# Regional collapsing of rare variation implicates specific genic regions in ALS

— **Source link** ⧉

Sahar Gelfman, Sarah A. Dugger, Cristiane de Araújo Martins Moreno, Zhong Ren ...+15 more authors

**Institutions:** Columbia University Medical Center, McGill University, University of Massachusetts Medical School, Stanford University ...+2 more institutions

Related papers:

- A new approach for rare variation collapsing on functional protein domains implicates specific genic regions in ALS.

- A cross-disorder dosage sensitivity map of the human genome

- The genomic and functional characteristics of disease genes

- Network properties of genes harboring inherited disease mutations

- Linking common and rare disease genetics through gene regulatory networks

# Regional Collapsing of Rare Variation Implicates Specific Genic Regions in ALS

Sahar Gelfman[1], Sarah Dugger[1], Cristiane de Araujo Martins Moreno[2], Zhong Ren[1], Charles J. Wolock[1], Neil A. Shneider[2,3], Hemali Phatnani[1,2,4], Elizabeth T. Cirulli[5], Brittany N. Lasseigne[6], Tim Harris[7], Tom Maniatis[8], Guy A. Rouleau[9], Robert H. Brown Jr.[10], Aaron D. Gitler[11], Richard M. Myers[6], Slavé Petrovski[12], Andrew Allen[13], Matthew B. Harms[1,2,3¥*] and David B. Goldstein[1,14¥*]

[1]Institute for Genomic Medicine, Columbia University Irving Medical Center, New York, NY, 10032, USA

[2]Department of Neurology, Columbia University Irving Medical Center, New York, NY 10032 USA

[3]Motor Neuron Center, Columbia University Irving Medical Center, New York, NY 10032 USA

[4]New York Genome Center, New York, NY 10013 USA

[5]Human Longevity INC, San Diego, CA 92121 USA

[6]HudsonAlpha Institute for Biotechnology, Huntsville, AL 35806 USA

[7]SV Health Investors, Boston, MA 02108 USA

[8]Department of Biochemistry and Molecular Biophysics, Columbia University Irving Medical Center, New York, NY 10032 USA

[9]Department of Neurology and Neurosurgery, McGill University, Montreal, H3A 2B4 Canada

[10]Department of Neurology, University of Massachusetts Medical School, Worcester, MA 01655 USA

[11]Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305 USA

29   [12]Department of Medicine, Austin Health and Royal Melbourne Hospital, University of

30   Melbourne, Melbourne, Australia

31   [13]Department of Biostatistics and Bioinformatics, Duke University, Durham, NC 27708,

32   USA

33   [14]Department of Genetics and Development, Columbia University Irving Medical Center,

34   New York, NY 10032, USA

35   [¥]Authors contributed equally to this work

36   [*]Co-corresponding authors

37

## Abstract

39

40   Large-scale sequencing efforts in amyotrophic lateral sclerosis (ALS) have

41   implicated novel genes using gene-based collapsing methods. However,

42   pathogenic mutations may be concentrated in specific genic regions. To

43   address this, we developed two collapsing strategies, one focuses rare

44   variation collapsing on homology-based protein domains as the unit for

45   collapsing and another gene-level approach that, unlike standard methods,

46   leverages existing evidence of purifying selection against missense

47   variation on said domains. The application of these two collapsing methods

48   to 3,093 ALS cases and 8,186 controls of European ancestry, and also

49   3,239 cases and 11,808 controls of diversified populations, pinpoints risk

50   regions of ALS genes including *SOD1*, *NEK1*, *TARDBP* and *FUS*. While

51   not clearly implicating novel ALS genes, the new analyses not only pinpoint

52   risk regions in known genes but also highlight candidate genes as well.

53

## Introduction

55

56  Amyotrophic lateral sclerosis (ALS) is an adult-onset neurodegenerative

57  disease characterized by progressive motor-neuron loss leading to

58  paralysis and death, most often from respiratory failure. Roughly 60-70% of

59  familial and 10% of sporadic cases have an identifiable mutation in a

60  known causal ALS gene, the majority of which are repeat expansions in

61  *C9ORF72* and point mutations in *SOD1*[1]. Recent efforts in gene discovery,

62  largely driven by advances in sequencing and identification of rare variants,

63  have implicated and confirmed several new genes in ALS pathogenesis

64  including *TBK1*, *NEK1*, *ANXA11* and *CCNF*[2-10]. Despite this progress, the

65  majority of sporadic cases still remain to be resolved genetically.

66

67  The now established paradigm for case-control analyses of exome or

68  genome sequencing data of complex diseases and traits involves a gene-

69  based collapsing framework in which all qualifying variants in a gene are

70  treated as equivalent. Genes are associated with the trait when they exhibit

71  a significant excess of qualifying variants occurring anywhere in the gene.

72  This approach has implicated disease genes in a growing number of other

73  complex conditions beyond ALS, including idiopathic pulmonary fibrosis

74  (IPF), myocardial infarction (MI) and Alzheimer's disease [11-13].

75

76  While clearly effective, the power of this approach is limited by the inclusion

77  of benign variants that reduce statistical power. However, for genes where

78  pathogenic mutations are localized to specific regions, such as functional

79  domains, power can be increased by using these regions as the unit for the

80  collapsing analysis.  In ALS-associated genes, there are several examples

81  of genes that show regionally localized pathogenic variation. For example,

82  in *TARDBP*, highly-penetrant ALS variants are concentrated in a glycine-

83  rich domain near the C-terminus [14]. Furthermore, the gene *FUS*, which has

84  a similar structure as *TARDBP,* has pathogenic mutations clustering in two

85  regions: exons 13–15 (encoding an Arg-Gly-Gly-rich domain and the

86  nuclear localization signal) and exons 3, 5-6 (encoding Gln-Gly-Ser-Tyr-rich

87  and Gly-rich domains) [15].

88

89  Recognizing that undiscovered ALS-associated genes might similarly have

90  specific domains where pathogenic variants cluster, we now apply two

91  complementary regional approaches to gene collapsing analyses to identify

92  localized signals of rare variation in a data set of 3,093 ALS cases of

93  European ancestry (2,663 exomes and 430 whole genomes) compared

94  with 8,186 controls of matched ancestry (7,612 control exomes and 574

95  whole genomes).  We further apply these analyses to a set of samples of

96  diversified ancestry origins, consisting of 3,239 cases and 11,808 controls.

97  We compare the regional approaches to the standard gene collapsing

98  analysis and highlight the importance of a regional view specifically for ALS

99  genetics.

100

101  **Results**

102

103  **Collapsing analyses using homology-defined protein domains**

104

105  The standard approach to gene discovery focuses on the burden of rare

106  variants across an entire gene by comparing the frequency of qualifying

107  variants in cases and controls. The qualifying variants can be defined by

108  various criteria such as function and allele frequency (illustrated in figure

109  1A).

110

111 In this study, we describe two additional approaches to rare variant
112 collapsing: 1) a regional approach, where the unit for collapsing is not the
113 gene, but rather the functional domains within the gene (figure 1B), and 2)
114 a gene-based approach where the definition of qualifying variants is
115 informed by regional intolerance to missense variation (figure 1C).

116

117 *Figure 1*

118

119 We first utilized the standard gene collapsing approach (figure 1A) to
120 identify the burden of rare variants in a set of 3,093 ALS cases and 8,186
121 controls of European ancestry. The demographic features of our cohort
122 reflect known epidemiological features of ALS, including male
123 predominance and the distributions of age at onset and survival
124 (supplementary Table S1). Qualifying variants were defined as non-
125 synonymous coding or canonical splice variants that have a minor allele
126 frequency (MAF) ≤0.1% in cases and controls (internal MAF) and also a
127 ≤0.1% MAF imposed for each population represented in the ExAC
128 database[16]. High quality control (QC) metrics were further imposed on the
129 variants (see Methods).

130 Comparing genetic variation across 18,653 protein-coding genes found a
131 genome-wide and study-wide significant (p<4.5x10-7) case-enrichment
132 only for *SOD1* ( p=$1.23\times10^{-18}$, figure 2A), with qualifying variants
133 identified in 43 cases (1.39%) and only 6 controls (0.07%; OR=19.2).
134 *TARDBP* showed the second strongest enrichment (OR=3.6,
135 p=$1.02\times10^{-4}$), but with 23 cases (0.74%) and 17 controls (0.21%) it did not

136   achieve genome-wide significance. *FUS* harbored qualifying variants in 20

137   cases and 37 controls (OR=1.43, p=0.23, figure 2A).

138

139   A gene-based analysis evaluating only rare loss of function (LoF) variants

140   was also performed, identifying a genome-wide and study-wide significant

141   case-enrichment of *NEK1* variants (OR=7.35, p=$1.85\times10^{-10}$), with 33

142   cases (1.07%) compared to 12 controls (0.15%, supplementary figure S1).

143   As a negative control, we included a model for rare synonymous variants,

144   and did not observe any genes with significant enrichment. The genomic

145   inflation factor, lambda (λ) for this model was 1.03 (supplementary figure

146   S2).

147

148   We hypothesized that genes with clustered mutations that had weak

149   enrichments using this standard gene-based collapsing approach, such as

150   *TARDBP* and *FUS*, could be identified by a collapsing method that uses

151   functional gene regions (i.e. domains) as the unit for collapsing (figure 1B).

152   For this analysis, we utilized a list of 89,522 gene domains covering the

153   human coding sequence, as described previously[17]. In short, the coding

154   sequence of each gene was aligned to a set of conserved protein domains

155   based on the Conserved Domain Database (CDD)[18]. The final domain

156   coordinates for each gene were defined as the regions within the gene that

157   aligned to the CDD and the unaligned regions between each CDD

158   alignment.   These domains were then used as the unit for collapsing

159   compared with a standard gene-based collapsing approach (figure 1A and

160   1B).

161

162   This domain-based analysis was performed using the same cohort and

163  coding model as the standard approach (European ancestry, non-
164  synonymous and canonical splice variants, internal and population MAF
165  ≤0.1%). As hypothesized, the top three case-enriched domains reside in
166  ALS genes: *SOD1*, *TARDBP* and *FUS* (figure 2B). For *SOD1*, a long
167  domain spanning the majority of the coding sequence contains most of the
168  variation found in 1.29% of cases and 0.07% of controls (OR=17.9;
169  $p=4.1\times10^{-17}$, figure 2B).

170

171  Strikingly, the glycine-rich *TARDBP* domain where known mutations cluster
172  is now identified with genome-wide and study-wide significance (OR=7;
173  $p=5.84\times10^{-7}$). Of note, this glycine-rich domain covers exon 6 of *TARDBP*
174  and was not mapped to a conserved domain from the CDD.

175

176  The same trend was observed for *FUS*, which shows the third strongest
177  enrichment in this analysis (OR=8.6; $p=3.6\times10^{-5}$, figure 2B). Specifically,
178  qualifying variants were identified in 13 cases (0.42%) and 4 controls
179  (0.05%) in the previously reported Arg-Gly rich domain covering exons 13-
180  15, which is also not considered a conserved CDD domain [14]). Although not
181  at genome-wide or study-wide significance, this represents a substantial
182  improvement over the gene-based collapsing approach (OR=1.43,
183  uncorrected p=0.23).

184

185  The fourth most case-enriched domain was a conserved armadillo repeat
186  domain spanning exons 12-14 of *PKP4* (plakophilin 4, also known as
187  p0071). Qualifying variants occurred in 0.61% of ALS cases and 0.13% of
188  controls (OR=4.6, $p=4.1\times10^{-5}$). While not genome-wide significant, *PKP4*

189  is an interesting candidate gene that has been previously linked to various

190  ALS-related pathways (see Conclusions).

191

192  *Figure 2*

193

194  **Gene-wide collapsing analyses informed by regional intolerance to**
195  **missense variation**

196

197  As we have demonstrated, domain-based collapsing effectively identifies

198  genes where pathogenic variants are localized to single specific regions

199  (e.g *TARDBP* and *FUS*), and highlights suggestive candidates for further

200  study (*PKP4*). However, to identify haploinsufficient genes where truncating

201  variants and sufficiently damaging missense mutations could both

202  contribute to risk of disease, the difficulty lies in determining which

203  missense variants should qualify in the analysis. To address this challenge,

204  we implemented a collapsing approach that leverages regional patterns of

205  intolerance to missense variation (sub-RVIS[17,19]) as a way to prioritize

206  missense variants most likely to result in disease. In this 'intolerance-

207  informed' approach, rare missense alleles were counted as qualifying if

208  they resided in gene regions that are intolerant to missense variation,

209  whereas LoF variants were counted as qualifying regardless of location

210  within the gene (figure 1C).

211

212  As a measure of intolerance of gene regions, we applied a complementary

213  approach to subRVIS[17] for when there is limited resolution in the sequence

214  region of interest. This approach uses the observed to expected missense

215  ratio in a domain (OE-ratio), which is equivalent to a domain-based

216  missense tolerance ratio (MTR)[19]. In short, the expected rate leverages the

217 underlying sequence context in the domain, and the observed rate is based

218 on the rate of non-synonymous variants identified in the sub-region of

219 interest based on the ExAC database, release 0.3[16] (see Methods).

220

221 We focus our intolerance-informed gene collapsing approach on domains

222 that have intolerance below the exome-wide 50th percentile (OE-Ratio

223 percentile <50%), thus sub-selecting variants in genic regions that have

224 greater evidence of purifying selection acting against non-synonymous

225 variation. As mentioned earlier, for each gene, variants in these intolerant

226 regions are considered along with LoF variants independent of their

227 location within the gene.

228

229 Because intolerant coding regions are expected to have a lower rate of

230 common variation, we included samples from diversified ancestries when

231 applying intolerance-informed gene collapsing.  The diversified population

232 approach increased the total number of samples by 3,768, to 3,239 cases

233 and 11,808 controls, thereby increasing the power of the analysis. For this

234 approach, we applied similar rules for qualifying variants, including low

235 population frequency (MAF ≤0.1% imposed for each population

236 represented in ExAC), an internal MAF ≤0.04% (decreased from 0.1% due

237 to a larger control cohort), coding annotation (non-synonymous and splice

238 variants) and high QC metrics, with the additional criteria of residing in the

239 lower 50th percentile of OE-ratio domains.

240

241 The genomic inflation factor (λ) of the diversified populations intolerance-

242 informed analysis was 1.14, slightly higher than the European-only cohort

243 used for the standard gene level analyses (figure 2A, λ=1.1).  Yet, this

244  inflation is much lower than for the standard gene based analysis using a

245  diversified population (λ=1.25, supplementary figure S3), demonstrating the

246  advantage of an intolerance-informed approach for reducing the genomic

247  inflation due to variation in tolerant regions.

248

249  In this analysis, *SOD1* achieved a slightly better enrichment than in either

250  gene-based or domain-based analyses (OR=20.31; p=$4.13\times10^{-22}$, figure

251  3A).

252  *TARDBP* also had genome-wide and study-wide significant enrichment

253  (OR=4.95; p=$8.77\times10^{-8}$; figure 3A), which presents a considerably-

254  strengthened enrichment signal for *TARDBP* compared to the standard

255  gene collapsing analysis observed in figure 2A (OR=3.6, uncorrected

256  p=$1.02\times10^{-4}$).

257

258  *LGALSL* (Galectin-like or Lectin, Galactoside-Binding, Soluble, Like) was

259  the third gene to have a strong enrichment of qualifying variants in cases

260  (OR=14.63; p=$2.29\times10^{-6}$; figure 3A) that was not study-wide significant

261  given the models tested. The enrichment of this gene originates from one

262  specific domain that harbors variants for twelve cases (0.37%) and three

263  controls (0.025%) with the addition of an African American and a Latino

264  case over the European-only analysis (figure 2). The target domain

265  harboring all *LGALSL* case-variants is a region comprising 378bp that is

266  mapped to a conserved protein domain intolerant to variation.  Notably,

267  *LGALSL* LoF variants were only identified in cases and absent from nearly

268  12,000 controls.  To assess the rate of LoF variants in a larger control

269   population we looked at the ExAC cohort and found three LoF alleles in

270   60,033 individuals[16].

271

272   *Figure 3*

273

274   **Genome wide associations with age-at-onset**

275

276   We next examined whether qualifying variants in known ALS genes, or

277   candidate genes identified by our novel approaches, influence age at

278   symptom onset (AAO).  We found that *SOD1* variant carriers tended to be

279   younger than the rest of the cohort (52.2 vs. 57.1 years, p=0.059; MW-

280   test). Also, subjects harboring qualifying variants in *ANXA11* showed

281   delayed onset (63.8 years, p=0.037; MW-test), which is consistent with

282   prior studies[9]. No other known ALS genes showed significant influence on

283   AAO.

284   Interestingly, subjects harboring *LGALSL* qualifying variants showed a

285   mean AAO that is 13 years younger than the rest of the cohort (43.8 years

286   vs. 57.1, p=$8.1 \times 10^{-4}$; Mann-Whitney test). The AAO information was

287   available for 11/12 variant carriers and 2,767/3,239 non-carriers.

288

289   The early onset in cases carrying *LGALSL* variants was further validated by

290   a random sampling approach where *LGALSL* carriers' average AAO was

291   significantly lower than 9,983/10,000 randomly sampled sets of 11 cases

292   (p=0.0017, supplementary Methods).

293

294   *Figure 4*

295

## Discussion

297

Here we present a regional approach to rare variant collapsing analyses. This approach has two distinct forms: 1) aggregating rare variants on genic sub-regions defined using conserved protein domain annotations, and 2) aggregating rare variants on a gene unit but using the pattern of purifying selection to identify the most damaging missense variants and combine them with loss of function mutations occurring anywhere in the gene. Both approaches show improved sensitivity for known ALS genes, finding *SOD1*, *NEK1* and *TARDBP* as genome-wide significant. We also find *FUS's* Arg-Gly-rich domain within the top three associations in our domain-based regional collapsing, jumping from an insignificant OR=1.43 to a high OR=8.6. These findings underscore the utility of applying a regional approach to ALS genetics, especially in light of similar Gly-rich domains importance in mediating pathologic RNA-protein complexes [20].

311

This approach has also implicated a potential new candidate ALS gene, *LGALSL*, encoding the galectin-like protein GRP (galectin-related protein, also known as HSPC159 and lectin galactoside-binding-like protein). We identified a case-enriched intolerant galectin-binding domain (Figure 4A). Interestingly, while the functions of *LGALSL* remain largely unknown, members of the galectin family, including *LGALS1* and *LGALS3,* have been implicated in ALS disease processes and progression. Specifically, *LGALS1* has been identified as a component of sporadic and familial ALS-related neurofilamentous lesions[21], and is associated with early axonal degeneration in the *SOD1^{G93A}* ALS mouse model[22]. Furthermore,

322 homozygous deletion of *LGALSL3* reportedly led to accelerated disease

323 progression and reduced lifespan in *SOD1*$^{G93A}$ mice[23]. We also performed

324 an analysis of age at onset that is independent of the predisposition

325 analysis, showing a significant association between carriers of qualifying

326 variants in *LGALSL* and early age at onset. This independent analysis

327 supports *LGALSL* as a candidate ALS gene that might be responsible for a

328 form of ALS with a younger age at onset. Because *LGALSL* carriers were

329 contributed by five different sites and sequenced at two sequencing

330 centers, we cannot exclude hidden variable stratification that might explain

331 the low AAO. Further study of additional *LGALSL* mutation carriers will be

332 required to confirm this observed genotype-phenotype correlation.

333

334 Regional collapsing analyses also highlighted *PKP4* as a new candidate

335 gene, with a single armadillo repeat domain strongly enriched for qualifying

336 variants in cases (Figure 4B). Evidence supporting *PKP4*'s role in ALS-

337 linked processes including microtubule transport and endosomal

338 processing, in addition to its local translation in ALS-mutant *FUS* granules,

339 all provide evidence in favor of *PKP4* as a risk factor for ALS[24-27].

340

341 This study incorporated both exome and whole genome samples from a

342 large cohort of over 3,000 cases and close to 12,000 controls. Yet, despite

343 these large cohorts, the standard gene collapsing approach identified only

344 *SOD1* and *NEK1* (loss-of-function specific model) as achieving genome-

345 wide significance, and failed to uncover other known signals for ALS risk

346 factors. We were able to capture these signals, along with candidate novel

347 signals, using a regional approach that is informed by missense variation

348 intolerance. That being said, while confirming *TARDBP*, and suggesting

349  *LGALSL* (0.42% of cases) as a candidate gene, the regional approach was

350  still underpowered with the current sample size to show genome-wide

351  significance for *FUS* and *PKP4* that might reflect true associations. This

352  suggests that even with these signal optimization approaches, larger

353  sequencing studies are required in ALS. We are, however, confident that

354  the continued application of regional approaches to collapsing analyses in

355  ALS and other rare disorders will enable the identification of novel risk

356  factors with small proportions in patient populations, that were previously

357  difficult to identify due to being masked by benign variants in regions that

358  are tolerant to variation.

359

360  **Methods**

361

362  **Subject sources**

363

364  ALS samples analyzed by whole exome or genome sequencing came from

365  the Genomic Translation for ALS Care (GTAC study), the Columbia

366  University Precision Medicine Initiative for ALS, the New York Genome

367  Consortium, and the ALS Sequencing Consortium (IRB-approved genetic

368  studies from Columbia University Medical Center (including the Coriell

369  NINDS repository), University of Massachusetts at Worchester, Stanford

370  University (including samples from Emory University School of Medicine,

371  the Johns Hopkins University School of Medicine, and the University of

372  California, San Diego), Massachusetts General Hospital Neurogenetics

373  DNA Diagnostic Lab Repository, Duke University, McGill University

374  (including contributions from Saint-Luc and Notre-Dame Hospital of the

375  Centre Hospitalier de l'Université de Montréal (CHUM) (University of

376 Montreal), Gui de Chauliac Hospital of the CHU de Montpellier (Montpellier

377 University), Pitié Salpe☐trière Hospital, Fleurimont Hospital of the Centre

378 Hospitalier Universitaire de Sherbrooke (CHUS) (University of Sherbrooke),

379 Enfant- Jésus Hospital of the Centre hospitalier affilié universitaire de

380 Québec (CHA) (Laval University), Montreal General Hospital, Montreal

381 Neurological Institute and Hospital of the McGill University Health Centre),

382 and Washington University in St. Louis (including contributions from

383 Houston Methodist Hospital, Virginia Mason Medical Center, University of

384 Utah, and Cedars Sinai Medical Center).

385

386 **Subject selection criteria**

387

388 ALS subjects were diagnosed according to El Escorial revised criteria as

389 suspected, possible, probable, or definite ALS by neuromuscular

390 physicians at submitting centers. Subjects were considered sporadic if no

391 first or second-degree relatives had been diagnosed with ALS or died of an

392 ALS-like syndrome. Because screening for known ALS gene mutations

393 prior to sample submission was highly variable across the cohort, gene

394 status was not considered *a priori*. Controls were selected from >45,000

395 whole exome or genome sequenced individuals housed in the IGM Data

396 Repository. We excluded all individuals with a known diagnosis or family

397 history of neurodegenerative disease, but not all had been specifically

398 screened for ALS.

399

400 **Sequencing**

401

402 Sequencing of DNA was performed at Columbia University, the New-York

403 Genome Center, Duke University, McGill University, Stanford University,

404 HudsonAlpha, and University of Massachusetts, Worcester. Whole exome

405 capture used Agilent All Exon kits (50MB, 65MB and CRE), Nimblegen

406 SeqCap EZ Exome Enrichment kits (V2.0, V3.0, VCRome and

407 MedExome), IDT Exome Enrichment panel and Illumina TruSeq kits.

408 Sequencing occurred on Illumina GAIIx, HiSeq 2000, HiSeq 2500, or HiSeq

409 X sequencers according to standard protocols.

410

411 Illumina lane-level fastq files were aligned to the Human Reference

412 Genome (NCBI Build 37) using the Burrows-Wheeler Alignment Tool

413 (BWA)[28]. Picard software (http://picard.sourceforge.net) removed duplicate

414 reads and processed lane-level SAM files to create a sample-level BAM

415 file. Genomes (n=402) from the New York Genome Center were

416 transferred as sample-level BAM files. We used GATK to recalibrate base

417 quality scores, realign around indels, and call variants[29].

418

419 **Samples quality control**

420

421 The initial sample consisted of 4,149 ALS cases and 15,107 controls.

422 Samples reporting >8% contamination according to VerifyBamID[30] were

423 excluded. KING[31] was used to ensure only unrelated (up to third-degree)

424 individuals contributed to the analysis. For controls, where sample

425 collection methods were not known, we excluded samples where X:Y

426 coverage ratios did not match expected sex. For studies where sample

427 collection and processing involved only ALS patients, mismatches were not

428 exclusionary. Further, to be eligible, samples were further subjected to a

429  CCDS 10-fold coverage principal components analysis (PCA) and an
430  ancestry prediction filter (for European ancestry analysis).

431

432  **Cohort construction: ancestry prediction**

433

434  The ancestry classification model was trained using genotyped data from
435  5,287 individuals of known ancestry and 12,840 well-genotyped and
436  ancestry-informative markers that were limited to the human exome. The
437  model was trained, tested and validated on a set of individuals with
438  ancestry as follows: non-Finnish European(N=2911), Middle Eastern
439  (N=184), Hispanic (N=368), East Asian (N=539), South Asian (N=529) and
440  African (N=756). Briefly, the sample × genotype matrix was scaled to have
441  unit mean and standard deviation along each SNV and subjected to a
442  principal component analysis (PCA). For training the classifier, the
443  genotypes were projected onto the top 6 PCs and used as feature vectors.
444  The classifier is a Multi-layer perceptron with 1 hidden layer, a logistic
445  activation function, L2 regularization term, alpha = 1e-05, size of hidden
446  layer = 6 and a L-BFGS solver. The classifier was implemented using the
447  scikit-learn API in Python. A stratified 10-Fold CV with 80:20 split of the
448  training data was used to tune parameters using a grid search. Cross
449  validation performance on the cohort yielded precision/recall scores as
450  follows: NFE: 0.99/1, AFR: 0.99/1, SAS 0.99/1, EAS: 0.99/1, HIS:0.93/0.97,
451  ME: 0.93/0.77.  Samples in this study were subjected to ancestry prediction
452  using the model trained above by projecting their genotype vector to the
453  training PCA model and running the classifier to obtain a given sample's
454  ancestry probabilities for each of the trained population.

455

456 For samples to qualify for a European ancestry analysis, they were
457 required to have a European probability greater than 0.5 and an overall
458 genotyping rate of 0.87 across the 12,840 well-genotyped and ancestry
459 informative markers. Lower genotyping rates were considered as
460 uninformative for ancestry prediction. In the case of low genotyping rate,
461 we considered self-declared ethnicity of 'White' as qualifying for the
462 European based analysis.

463

464 Furthermore, once the final list was constructed, we applied an additional
465 analysis to control for population stratification by using EIGENSTRAT [32] to
466 remove samples that were considered as genetic outliers, this ensured that
467 the main cluster of samples was of European origins (see below).

468

469 **Cohort construction: CCDS coverage PCA**

470

471 To account for the variability of the coverage over the CCDS between
472 samples originating from various sequencing kits and platforms, we
473 developed a method to remove samples that are considered outliers due to
474 coverage. This step was performed for samples that passed QC and
475 ancestry prediction filters (if applied), and allowed for maximizing the
476 coding region available for the analysis when harmonizing variant level
477 coverage between cases and controls.

478

479 We first randomly selected a set of 1,000 CCDS genes for a coverage test.
480 We next constructed a table where the rows are the samples used for the
481 analysis and the columns are the number of based covered at 10x in each
482 of the 1,000 random genes. Finally, we used the coverage table in a

principal-component analysis. Outliers were identified as being further than three standard deviations away from the center of the first four principal components (PCs).

In the Caucasian analysis 3,866 cases and 9,426 controls passed initial QC and ancestry filters and were subjected to the coverage PCA filter. The coverage PCA maintained 3,314 cases and 9,214 controls.

In the diversified population analysis 4,075 cases and 14,494 controls passed initial QC and were subjected to the coverage PCA filter. The coverage PCA maintained 3,468 cases and 13,957 controls.

**Cohort construction: Eigenstrat PCA threshold adjustment**

EIGENSTRAT[32] PCA was used for removing genotypic outlier samples as a final cohort pruning step before running the collapsing analysis. The default threshold for removing outliers is six standard deviations from mean over the top ten PCs. This process, including recalculation of the PCs, was repeated five times.

In the Caucasian analysis 3,208 cases and 8,821 controls passed initial QC, ancestry, coverage PCA and kinship filters and were subjected to the final EIGENSTRAT PCA filter. The EIGENSTRAT PCA maintained the final 3,093 cases and 8,186 controls used for the collapsing analysis, including 383 out of 420 whole genome cases that were mapped by the New-York Genome Center (NYGC).

In the diversified analysis 3,353 cases and 13,373 controls passed initial QC, coverage PCA and kinship filters and were subjected to the final

510    EIGENSTRAT PCA filter. The default EIGENSTRAT PCA threshold

511    removed all 420 NYGC whole genomes. This was the result of the addition

512    to the Caucasian analysis of over 3,000 exomes, which reduced the

513    standard deviation and resulted in the exclusion of NYGC whole genomes

514    in the third PC. As these samples were very high quality and were included

515    in the Caucasian only analysis, we adjusted the threshold of the third PC to

516    seven standard deviations, thus maintaining 402 out of 420 NYGC whole

517    genomes. In total, following stratification phase, we maintained 3,239 cases

518    and 11,808 controls for the collapsing analysis.

519

520    **Variant-level quality control**

521

522    Quality thresholds were set based on previous studies [3,33]. Variants were

523    required to have a quality score of at least 30, quality by depth score of at

524    least 2, genotype quality score of at least 20, read position rank sum of at

525    least -3, mapping quality score of at least 40, mapping quality rank sum

526    greater than -10, and a minimum coverage of at least 10. SNVs had a

527    maximum Fisher's strand bias of 60, while indels had a maximum of 200.

528    For heterozygous genotypes, the alternative allele ratio was required to be

529    greater than or equal to 25%. Variants were excluded if they were marked

530    by EVS as being failures [34]. Variants were annotated to Ensembl 73 using

531    SnpEff [35].

532

533    **Variant-level statistical analysis**

534

535    Our primary model was designed to search for non-synonymous coding or

536    canonical splice variants that have a less than 12 cases with a recurring

537   variant in cases and controls (internal MAF) and also a ≤0.1% MAF

538   imposed for each population represented in the ExAC database[16].

539   We've tested this model in three forms: a standard gene-unit collapsing

540   analysis, a domain-unit analysis and an intolerance-informed gene

541   collapsing analysis. A further gene-based analysis evaluating only rare loss

542   of function (LoF) variants was also performed.

543   For each of the four models we tested the list of 18,653 CCDS genes. For

544   each gene, we counted the presence of at least one qualifying variant in

545   the gene. A two-tailed Fisher's exact test (FET) was performed for each

546   gene to compare the rate of cases carrying a qualifying variant compared

547   to the rate in controls. For our study-wide significance threshold, after

548   Bonferroni correction for the number of genes tested across the four non-

549   synonymous models, the study-wide multiplicity-adjusted significance

550   threshold $\alpha = (0.05/ [4*18653]) = 6.7 \times 10^{-7}$. We did not correct for the

551   synonymous (negative control) model.

552

553   **OE-ratio intolerance for coding domains**

554

555   The OE-ratio is calculated using the same approach as the missense

556   tolerance ratio (MTR) that is described by Traynelis et al [19]. This approach

557   uses the observed to expected missense ratio for the 89,522 domain

558   coordinates that are described by Gussow et al.[17].

559   For calculating a domain OE-ratio, the following requirements are applied:

560   1) adequate coverage - at least 50% of the bases within the domain must

561   have at least a 10-fold coverage in the ExAC database, release 0.3[16]. 2) At

562   least five distinct variants (of any annotation) are required to perform a

563   binomial exact test depletion of missense at uncorrected alpha of p<0.05.

There were 67,890 domains that passed the above requirements and were scored for their OE-ratio. The average size of the remaining unscored domains was usually very short (mean=21bp; median=12) and they accounted for 0.77% of the protein-coding exome. Unscored domains were considered as below the intolerance ratio required for the intolerance-informed analysis (figure 3) to prevent loss of gene level information. Once a domain lacking OE-ratio is implicated in an analysis, its intolerance is examined using the average missense intolerance ratio (MTR)[19] of the domain in question (http://mtr-viewer.mdhs.unimelb.edu.au).

In the case of *LGALSL*, the last three codons of the coding transcript are a short independent domain that was not mapped to a conserved domain from the CDD. However, this small region is still considered part of the gal-binding domain by other databases [36]. The MTR score for these three codons is below the 30th percentile of intolerance, marking this region at least as intolerant as the implicated galectin binding domain (OE-ratio percentile of 37).

**Figure and table legend**

**Figure 1. Gene and Regional Collapsing. (A)** A standard gene-based approach for collapsing analysis of non-synonymous and canonical splice rare variants in cases (green) and controls (black) on example *Gene A*. **(B)** A domain-unit based regional approach where only the domains that are intolerant to functional variation are considered as units for collapsing. **(C)** Intolerance informed gene collapsing: a regional approach to gene

592 collapsing where the unit for collapsing is the entire gene, yet missense
593 variants only qualify for the analysis if they reside in domains that are
594 intolerant to variation (domain 2). Loss-of-function variants (big circles)
595 continue to qualify regardless of whether they reside in a tolerant or
596 intolerant domains of the gene. Bright blue background marks qualifying
597 variants.

598

599 **Figure 2. Q-Q plots of gene and domain level collapsing. (A)** The
600 results for a standard gene level collapsing of 3,093 cases and 8,186
601 controls. 18,065 covered genes passed QC with more than one case or
602 control carrier for this test. The genes with the top associations and *FUS*
603 gene are labeled. The genomic inflation factor, lambda ($\lambda$), is 1.10. **(B)** The
604 results for the domain-based collapsing of 3,093 cases and 8,186 controls.
605 70,603 covered domains passed QC with more than one case or control
606 carrier for this test. The genes with the top associations are labeled and
607 genome-wide significant genes are in bold. $\lambda=1.046$.

608

609 **Figure 3. Intolerance informed gene level collapsing with**
610 **unified/diversified ancestry samples. (A)** A q-q plot presenting the
611 results of the gene-based intolerance-informed collapsing of 3,239 cases
612 and 11,808 controls from diversified ancestries. Missense variants are
613 aggregated only if they reside in an intolerant domain that is lower than $50^{th}$
614 percentile OE-ratio score, while loss-of-function variants are aggregated
615 independent of location. 17,795 genes passed QC with more than one case
616 or control carrier for this test. The genes with the top associations are
617 labeled. $\lambda=1.14$. (**B**) A q-q plot of a gene-based intolerance-informed

618    collapsing of 3,093 cases and 8,186 controls of European ancestry. 18,135

619    genes passed QC with more than one case or control carrier for this test.

620    The genes with the top associations are labeled and genome-wide

621    significant genes are in bold. $\lambda = 1.073$.

622

623    **Figure. 4. Distribution of functional coding variants across *LGALSL***

624    **and *PKP4*.** The distribution of *LGALSL* **(A)** and *PKP4* **(B)** coding variants

625    across domains (*LGALSL* transcript ENST00000238875 and *PKP4*

626    transcript ENST00000389757). The y-axis corresponds to the total number

627    of variants identified at a specific location. The blue boxes highlight the **(A)**

628    *LGALSL* carbohydrate binding domain and **(B)** *PKP4* armadillo repeat

629    domain 2 (ARM2) found to be enriched for variants in cases (green)

630    compared to controls (black). Each domain's OE-ratio percentile is marked

631    above for both tolerant (bright blue) and intolerant (orange) domains.

632

## Acknowledgements

633

634

715

## References

717

718
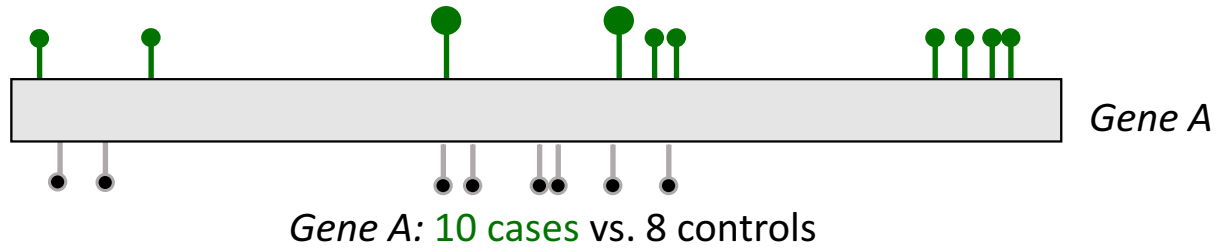
719    1.    Renton, A.E., Chio, A. & Traynor, B.J. State of play in amyotrophic
720          lateral sclerosis genetics. *Nat Neurosci* **17**, 17-23 (2014).
721    2.    Bannwarth, S. *et al.* A mitochondrial origin for frontotemporal
722          dementia and amyotrophic lateral sclerosis through CHCHD10
723          involvement. *Brain* **137**, 2329-45 (2014).
724    3.    Cirulli, E.T. *et al.* Exome sequencing in amyotrophic lateral sclerosis
725          identifies risk genes and pathways. *Science* **347**, 1436-41 (2015).
726    4.    Johnson, J.O. *et al.* Mutations in the Matrin 3 gene cause familial
727          amyotrophic lateral sclerosis. *Nat Neurosci* **17**, 664-666 (2014).
728    5.    Kenna, K.P. *et al.* NEK1 variants confer susceptibility to amyotrophic
729          lateral sclerosis. *Nat Genet* **48**, 1037-42 (2016).

6.  Mackenzie, I.R. *et al.* TIA1 Mutations in Amyotrophic Lateral Sclerosis and Frontotemporal Dementia Promote Phase Separation and Alter Stress Granule Dynamics. *Neuron* **95**, 808-816 e9 (2017).

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

8.  Smith, B.N. *et al.* Exome-wide rare variant analysis identifies TUBA4A mutations associated with familial ALS. *Neuron* **84**, 324-31 (2014).

9.  Smith, B.N. *et al.* Mutations in the vesicular trafficking protein annexin A11 are associated with amyotrophic lateral sclerosis. *Sci Transl Med* **9**(2017).

10. Williams, K.L. *et al.* CCNF mutations in amyotrophic lateral sclerosis and frontotemporal dementia. *Nat Commun* **7**, 11253 (2016).

11. Cruchaga, C. *et al.* Rare coding variants in the phospholipase D3 gene confer risk for Alzheimer's disease. *Nature* **505**, 550-554 (2014).

12. Do, R. *et al.* Exome sequencing identifies rare LDLR and APOA5 alleles conferring risk for myocardial infarction. *Nature* **518**, 102-6 (2015).

13. Petrovski, S. *et al.* An Exome Sequencing Study to Assess the Role of Rare Genetic Variation in Pulmonary Fibrosis. *Am J Respir Crit Care Med* **196**, 82-93 (2017).

14. Pesiridis, G.S., Lee, V.M. & Trojanowski, J.Q. Mutations in TDP-43 link glycine-rich domain functions to amyotrophic lateral sclerosis. *Hum Mol Genet* **18**, R156-62 (2009).

15. Mackenzie, I.R., Rademakers, R. & Neumann, M. TDP-43 and FUS in amyotrophic lateral sclerosis and frontotemporal dementia. *Lancet Neurol* **9**, 995-1007 (2010).

16. Exome Aggregation Consortium (ExAC), C., MA. . (Accessed February 2016).

17. Gussow, A.B., Petrovski, S., Wang, Q., Allen, A.S. & Goldstein, D.B. The intolerance to functional genetic variation of protein domains predicts the localization of pathogenic mutations within genes. *Genome Biol* **17**, 9 (2016).

18. Marchler-Bauer, A. *et al.* CDD: conserved domains and protein three-dimensional structure. *Nucleic Acids Res* **41**, D348-52 (2013).

19. Traynelis, J. *et al.* Optimizing genomic medicine in epilepsy through a gene-customized approach to missense variant interpretation. *Genome Res* **27**, 1715-1729 (2017).

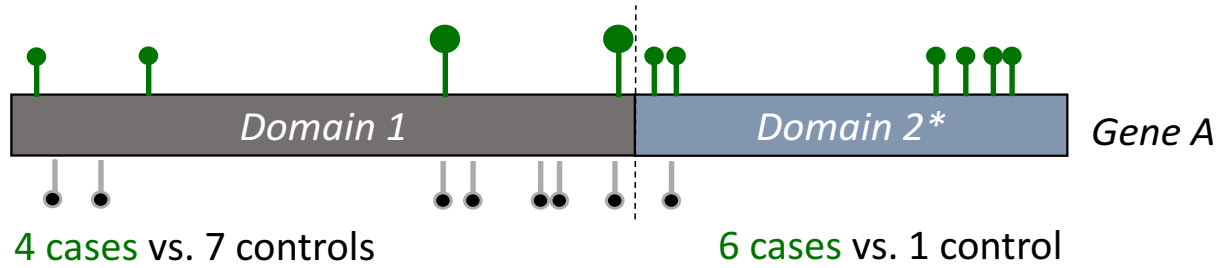20. Rogelj, B., Godin, K.S., Shaw, C.E. & Ule, J. The functions of glycine‐rich regions in TDP‐43, FUS and related Rna‐binding proteins. *RNA Binding Proteins*, 1-17 (2011).

21. Kato, T. *et al.* Galectin-1 Is a Component of Neurofilamentous Lesions in Sporadic and Familial Amyotrophic Lateral Sclerosis. *Biochemical and Biophysical Research Communications* **282**, 166-172 (2001).

22. Kobayakawa, Y. *et al.* Galectin-1 deficiency improves axonal swelling of motor neurones in SOD1(G93A) transgenic mice. *Neuropathology and Applied Neurobiology* **41**, 227-244 (2015).

23. Lerman, B.J. *et al.* Deletion of galectin-3 exacerbates microglial activation and accelerates disease progression and demise in a SOD1(G93A) mouse model of amyotrophic lateral sclerosis. *Brain and Behavior* **2**, 563-575 (2012).

24. Becher, A. *et al.* The armadillo protein p0071 controls KIF3 motor transport. *Journal of Cell Science* **130**, 3374-3387 (2017).

25. Keil, R., Schulz, J. & Hatzfeld, M. p0071/PKP4, a multifunctional protein coordinating cell adhesion with cytoskeletal organization. *Biol Chem* **394**, 1005-17 (2013).

26. Keil, R. & Hatzfeld, M. The armadillo protein p0071 is involved in Rab11-dependent recycling. *Journal of Cell Science* **127**, 60-71 (2014).

27. Yasuda, K. *et al.* The RNA-binding protein Fus directs translation of localized mRNAs in APC-RNP granules. *J Cell Biol* **203**, 737-46 (2013).

28. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-60 (2009).

29. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297-303 (2010).

30. Jun, G. *et al.* Detecting and estimating contamination of human DNA samples in sequencing and array-based genotype data. *Am J Hum Genet* **91**, 839-48 (2012).

31. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867-73 (2010).

32. Price, A.L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet* **38**, 904-9 (2006).

807   33.   Epi, K.c. & Epilepsy Phenome/Genome, P. Ultra-rare genetic
808         variation in common epilepsies: a case-control sequencing study.
809         *Lancet Neurol* **16**, 135-143 (2017).
810   34.   (ESP), N.G.E.S.P. Exome Variant Server.
811   35.   Cingolani, P. *et al.* A program for annotating and predicting the
812         effects of single nucleotide polymorphisms, SnpEff: SNPs in the
813         genome of Drosophila melanogaster strain w1118; iso-2; iso-3. *Fly*
814         *(Austin)* **6**, 80-92 (2012).
815   36.   Finn, R.D. *et al.* The Pfam protein families database: towards a more
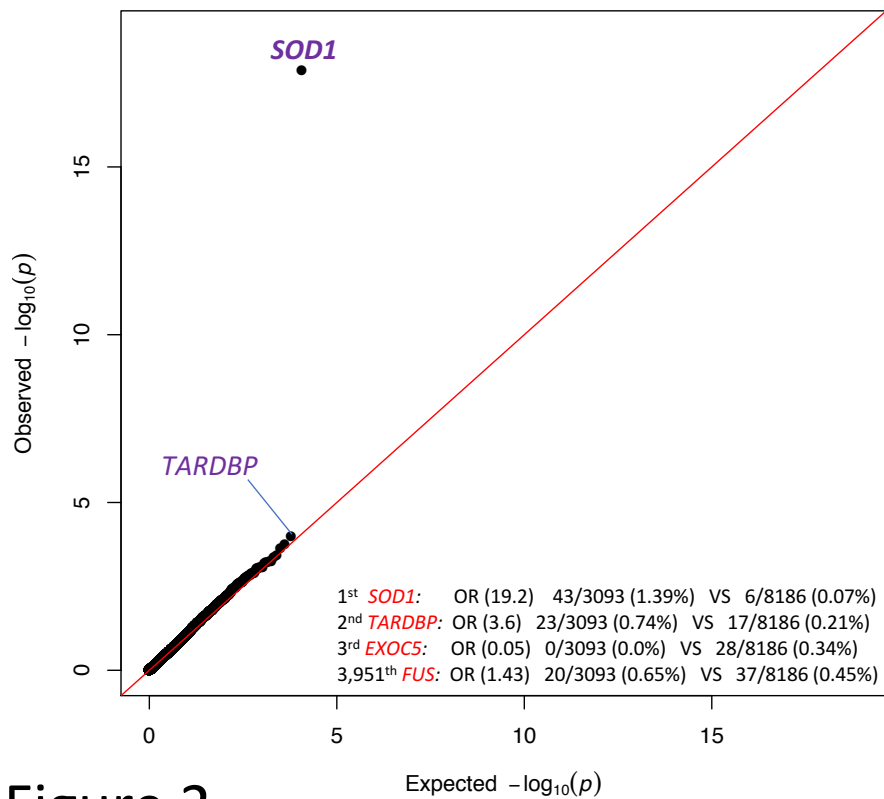816         sustainable future. *Nucleic Acids Res* **44**, D279-85 (2016).
817

**A**

Gene A

*Gene A:* 10 cases vs. 8 controls

**B**

Domain 1     Domain 2*     Gene A

4 cases vs. 7 controls     6 cases vs. 1 control

**C**

Domain 1     Domain 2*     Gene A
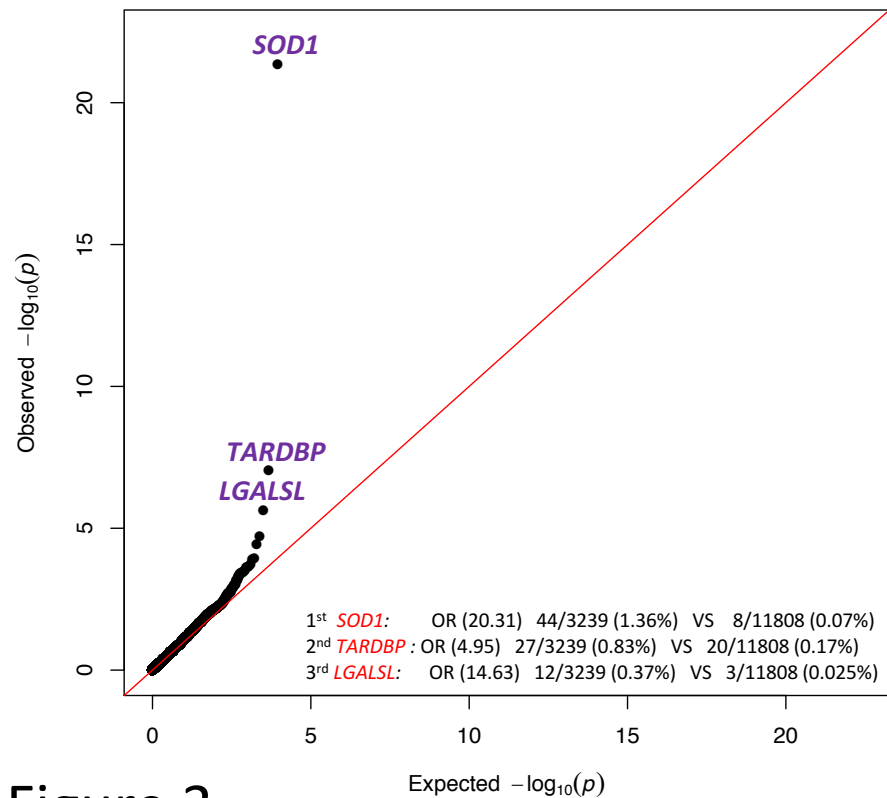
*Intolerant Gene A :* 8 cases vs. 1 controls

Figure 1

*intolerant to variation

Figure 2

Figure 3

Figure 4