other related neurodegenerative**
 535 **disease genes, or brain specific genes**
 536 **a, Analysis of ALS genes.** Synonymous variants, benign missense variants, damaging missense
 537 variants, and protein-truncating variants within singletons, doubletons, ultra-rare singletons, and
 538 rare variants are shown. Odds ratios and 95% confidence intervals for each class of variation are
 539 depicted by different colors. P-values from firch logistic regression test are also displayed.
 540 Multiple test correction P-value: 0.0125. N=3,864 ALS cases; N=7,839 controls. The graphs
 541 display the mean and standard deviation.

542 **b**, Analysis of other neurodegenerative disease genes.

543 **c**, Analysis of brain specific genes.

544



545

546 **Fig. 4. Quantile-quantile plot of discovery results for rare variants**

547 **a**, Rare protein-truncating variants in ALS dataset. X and Y axis represent the negative logarithm

548 P-value. N=3,864 ALS cases; N=7,839 controls. The top 10 genes with their P-values are

549 displayed. Genes in red and blue pass or approach exome-wide significance, respectively. The

550 results displayed are from a burden analysis using Fisher's exact test as well as SKAT, with

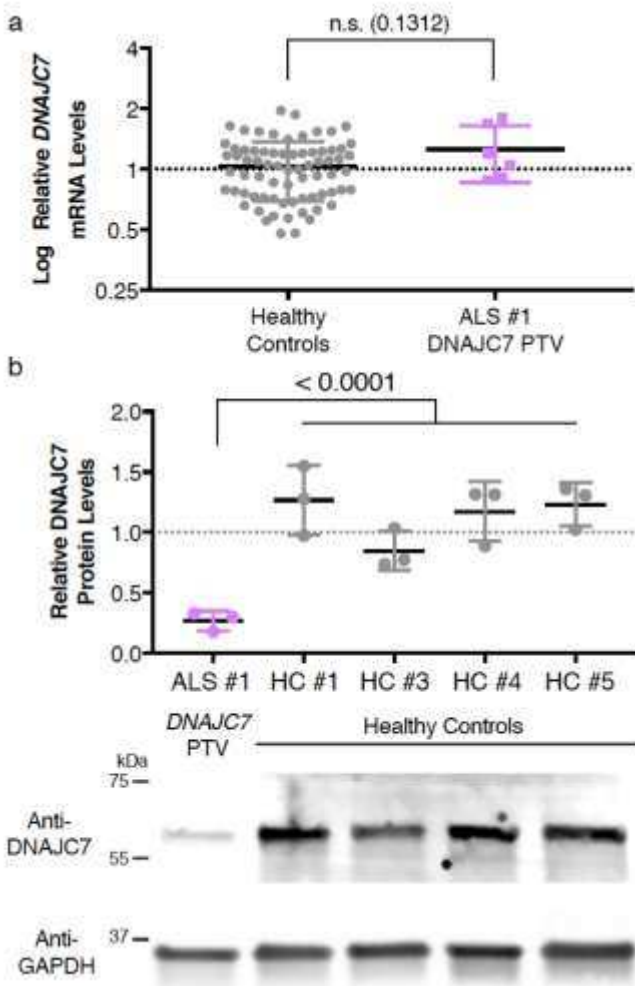
551 previously defined covariates (sample sex, PC1-PC10, and total exome count). Exome-wide

552 correction for multiple testing was set at ($P < 2.5 \times 10^{-6}$), which was the 5% type-I error rate

553 multiplied by the number of genes tested.

554 **b**, Rare damaging missense variants in ALS dataset.

555 **c**, Rare protein-truncating variants in ALS cases with an additional 21,071 non-Finnish European
 556 controls for a total of 28,910 controls. Genes in blue were the most significant genes in the
 557 discovery analysis. Genes in green were the most significant genes in the secondary analysis.
 558 **d**, Rare damaging missense variants in ALS cases and 28,910 controls. The top 10 genes with
 559 their P-values are displayed.
 560



561
 562 **Fig. 5. Effects of DNAJC7 protein-truncating variant p.Arg156Ter on transcript and**
 563 **protein levels.**

564 **a**, qRT-PCR analysis of DNAJC7 transcripts in human fibroblasts from healthy controls or a
 565 patient harboring a DNAJC7 protein-truncating variant p.Arg156Ter. Data were normalized to
 566 GAPDH and displayed as mean of 3 technical replicates with s.d. from two independent
 567 experiments with n=12 control and 1 patient lines (unpaired t test, two-sided, P<0.05). P-value is
 568 displayed, 0.1312.

569 **b**, Immunoblot analysis for DNAJC7 protein levels in human fibroblast lysates. Protein levels
 570 were normalized to GAPDH and displayed relative to the average levels in healthy controls. Data
 571 are displayed as mean with s.d. of n=3 technical replicates (unpaired t test, two-sided, P< 0.05).
 572 P-value is displayed, <0.0001. The blot image was cropped to make this figure, for the full scan
 573 of the blot, please see Supplementary Fig. 10.

574

575 TABLES

576 **Table 1. Top hits in protein-truncating variants model in initial (3,864 cases and 7,839**
 577 **controls), and secondary datasets (3,864 cases and 28,910 controls).**

578

Gene	Initial OR	Initial P-value	Secondary OR	Secondary P-value
Protein truncating variants model				
NEK1	12.8 (25, 4)	4.6×10⁻⁹*	6.5 (25, 29)	3.0×10⁻¹⁰#
DAK	42.7 (10, 0)	1.5×10 ⁻⁵	5.8 (10, 13)	1.4×10 ⁻⁴
IRAK3	11.2 (11, 0)	2.0×10 ⁻⁴	2.1 (11, 40)	0.05
OPTN	6.6 (13, 4)	3.0×10⁻⁴	2.6 (13, 38)	6.9×10⁻³
SEC14L3	18.3 (9, 1)	3.0×10 ⁻⁴	1.4 (9, 48)	0.31
VWA3B	5.7 (14, 5)	3.0×10 ⁻⁴	3.0 (14, 50)	0.02
CDHR3	9.1 (9, 2)	1.3×10 ⁻³	1.4 (9, 70)	0.9
ABCC2	3.1 (20, 13)	1.3×10 ⁻³	1.8 (20, 84)	0.03
LRRC6	26.4 (6, 0)	1.3×10 ⁻³	1.2 (6, 38)	0.64
FUS	26.4 (6, 0)	1.3×10⁻³	97.4 (6, 0)	2.7×10⁻⁶#
GRIN3B	0.4 (0, 10)	0.04	0.05 (0, 80)	7.7×10 ⁻⁵ #
HRCT1	0.2 (0, 15)	4.1×10 ⁻³	0.05 (0, 78)	1.2×10 ⁻⁴ #
IL3	14.2 (7, 1)	2.4×10 ⁻³	10.5 (7, 5)	1.8×10 ⁻⁴ #
DNAJC7	18.3 (4, 0)	0.01	67.4 (4, 0)	1.9×10 ⁻⁴ #
KRT74	3.0 (9, 6)	0.05	8.7 (9, 11)	2.2×10 ⁻⁴ #
SLC26A10	0.1 (0, 9)	0.03	0.1 (0, 70)	2.7×10 ⁻⁴ #
FAM206A	4.1 (6, 3)	0.07	11.2 (6, 4)	3.7×10 ⁻⁴ #
TBK1	22.3 (5, 0)	3.9×10⁻³	12.5 (5, 3)	9.3×10⁻⁴#
KLHDC4	0.1 (0, 14)	7.4 ×10 ⁻³	0.1 (0, 61)	9.8×10 ⁻⁴ #
DUOXA2	2.4 (7, 6)	0.14	5.8 (7, 9)	1.4×10 ⁻³ #

Damaging missense variants model				
SOD1	87.7 (21, 0)	7.5×10⁻¹¹*	79.0 (21, 2)	6.0×10⁻¹⁸#
ATRNL1	5.8 (20, 7)	1.4×10 ⁻⁵	2.0 (20, 74)	9.2×10 ⁻³
HRASLS5	24.4 (12, 1)	1.5×10 ⁻⁵	1.6 (12, 56)	0.13
TMEM79	5.8 (17, 6)	7.0×10 ⁻⁵	1.8 (17, 70)	0.03
SYNE3	12.2 (12, 2)	7.3×10 ⁻⁵	1.1 (12, 79)	0.62
LCP1	18.3 (9, 1)	3.3×10 ⁻⁴	3.7 (9, 18)	2.8×10 ⁻³
SCAND3	5.7 (14, 5)	3.3×10 ⁻⁴	1.8 (14, 58)	0.06
C6orf89	0.05 (0, 20)	4.7×10 ⁻⁴	0.05 (0, 47)	5.2×10 ⁻³
EGF	6.1 (12, 4)	7.1×10 ⁻⁴	3.1 (12, 29)	2.1×10 ⁻³
ZNF679	16.3 (8, 1)	8.9×10 ⁻⁴	4.6 (8, 13)	1.8×10 ⁻³
THUMPD3	0.07 (0, 15)	4.1×10 ⁻³	0.08 (0, 76)	1.2×10 ⁻⁴ #
MEPCE	9.1 (9, 2)	1.3×10 ⁻³	6.7 (9, 10)	1.3×10 ⁻⁴ #
REG3G	5.4 (8, 3)	8.2×10 ⁻³	7.5 (8, 8)	2.0×10 ⁻⁴ #
USHBP1	0.09 (1, 22)	1.6×10 ⁻³	0.08 (1, 92)	2.6×10 ⁻⁴ #
WSCD2	3.6 (14, 8)	4.8×10 ⁻³	3.6 (14, 29)	2.7×10 ⁻⁴ #
WWTR1	26.4 (6, 0)	1.3×10 ⁻³	11.2 (6, 4)	3.7×10 ⁻⁴ #
FAM63B	0.08 (0, 12)	0.01	0.06 (0, 62)	6.2×10 ⁻⁴ #

579 *Passed exome-wide significance (P-value <2.5×10⁻⁶) in first analysis (3,864 cases and 7,839
580 controls.#OR direction is maintained in secondary analysis (3,864 cases and 28,910 controls) and
581 P-value is lower. Bolded genes have been previously reported in ALS. The results displayed are
582 from a burden analysis using Fisher's exact test as well as SKAT, with previously defined
583 covariates (sample sex, PC1-PC10, and total exome count). Exome-wide correction for multiple
584 testing was set at (P<2.5×10⁻⁶), which was the 5% type-I error rate multiplied by the number of
585 genes tested.

586
587

588 **Table 2. Protein-truncating variants and 'damaging' missense variants in DNAJC7.**

589

Variant type	Variant location	cDNA change	Protein change	Cases (n=5,095)	Controls (n=28,910)	gnomAD (non-neuro) AF	CADD	MPC
Stop gain	17:g.40152569C>A	c.97G>T	p.E33X	1	0	0	39	
Stop gain	17:g.40148376G>A	c.358C>T	p.Q120X	1	0	0	37	
Stop gain	17:g.40146902G>A	c.466C>T	p.R156X	2	0	0	41	
Frameshift	17:g.40142393delA	c.488delT	p.F163fs	1	0	0		
Stop gain	17:g.40141529G>A	c.646C>T	p.R216X	2	0	0	40	
Essential splice site	17:g.40135656T>C	c.1011-2A>G		1	0	0	26.3	
Missense	17:g.40169413C>G	c.22G>C	p.D8H	1	0	1.985×10 ⁻⁵	25	0.78
Missense	17:g.40149189G>A	c.235C>T	p.R79W	0	1	1.204×10 ⁻⁵	35	1.58
Missense	17:g.40141544C>T	c.631G>A	p.D211N	1	0	0	26.4	0.94
Missense	17:g.40134023G>A	c.1234C>T	p.R412W	1	0	4.029×10 ⁻⁶	34	1.66
Missense	17:g.40133984C>T	c.1273G>A	p.E425K	2	0	0	35	1.69

590 AF, allele frequency; empty cell denotes inapplicable information.

591

592

593 **METHODS**

594 **Study overview**

595 The familial ALS (FALS) and the ALS Genetics (ALSGENS) consortia were assembled to
596 aggregate the existing ALS sequencing data in the community to improve the power to discover
597 novel genetic risk factors for ALS. Herein, we describe our approach of assembling the largest
598 ALS exome case-control study to date.

599

600 **Sample acquisition**

601 Blood samples were collected from subjects following appropriate and informed consent in
602 accordance with the Research Ethics Board at each respective recruiting site within the CReATe,
603 FALS, and ALSGENS consortia. All samples known to be carriers of the C9orf72
604 hexanucleotide expansion (G4C2) were excluded from the study. Additionally, prior to exome
605 sequencing, a subset of the samples (approximately 2,000) were genotyped and screened for
606 known variants in known ALS genes, SOD1, FUS, and TARDBP; and were only included in our
607 study if they were found to be negative for the variants tested.

608 Exome sequencing data for control and a subset of case samples were downloaded from
609 dbGAP and were not enriched for (but not specifically screened for) ALS or other
610 neurodegenerative disorders. Control samples were matched to case samples with respect to
611 similar capture kits and coverage levels. The age of control samples was not provided for all
612 samples but in general, controls were older than typical age of onset of ALS. The data are
613 available under the following accession codes: MIGen Exome Sequencing: Ottawa Heart
614 (phs000806.v1.p1); MIGen Exome Sequencing: Leicester UK Heart Study (phs001000.v1.p1);
615 Swedish Schizophrenia Population-Based Case-control Exome Sequencing (phs000473.v2.p2);
616 Genome-Wide Association Study of Amyotrophic Lateral Sclerosis (phs000101.v5.p1).

617 No statistical methods were used to pre-determine sample sizes but our sample sizes are
618 similar to those reported in previous publications⁹. Randomization of experimental groups was
619 not applicable to this study. The experimental conditions are determined by each individual's
620 genetics, which are fixed at conception. This reflects a randomization of the alleles inherited
621 from each individual's parents (i.e. mendelian randomization), but it does not involve
622 randomization of experimental parameters. Blinding was not relevant to the study as this study
623 was composed of cases and controls. Therefore, the analyst needed to know the case-control
624 status of every participant.

625

626 **Whole exome sequencing**

627 15,722 DNA samples were sequenced at the Broad Institute, Guy's Hospital, McGill University,
628 Stanford University, HudsonAlpha, and University of Massachusetts, Worcester. Samples were
629 sequenced using the exome Agilent All Exon (37MB, 50MB, or 65MB), Nimblegen SeqCap EZ
630 V2.0 or 3.0 Exome Enrichment kit, Illumina GAIIX, HiSeq 2000, or HiSeq 2500 sequencers
631 according to standard protocols.

632 All samples were joint called together and were aligned to the consensus human genome
633 sequence build GRCh37/hg19; and BAM files were processed using BWA Picard. Genotype
634 calling was performed using the Genome Analysis Toolkit's (GATK) HaplotypeCaller and was
635 performed at the Broad Institute as previously described^{50,51}.

636

637 **Hail software and quality control**

638 Code availability: we used Hail, an open-source, scalable framework for exploring and analyzing
639 genomic data <https://hail.is/> to process the data. All quality control steps were performed using
640 Hail 0.1 (Supplementary Table 1).

641

642 (1) Sample QC and Variant QC

643 Samples with high proportion of chimeric reads (>5%) and high contamination (>5%) were
644 excluded. Samples with poor call rates (<90%), mean depth <10x, or mean genotype-quality <65
645 were also eliminated from further analysis.

646 For variant QC, we restricted to GENCODE coding regions, independent of capture
647 interval, where both Agilent and Illumina exomes surpass 10x mean coverage. We restricted to
648 ‘PASS’ variants in GATK’s Variant Quality Score Recalibration (VQSR) filter. Individual
649 genotypes were filtered (set to missing) if they did not meet the following criteria: 1) genotype
650 depth (g.DP) 10 or greater 2) Allele balance ≥ 0.2 in heterozygous sites or ≤ 0.8 for
651 homozygous reference and homozygous alternate variants 3) Genotype quality (GQ) > 20 .
652 Finally, we selected variants with call rate >90% and Hardy-Weinberg equilibrium test P-value
653 $> 1 \times 10^{-6}$. For quality control analysis, see Supplementary Table 1 and Supplementary Fig. 1.

654

655 (2) Sex imputation

656 We used the X chromosome inbreeding coefficient to impute sample sex. Samples with an X
657 chromosome inbreeding coefficient >0.8 were classified as males and samples with an X
658 chromosome inbreeding coefficient <0.4 were classified as females. Samples within <0.8 and
659 >0.4 were classified as having ambiguous sex status, and therefore were excluded from the
660 dataset (Supplementary Table 1).

661

662 (3) Principal component analysis

663 Principal component analysis (PCA) was performed using Hail. We used a subset of high
664 confidence SNPs in the exome capture region to calculate the principal components. We used
665 only ancestry-matched cases and controls as indicated by overlapping population structure.
666 Furthermore, we used 1000 Genomes samples to determine the general ethnicity of the ALS
667 dataset. The majority of the samples in the ALS dataset were reported to be of European descent
668 and this was confirmed by PCA with 1000 Genomes samples (Supplementary Fig. 2,
669 Supplementary Table 1).

670

671 (4) Relatedness check

672 We included only unrelated individuals (IBD proportion < 0.2) (Supplementary Table 1).

673

674 (5) Variant annotation

675 We annotated protein-coding variants into four classes: (1) synonymous; (2) benign missense;
676 (3) damaging missense; and (4) protein-truncating variants (PTV). Using VEP annotations
677 (Version 85)⁵², we classified synonymous variants as: "synonymous_variant",
678 "stop_retained_variant", and "incomplete_terminal_codon_variant". Missense variants were
679 classified as: "inframe_deletion", "inframe_insertion", "missense_variant", "stop_lost",
680 "start_lost", and "protein_altering_variant". Furthermore, benign missense variants were
681 predicted as "tolerated" and "benign" by PolyPhen-2 and SIFT, respectively; whereas damaging
682 missense variants were predicted as "probably damaging" and "deleterious". Finally, protein-

683 truncating variants were classified as: "frameshift_variant", "splice_acceptor_variant",
684 "splice_donor_variant", and "stop_gained".

685

686 (6) Allele frequency categorization

687 Allele frequencies were estimated within our case-control sample, and from two external exome
688 sequence databases, DiscovEHR and ExAC⁵³. DiscovEHR is a publicly available database with
689 >50,000 exomes of participants who may have some health conditions however, they do not have
690 ALS. ExAC is a mixture of healthy controls and complex disease patients, and we restricted to
691 the non-psychiatric subset of ExAC for allele frequency estimation. Of note, many of our
692 controls are present in the ExAC database, so we restricted to the DiscovEHR cohort to
693 determine ultra-rare singletons. We did not use gnomAD for this analysis as our cases and our
694 controls have been deposited into this resource.

695 We classified variant allele frequency using the following criteria: (1) singletons, which
696 are variants present in a single individual in our dataset (allele count, AC = 1); (2) doubletons,
697 which are present in two individuals in our dataset (AC = 2); (3) ultra-rare singletons, which are
698 singletons in our dataset and are absent in DiscovEHR (AC = 1, 0 in DiscovEHR); and finally,
699 (4) rare variants, which have a MAF of <0.01% in our dataset (11,703 samples), in ExAC (non-
700 psychiatric studies, >45,000 samples) and in DiscovEHR (>50,000 samples).

701

702 **Multivariate models used for analysis**

703 To determine whether an enrichment of a specific class of variation was present in ALS cases
704 versus controls, we ran multiple Firth logistic regression models. The Firth penalization is used
705 in the likelihood model due to the low counts in many tests, and helps to minimize the type I

706 error rate when multiple covariates are included in the model⁵⁴. Model 1 predicted ALS case-
707 control status solely from variant count; Model 2 incorporated multiple covariates: (1) sample
708 sex, (2) sample population structure from the first 10 principal components; Model 3
709 incorporated all covariates used the second model along with (3) sample total exome count,
710 which is the exome-wide count of variants in the specific frequency class tested. Finally, Model
711 4 is similar to Model 3, but instead uses the “benign variant” count as a covariate, which is the
712 exome-wide count of synonymous variants and benign missense variants only, rather than total
713 exome count. Model 3, which we considered to be the most conservative model to represent the
714 dataset, was used as the preferred model for our analysis (Supplementary Fig. 3).

715

716 **Exome-wide burden**

717 The four Firth logistic regression models above were used to predict case-control status from
718 exome-wide counts of synonymous, missense, and protein-truncating variants. Given that
719 sequencing errors are more prevalent when calling insertions or deletions (indels)^{55,56}, we
720 divided variants within the protein-truncating variants category as either 1) SNV-based protein-
721 truncating variants or 2) indel-based protein-truncating variants, due to single nucleotide variants
722 (SNVs) or indels, respectively. This ensures that any enrichment observed in protein-truncating
723 variants is not solely from indel-based protein-truncating variants.

724

725 **Gene sets**

726 (1) Constrained genes (pLI genes: 3,488, constrained missense genes: 1,730)

727 We evaluated whether variation in loss of function intolerant (pLI) genes are associated with

728 ALS using the same approach as described in the exome-wide approach however, we extracted

729 only high pLI genes from the exome. We obtained the genic pLI intolerance metrics from Lek et
730 al., 2016 available online:
731 (ftp://ftp.broadinstitute.org/pub/ExAC_release/release0.3/functional_gene_constraint/). For
732 protein-truncating variants, we used genes with a probability of loss-of-function intolerant (pLI)
733 >0.9. We also evaluated missense constrained genes generated by Samocha et al., 2014⁵⁷. For
734 missense variants, we used genes with a z-score of >3.09.

735

736 (2) ALS associated genes (38 genes)

737 We also examined exome-wide burden with known ALS genes removed. The list of ALS genes
738 are as follows: TARDBP, DCTN1, ALS2, CHMP2B, ARHGEF28, MATR3, SQSTM1, FIG4,
739 HNRNPA2B1, C9orf72, SIGMAR1, VCP, SETX, OPTN, PRPH, HNRNPA1, DAO, ATXN2, ANG,
740 FUS, PFN1, CENPV, TAF15, GRN, MAPT, PNPLA6, UNC13A, VAPB, SOD1, NEFH, ARPP21,
741 and UBQLN2. We did not remove TBK1, NEK1, KIF5A, C21orf2, MOBP, or SCFD1 as these
742 genes were discovered using datasets that contained a large subset of the same samples. We also
743 performed an analysis with all proposed ALS genes.

744

745 (3) Neurodegenerative disease genes (120 genes)

746 We investigated whether genes associated with other neurodegenerative phenotypes showed
747 enrichment in ALS cases. We included the following motor neuron diseases: primary lateral
748 sclerosis, progressive muscular atrophy, progressive bulbar palsy, and spinal muscular atrophy.
749 We also used genes associated with Parkinson's disease, frontotemporal dementia, Pick's
750 disease, and Alzheimer's disease as patients with ALS can also present with frontotemporal
751 dementia, cognitive impairment, or Parkinsonism (Supplementary Table 5).

752

753 (4) Brain expressed genes (2,650 genes)

754 We evaluated whether genes expressed specifically in the brain were enriched for variation in
755 our dataset. For this analysis, we used brain specific genes generated by Ganna et al., 2016.

756

757 **Single gene burden analysis**

758 (1) ALS dataset (3,864 cases and 7,839 controls)

759 To determine whether a single gene is enriched or depleted for rare protein-coding variation in
760 ALS cases, we performed a burden analysis using Fisher's exact test as well as SKAT, with
761 previously defined covariates (sample sex, PC1-PC10, and total exome count). Exome-wide
762 correction for multiple testing was set at ($P < 2.5 \times 10^{-6}$), which was the 5% type-I error rate
763 multiplied by the number of genes tested. We performed four different tests in ALS cases and
764 controls: (1) ultra-rare protein-truncating variants (AC=1 and absent in DiscovEHR); (2) ultra-
765 rare damaging missense variants (AC=1 and absent in DiscovEHR); (3) rare protein-truncating
766 variants (MAF < 0.001% in the dataset, DiscovEHR, and ExAC); and (4) rare damaging missense
767 variants (MAF < 0.001% in the dataset, DiscovEHR, and ExAC).

768

769 (2) ALS dataset and additional controls (3,864 cases and 28,910 controls)

770 We also included an additional 21,071 samples from ExAC that are of European descent (non-
771 Finnish) and were not a part of any psychiatric or brain related studies, to eliminate any sample
772 overlap. Furthermore, to mitigate against false discoveries, in addition to passing our QC filters,
773 we ensured each variant also passed gnomAD (123,136 exomes and 15,496 genomes) QC filters.
774 We included variants that were either a singleton (AC=1) in gnomAD or completely absent to

775 ensure we minimize the inclusion of an excess of variants that passed gnomAD QC, that were
776 rare (MAF <0.001%), yet were still observed in a very high number of individuals and were
777 likely, false positive variants. The additional 21,071 samples allowed us to perform a secondary
778 analysis of the genes that approached statistical significance ($P < 2.5 \times 10^{-6}$) and determine
779 whether their OR and P-values are maintained and exceed statistical significance, respectively.
780 Additionally, we also used the 21,071 controls to increase statistical power to detect any gene
781 discoveries not detected in the original dataset. Importantly, we did not perform a joint PCA on
782 the 21,071 non-Finnish European controls and our dataset, therefore, we are unable to
783 completely match the ancestry of our dataset.

784

785 **Cell acquisition culture and authentication**

786 The fibroblasts used in this study were previously approved by the institutional review boards
787 (IRBs) of Harvard University, Massachusetts General Hospital, and Columbia University.
788 Specific point mutations were confirmed by PCR amplification followed by Sanger sequencing.
789 Weekly, cultures were checked for mycoplasma contamination using the MycoAlert kit (Lonza)
790 with no cell lines used in this study testing positive. The use of these cells at Harvard was further
791 approved and determined not to constitute Human Subjects Research by the Committee on the
792 Use of Human Subjects in Research at Harvard University. Human fibroblasts were grown with
793 DMEM (Invitrogen) supplemented with 15% fetal bovine serum (VWR), 10 mM MEM Non-
794 essential amino acid (Millipore), and B-mercaptoethanol 55 μ M (Invitrogen), and cultured on
795 tissue culture dishes maintained in 5% CO₂ incubators at 37°C. Fibroblasts were passaged after
796 reaching confluency using trypsin (Invitrogen).

797

798 **Immunoblot assays**

799 For analysis of DNAJC7 protein expression levels, fibroblasts were lysed in RIPA buffer
800 (150mM Sodium Chloride; 1% Triton X-100; 0.5% sodium deoxycholate; 0.1% SDS; 50 mM
801 Tris pH 8.0) containing protease and phosphatase inhibitors (Roche) for 20 min on ice, and
802 centrifuged at high speed to remove insoluble components. 500 μ L of RIPA buffer per well of a
803 6-well plate were routinely used, which yielded \sim 20 μ g of total protein as determined by BCA
804 (Thermo Scientific). For immunoblot assays, 1 μ g of total protein was separated by SDS-PAGE
805 (BioRad), transferred to PDVF membranes (BioRad) and probed with antibodies against
806 DNAJC7 (1:1000, Abcam, Clone EPR13349) and GAPDH (1:1000, Millipore, Clone 6C5). LI-
807 COR software (Image Studios) was used to quantitate protein band signal, and GAPDH levels
808 were used to normalized each sample. Data are from three technical replicates with n=12 control
809 and 1 patient lines. To analyze the results from this experiment, we used an unpaired t test, two-
810 sided with a statistical threshold of $P < 0.05$.

811

812 **RNA preparation and qRT-PCR**

813 Total RNA was isolated from fibroblasts using Trizol (Invitrogen) according to manufacturer's
814 instructions. 500 μ L of Trizol were added per well of the 6-well cultures. A total of 300-1000ng
815 of total RNA was then used to synthesize cDNA by reverse transcription according to the
816 iSCRIPT kit (Bio-rad). Quantitative RT-PCR (qRT-PCR) was then performed using SYBR green
817 (Bio-Rad) and the iCycler system (Bio-rad). Quantitative levels for all genes assayed were
818 normalized using GAPDH expression. For comparison between control and patient lines,
819 normalized expression was displayed relative to the average of pooled data points from the
820 healthy controls. The primer sequences (forward, reverse) are for GAPDH

821 (AATGGTGAAGGTCGGTGTG, GTGGAGTCATACTGGAACATGTAG), DNAJC7 Exons 4-
822 6 (CAGTGAGGTTGGATGACAGTT, ACTCTTGTGTGCCTGAGC), DNAJC7 Exons 13-14
823 (TACTATCCTCTCTGATCCCAAGA, CCTTGTCTCCAGCTGAGAG). Data are from three
824 technical replicates with n=12 control and 1 patient lines. To analyze the results from this
825 experiment, we used an unpaired t test, two-sided with a statistical threshold of $P < 0.05$.

826

827 **Data presentation and statistical analysis**

828 In the figure elements, points and lines represent the median and standard deviation, respectively.
829 The plots display the minimum to maximum. Data distribution was assumed to be normal but
830 this was not formally tested. For the exome-wide and gene specific test, we build four models
831 that use firth logistic regression, please refer to ‘Multivariate models used for analysis’ in the
832 Materials and Methods section. Multiple test correction P -value < 0.0125 was considered
833 significant. For gene specific analyses, a multiple test correction P -value $< 2.5 \times 10^{-6}$ was
834 considered significant. For the immunoblot and qPCR assays, the statistical analyses were
835 performed using a two-tail unpaired Student’s t-test, with a P value of $*P < 0.05$ considered as
836 significant using Prism 7 (Graph Pad).

837

838 **Reporting Summary**

839 Further information on research design is available in the Nature Research Life Sciences
840 Reporting Summary linked to this article.

841

842 **Data availability**

843 The sequencing data discussed in this publication were obtained through dbGaP and are
844 available under the following accession codes: MIGen Exome Sequencing: Ottawa Heart
845 (phs000806.v1.p1); MIGen Exome Sequencing: Leicester UK Heart Study (phs001000.v1.p1);
846 Swedish Schizophrenia Population-Based Case-control Exome Sequencing (phs000473.v2.p2);
847 Genome-Wide Association Study of Amyotrophic Lateral Sclerosis (phs000101.v5.p1).

848

849 **Code availability**

850 Code used to conduct the analysis is provided online.

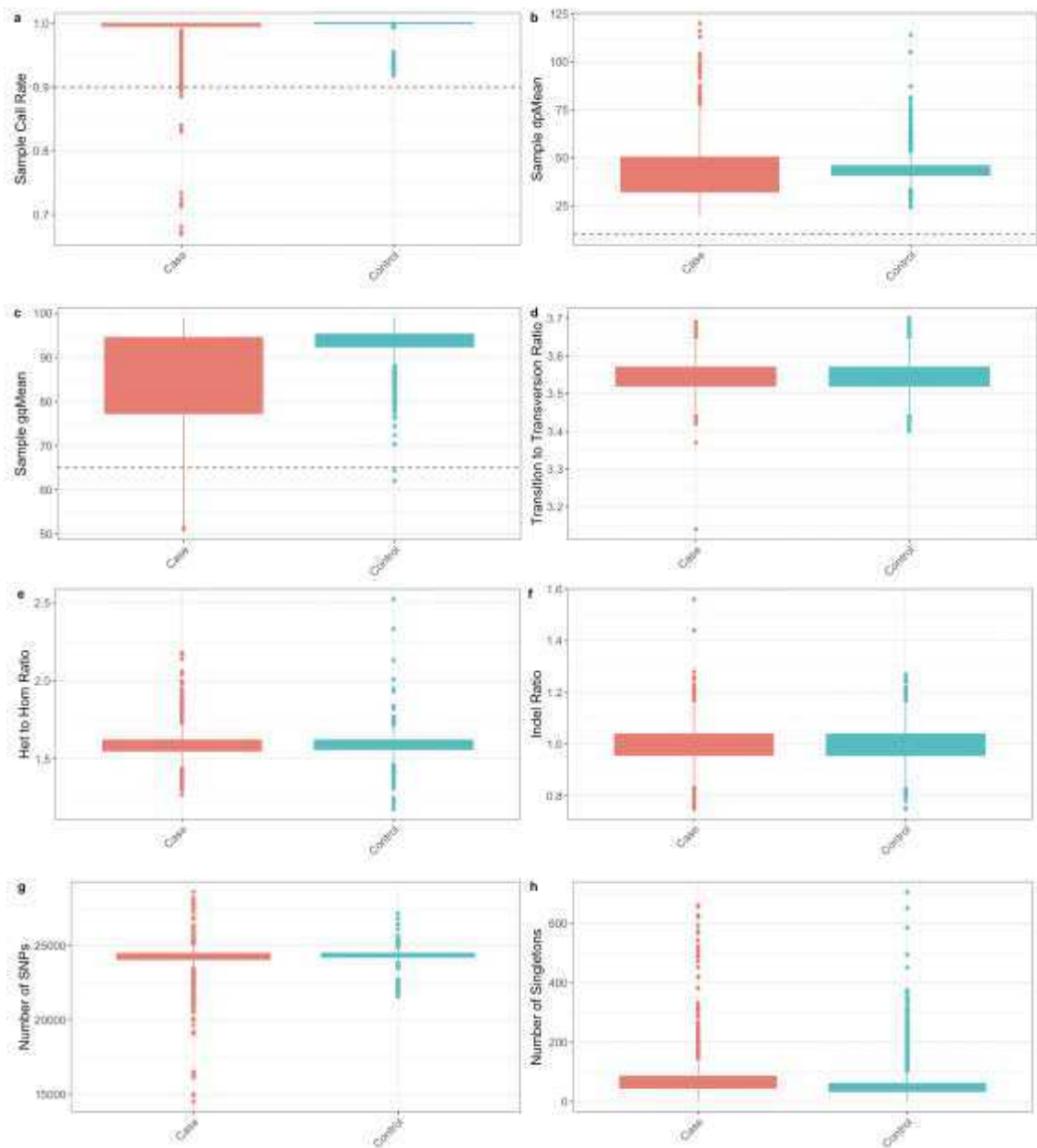
851

852 **METHODS-ONLY REFERENCES**

- 853 50. Ganna, A. et al. Ultra-rare disruptive and damaging mutations influence educational
854 attainment in the general population. *Nat Neurosci* **19**, 1563-1565 (2016).
- 855 51. Ganna, A. et al. Quantifying the Impact of Rare and Ultra-rare Coding Variation across
856 the Phenotypic Spectrum. *Am J Hum Genet* **102**, 1204-1211 (2018).
- 857 52. McLaren, W. et al. The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 122 (2016).
- 858 53. Lek, M. et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature*
859 **536**, 285-91 (2016).
- 860 54. Wang, X. Firth logistic regression for rare variant association tests. *Front Genet* **5**, 187
861 (2014).
- 862 55. Lam, H.Y. et al. Performance comparison of whole-genome sequencing platforms. *Nat*
863 *Biotechnol* **30**, 78-82 (2011).
- 864 56. O'Rawe, J. et al. Low concordance of multiple variant-calling pipelines: practical
865 implications for exome and genome sequencing. *Genome Med* **5**, 28 (2013).
- 866 57. Samocha, K.E. et al. A framework for the interpretation of de novo mutation in human
867 disease. *Nat Genet* **46**, 944-50 (2014).

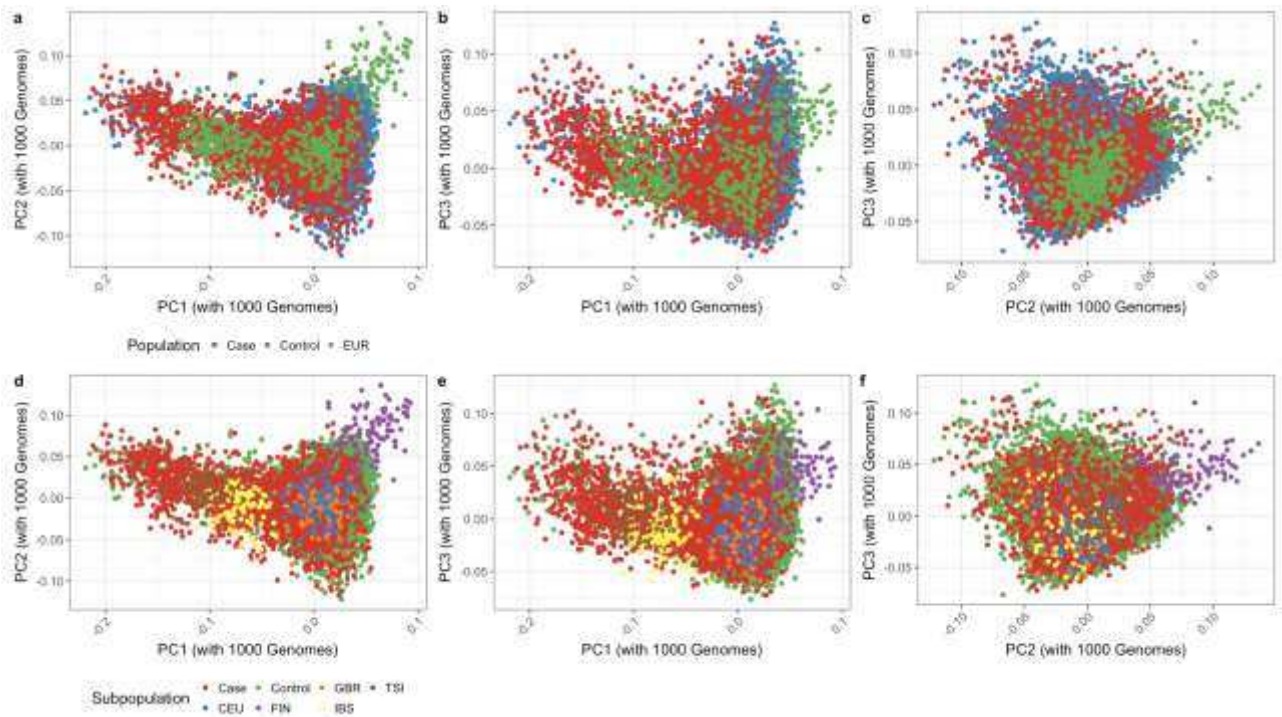
868

869



Supplementary Fig. 1. Initial sample quality control analysis

(A) Sample call rate. (B) Sample mean depth. (C) Sample mean genotype quality. (D) Sample transition to transversion ratio. (E) Sample heterozygous to homozygous ratio. (F) Sample insertion to deletion ratio. (G) Number of SNPs in each sample. (H) Number of singletons in each sample. N=3,864 ALS cases; N=7,839 controls. The box and whisker plots display the mean, minimum, and maximum.



Supplementary Fig. 2. Principal component analysis of ALS dataset with 1000 Genomes

(A) PC1 and PC2 of ALS dataset with 1000 Genomes. Cases, controls, and the European population is shown. N=3,864 ALS cases; N=7,839 controls. Each point represents one individual.

(B) PC1 and PC3 of ALS dataset with 1000 Genomes. Cases, controls, and the European population is shown.

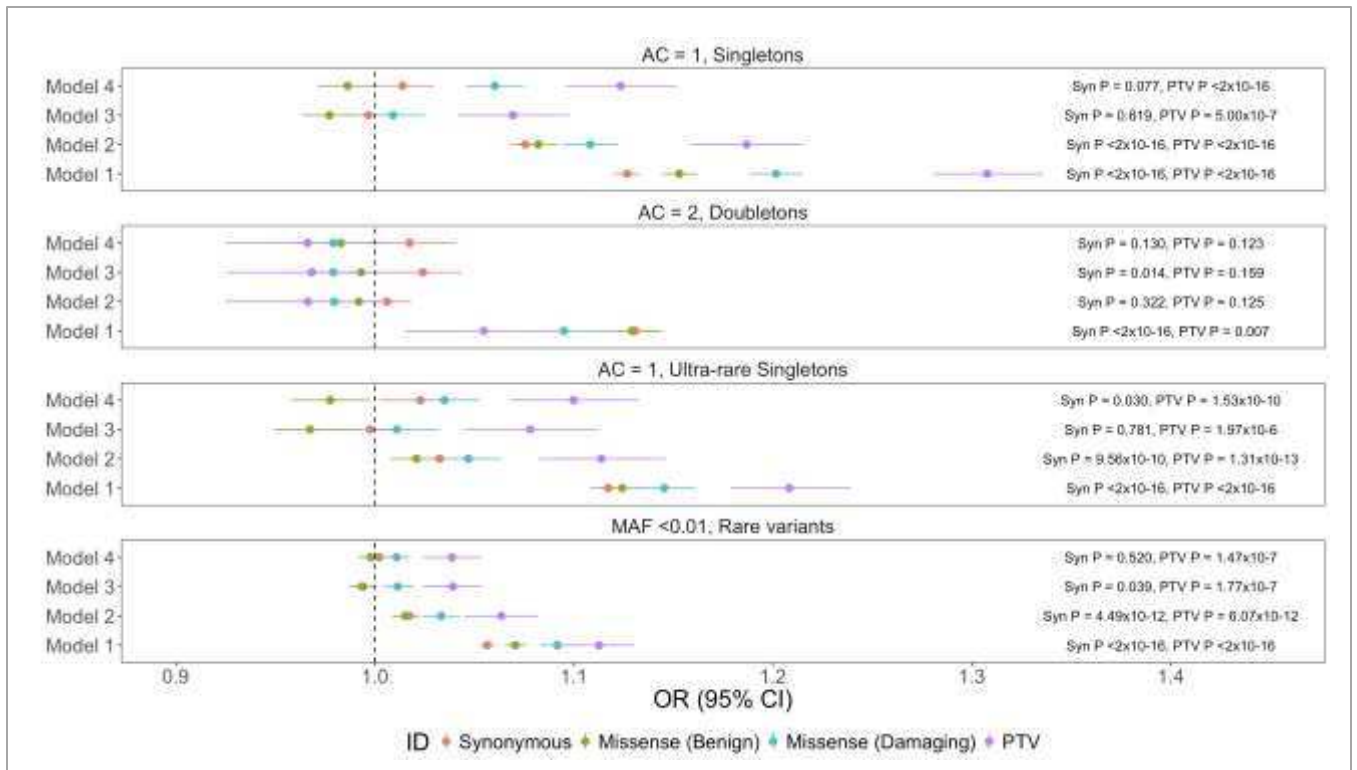
(C) PC2 and PC3 of ALS dataset with 1000 Genomes. Cases, controls, and the European population is shown.

(D) PC1 and PC2 of ALS dataset with 1000 Genomes. Cases, controls, and the European subpopulations are shown.

(E) PC1 and PC3 of ALS dataset with 1000 Genomes. Cases, controls, and the European subpopulations are shown.

(F) PC2 and PC3 of ALS dataset with 1000 Genomes. Cases, controls, and the European subpopulations are shown.

871
872



Supplementary Fig. 3. All models together

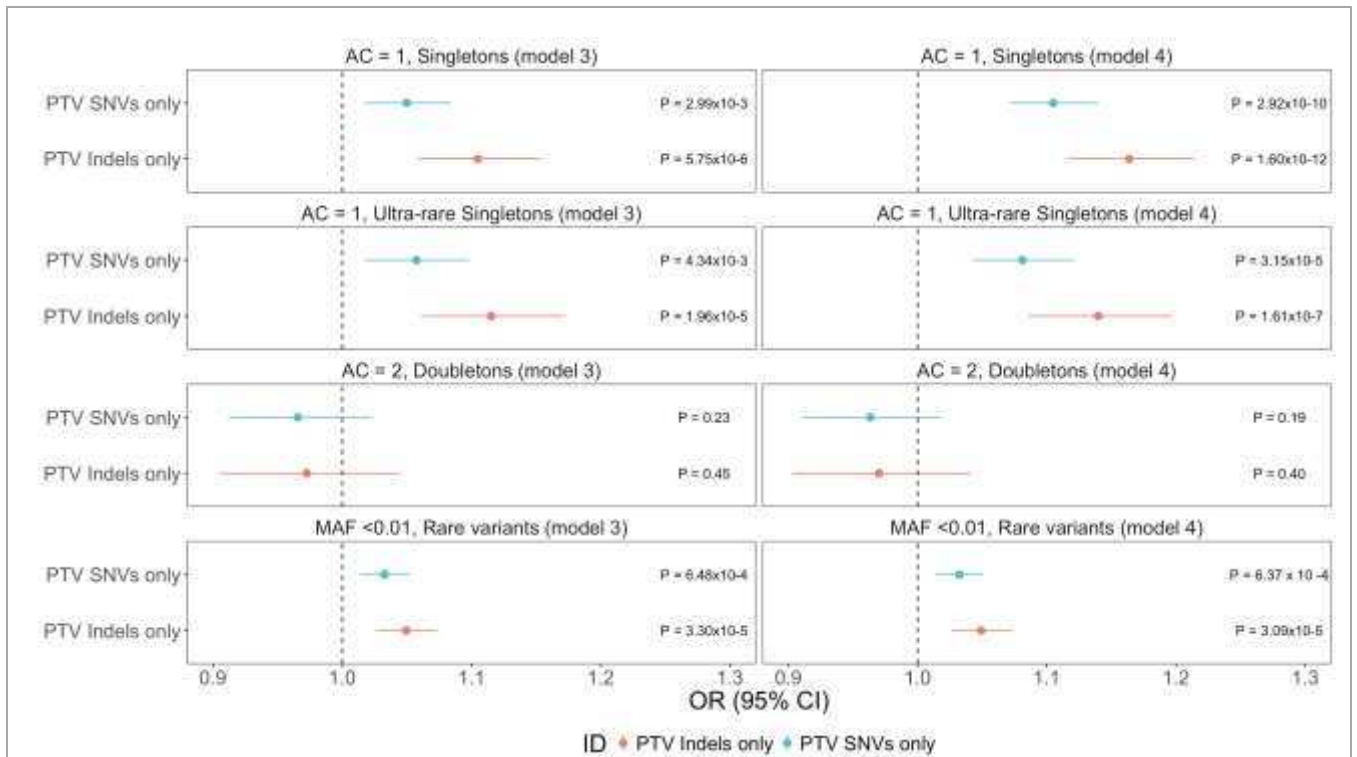
Model 1: Sample variation. The graph display the mean and standard deviation. P-values from firth logistic regression test are also displayed. Multiple test correction P-value: 0.0125. N=3,864 ALS cases; N=7,839 controls.

Model 2: Sample variation, sample sex, PC1-PC10.

Model 3: Sample variation, sample sex, PC1-PC10, and total exome count (summation of synonymous, benign missense, damaging missense, and PTV).

Model 4: Sample variation, sample sex, PC1-PC10, and benign variation count (summation of synonymous and benign missense variation).

873
874
875
876
877
878
879
880
881
882

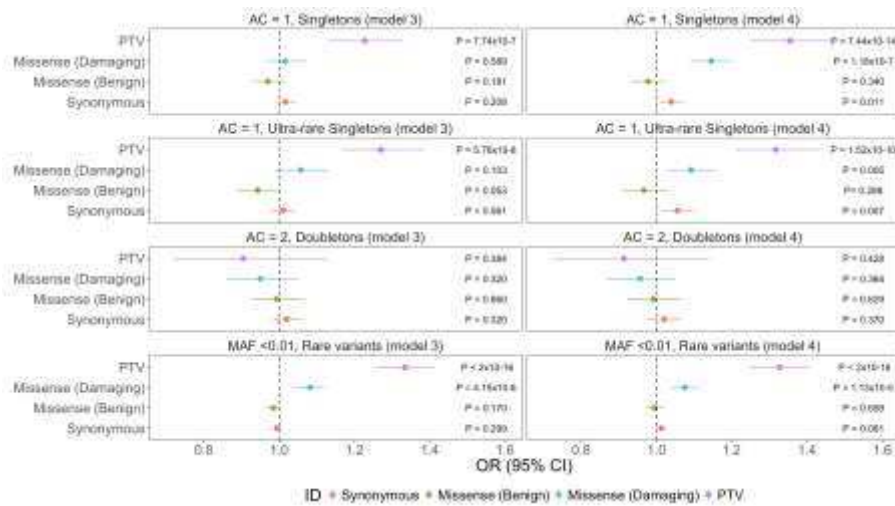


Supplementary Fig. 4. Exome wide enrichment of SNV-based PTVs and indel-based PTVs in ALS cases

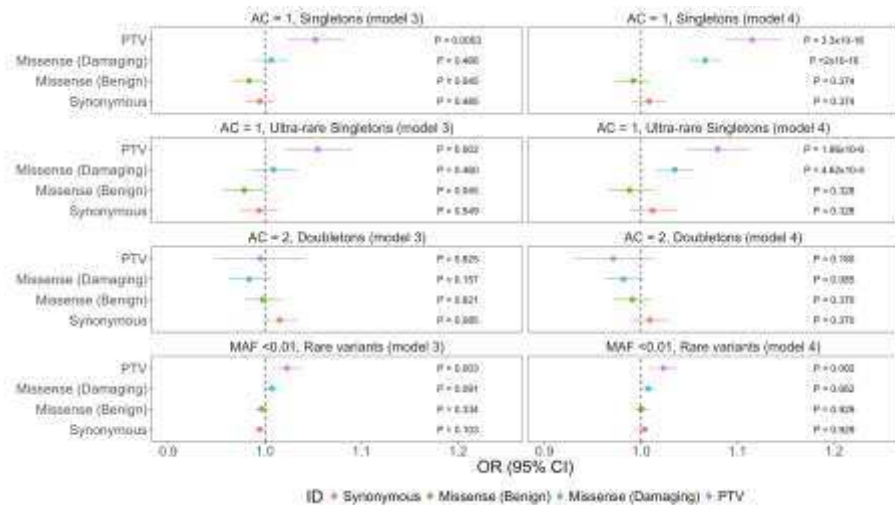
Extension of Fig. 1: Evaluating the effects of SNV-based and indel-based PTVs within singletons (AC=1), doubletons (AC=2), ultra-rare singletons (AC=1, 0 in DiscovEHR), and rare variants (MAF<0.01 in our dataset, DiscovEHR and ExAC). Odds ratios and 95% confidence intervals for each class of variation are depicted by different colors. P-values are also displayed. Model 3 evaluates sample variation with the covariates, sample sex, PC1-10, and total exome count (summation of synonymous variation, benign missense variation, damaging missense variation, and PTV SNV or PTV indel). Model 4 evaluates sample variation with the covariates, sample sex, PC1-10, and benign variation (summation of synonymous and benign missense variation). The graph display the mean and standard deviation. P-values from firth logistic regression test are also displayed. Multiple test correction P-value: 0.0125. N=3,864 ALS cases; N=7,839 controls.

883
884
885
886
887

a

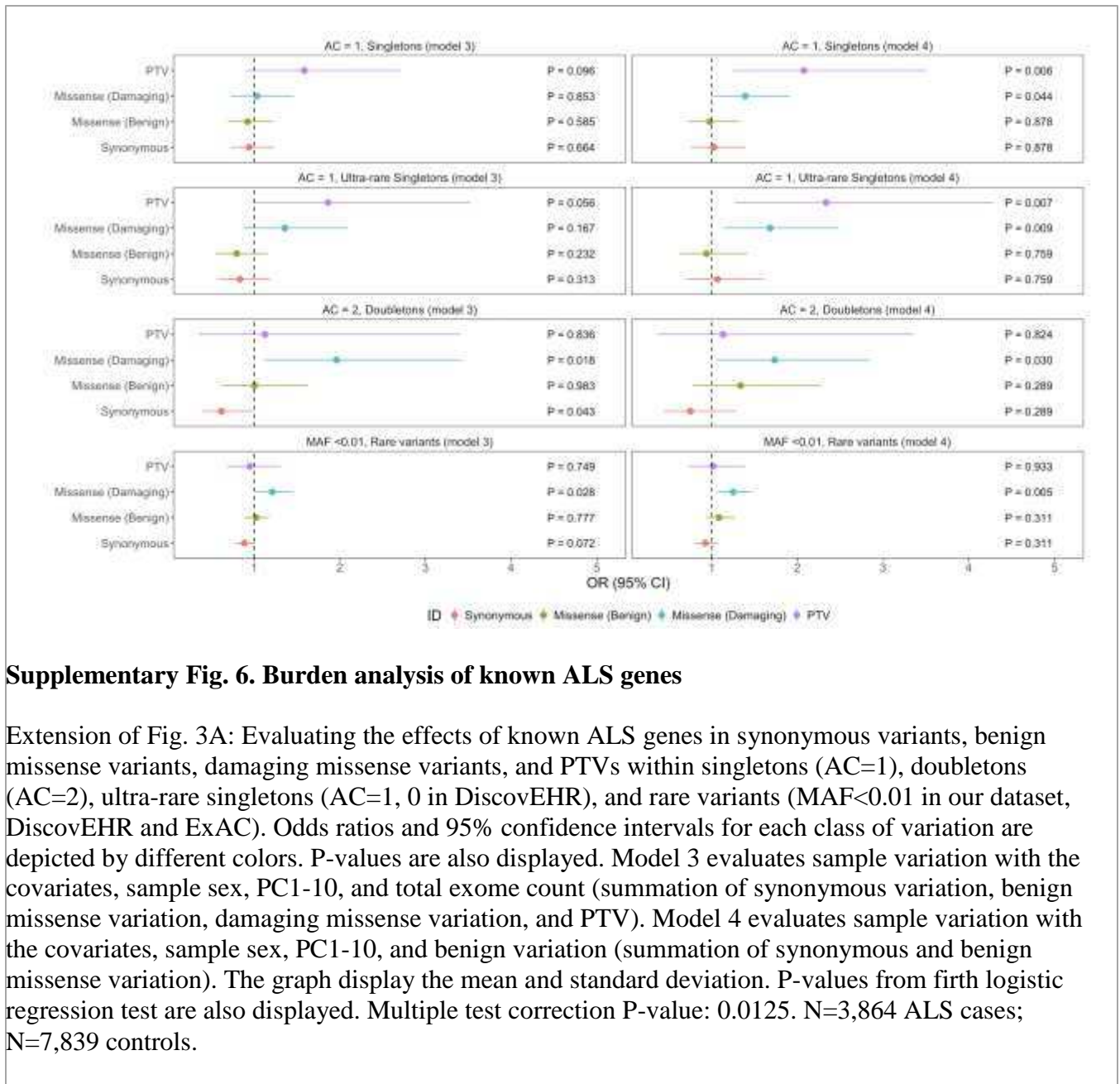


b

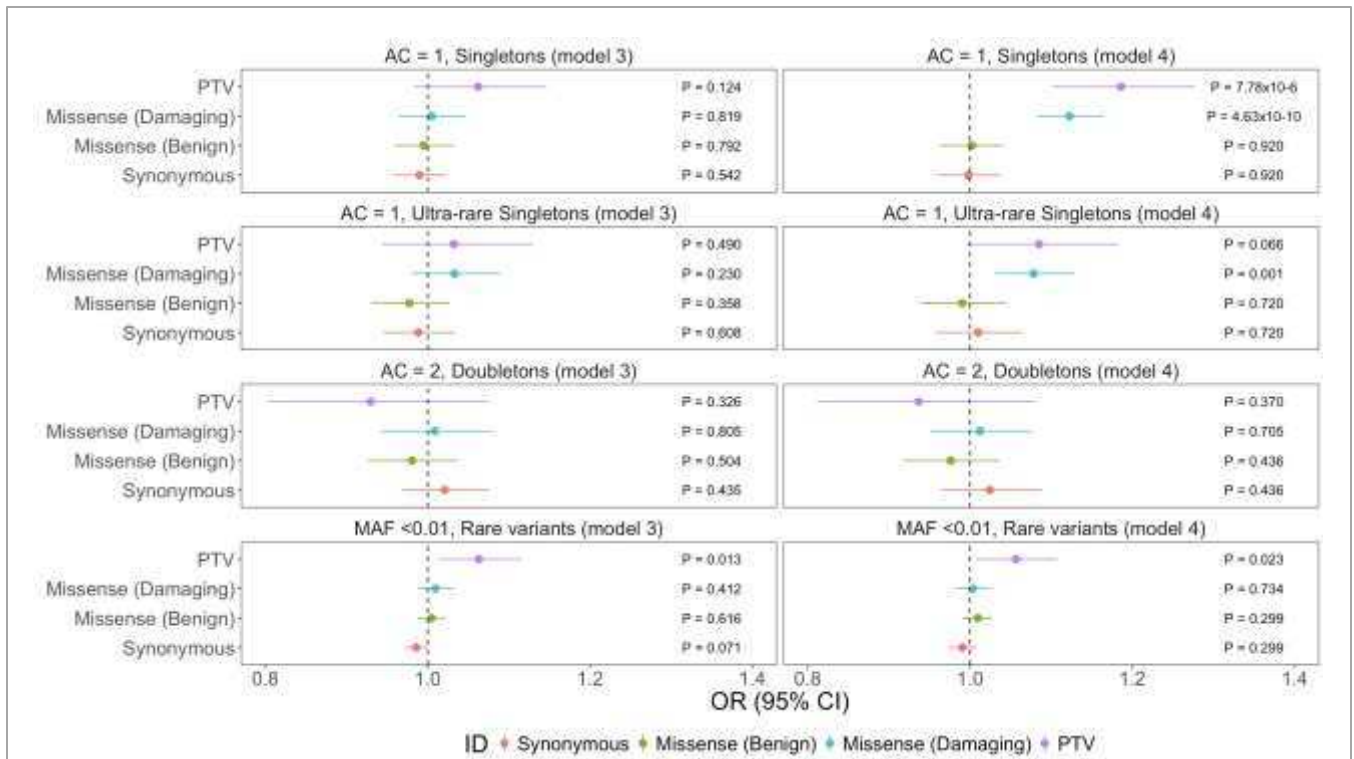


Supplementary Fig. 5. Enrichment of variants in constrained genes in ALS cases

Extension of Fig. 2A and 2B: (a) Evaluating the effects of constrained genes in synonymous variants, benign missense variants, damaging missense variants, and PTVs within singletons (AC=1), doubletons (AC=2), ultra-rare singletons (AC=1, 0 in DiscovEHR), and rare variants (MAF<0.01 in our dataset, DiscovEHR and ExAC). Odds ratios and 95% confidence intervals for each class of variation are depicted by different colors. P-values are also displayed. Model 3 evaluates sample variation with the covariates, sample sex, PC1-10, and total exome count (summation of synonymous variation, benign missense variation, damaging missense variation, and PTV). Model 4 evaluates sample variation with the covariates, sample sex, PC1-10, and benign variation (summation of synonymous and benign missense variation). (b) Evaluating the residual effects with constrained genes removed. The graph display the mean and standard deviation. P-values from firth logistic regression test are also displayed. Multiple test correction P-value: 0.0125. N=3,864 ALS cases; N=7,839 controls.



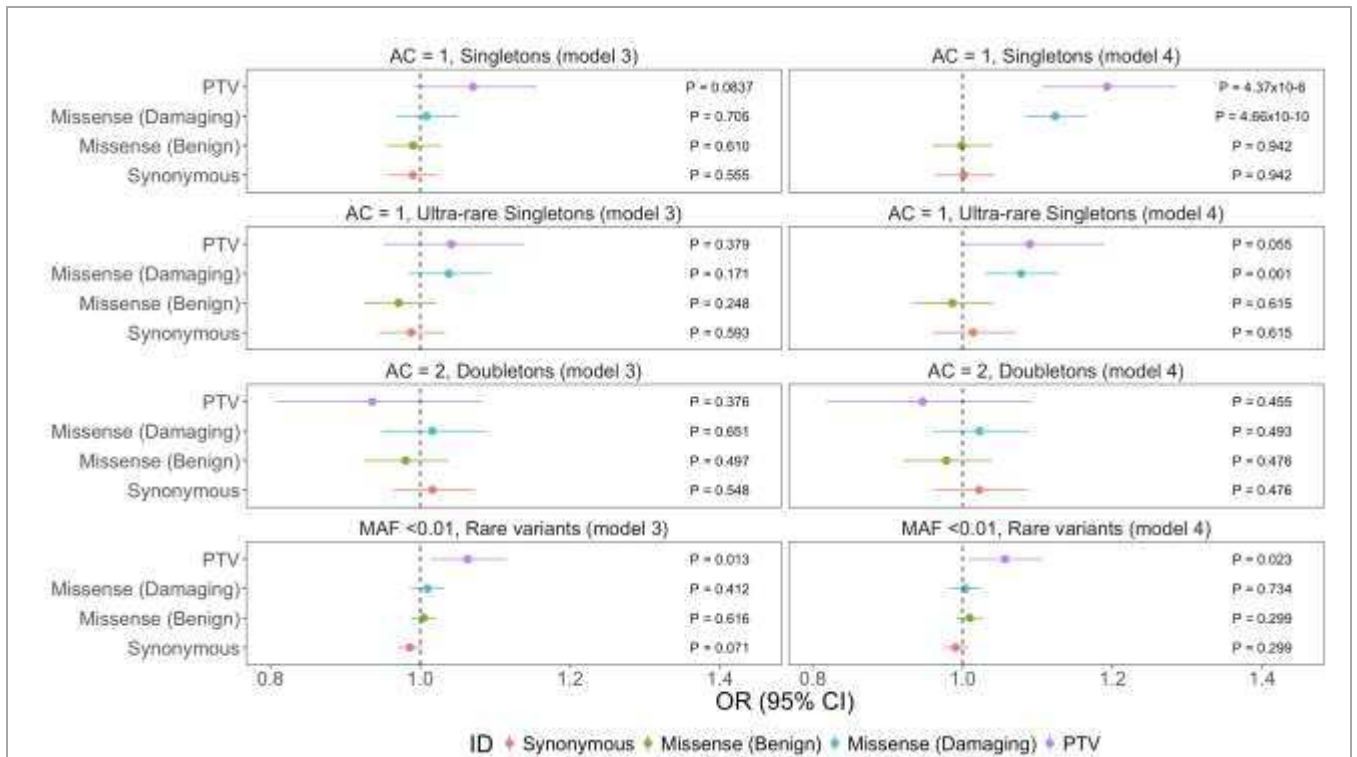
889
 890
 891
 892
 893
 894
 895



Supplementary Fig. 7. Analysis of other neurodegenerative disease genes

Extension of Fig. 3B: Evaluating the effects of genes associated with other neurodegenerative disease (motor neuron diseases: primary lateral sclerosis, progressive muscular atrophy, progressive bulbar palsy, and spinal muscular atrophy; diseases with overlapping phenotypes: frontotemporal dementia, Parkinson's disease, Pick's disease, and Alzheimer's disease) in synonymous variants, benign missense variants, damaging missense variants, and PTVs within singletons (AC=1), doubletons (AC=2), ultra-rare singletons (AC=1, 0 in DiscovEHR), and rare variants (MAF<0.01 in our dataset, DiscovEHR and ExAC). Odds ratios and 95% confidence intervals for each class of variation are depicted by different colors. P-values are also displayed. Model 3 evaluates sample variation with the covariates, sample sex, PC1-10, and total exome count (summation of synonymous variation, benign missense variation, damaging missense variation, and PTV). Model 4 evaluates sample variation with the covariates, sample sex, PC1-10, and benign variation (summation of synonymous and benign missense variation). The graph display the mean and standard deviation. P-values from firth logistic regression test are also displayed. Multiple test correction P-value: 0.0125. N=3,864 ALS cases; N=7,839 controls.

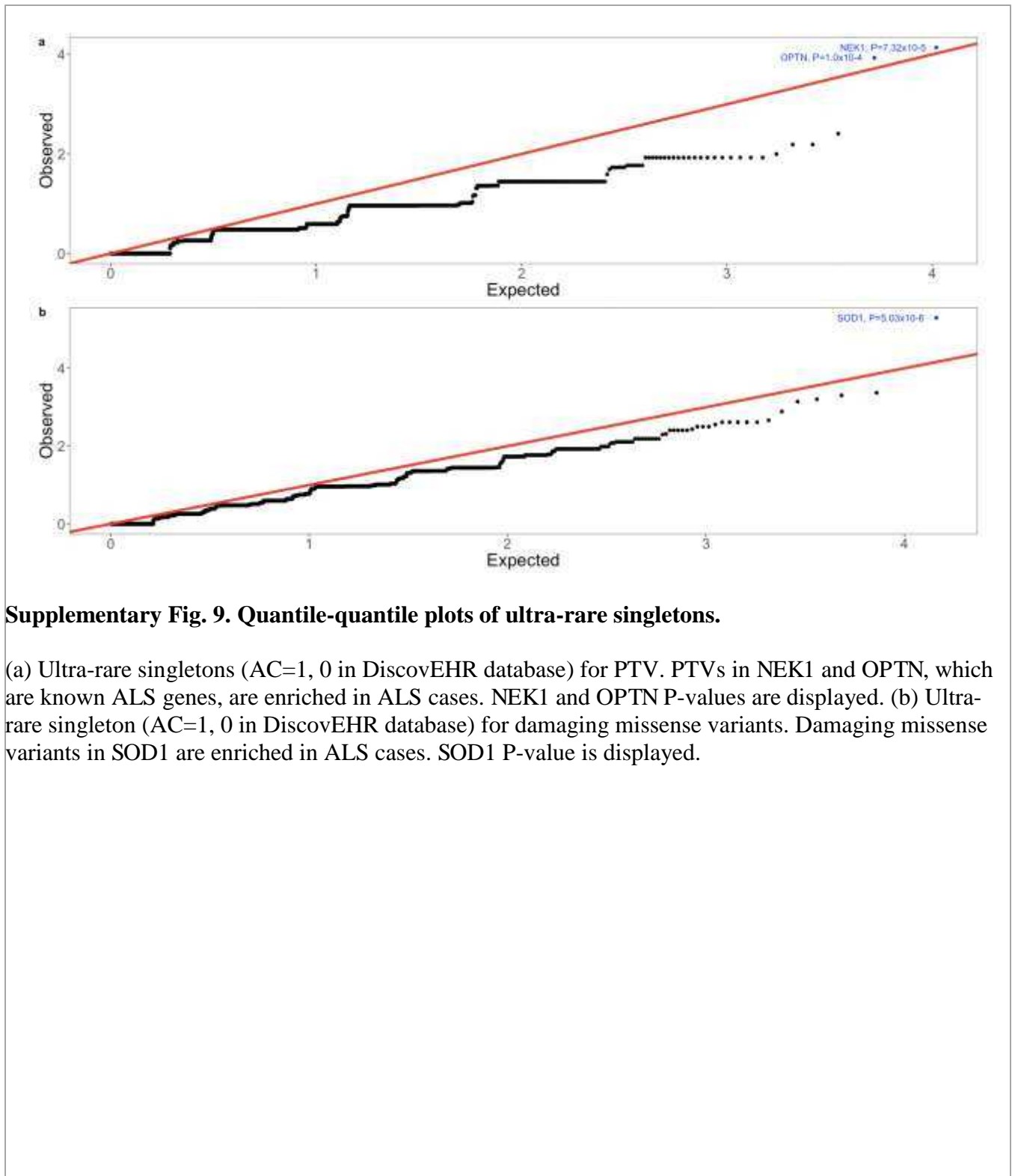
896
897
898
899

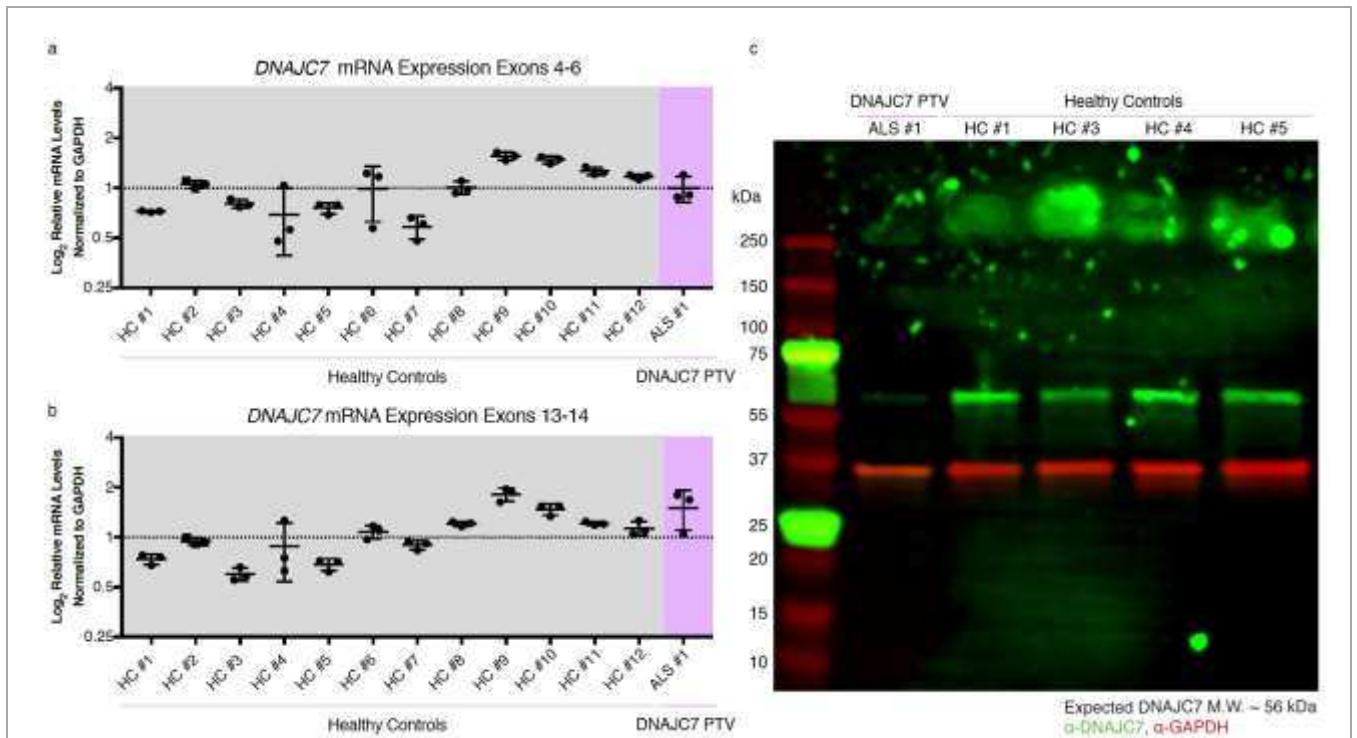


Supplementary Fig. 8. Analysis of brain specific genes

Extension of Fig. 3C: Analysis of brain specific genes in synonymous variants, benign missense variants, damaging missense variants, and PTVs within singletons (AC=1), doubletons (AC=2), ultra-rare singletons (AC=1, 0 in DiscovEHR), and rare variants (MAF<0.01 in our dataset, DiscovEHR and ExAC). Odds ratios and 95% confidence intervals for each class of variation are depicted by different colors. P-values are also displayed. Model 3 evaluates sample variation with the covariates, sample sex, PC1-10, and total exome count (summation of synonymous variation, benign missense variation, damaging missense variation, and PTV). Model 4 evaluates sample variation with the covariates, sample sex, PC1-10, and benign variation (summation of synonymous and benign missense variation). The graph display the mean and standard deviation. P-values from firch logistic regression test are also displayed. Multiple test correction P-value: 0.0125. N=3,864 ALS cases; N=7,839 controls.

900
901
902
903
904
905





Supplementary Fig. 10. DNAJC7 qPCR and immunoblot assays

(A-B) Relative levels of DNAJC7 mRNA in human fibroblasts using either primers recognizing exons 4 and 6 (A) or exons 13 and 14 (B). Levels for each sample were normalized to GAPDH and displayed relative to the average normalized levels of the healthy controls. Data are displayed as the mean of technical replicates with SD. (C) Uncropped immunoblot of human fibroblast protein lysates probed for the N-terminus of DNAJC7. Similar results were obtained in n=3 independent blots.

907
908
909
910
911