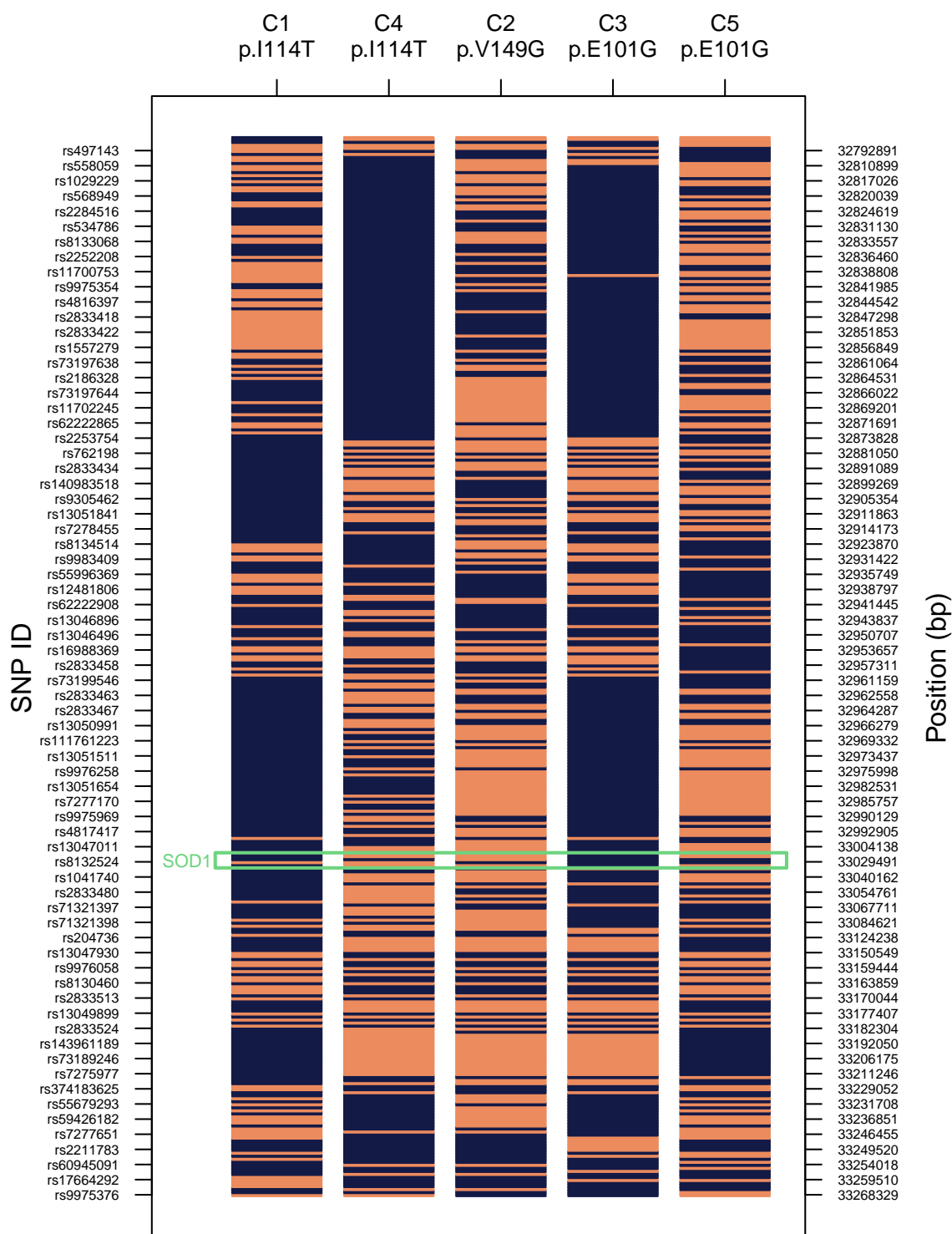


Supplementary Information

Identity by descent analysis identifies founder events and links *SOD1* familial and sporadic ALS cases

Lyndal Henden, Natalie A. Twine, Piotr Szul, Emily P McCann,
Garth A Nicholson, Dominic B Rowe, Matthew C Kiernan,
Denis C Bauer, Ian P Blair and Kelly L. Williams

Supplementary Figures



Supplementary Figure 1. Visualisation of the 350-SNP founder haplotypes for five clusters of Australian MND cases with one of three *SOD1* mutations.

Reference and alternative alleles are coloured blue and orange, respectively. SNP identifiers and their respective genomic position (hg19) are displayed for every 5th SNP beginning at SNP 1 (rs497143) in Supplementary Dataset 1.



Supplementary Figure 2. The process for extracting founder haplotypes over *SOD1*.

First, the founder haplotype is extracted for each cluster in the relatedness network. The endpoints of the founder haplotype for a cluster are taken as the intersection of the IBD segments inferred between all pairs of individuals within this cluster. Since the genotype data has been phased, the founder haplotype is the singular haplotype that appears in all samples within this cluster between the identified haplotype endpoints. Next, the interval that all of the founder haplotypes overlap is defined, and each founder haplotype is reported within this interval. Lastly, SNPs with an identical allele across all founder haplotypes are removed from the reported haplotypes.

Supplementary Data

Supplementary Data 1. Supplementary_data.xlsx

350-SNP founder haplotypes for five clusters of Australian MND cases with one of three *SOD1* mutations.

Supplementary Data 2. haplotype_code.R

R code to extract *SOD1* haplotypes from vcf data.

```
1 # Lyndal Henden
2 # 13/05/2020
3
4 # extract the 5x founder haplotypes over SOD1 for 3x SOD1 mutations
5
6 # load libraries
7 library(SeqArray)
8 library(data.table)
9
10 # load IBD segments and subset by chr 21
11 ibd_tribes <- readRDS("path/to/ibd.match")
12 ibd_tribes <- ibd_tribes[ibd_tribes[,"chr"] == 21,]
13
14 # load phased WGS data for chr 21
15 vcf.fn <- "path/to/chr21_phased.vcf.gz"
16 seqVCF2GDS(vcf.fn, "path/to/chr21_phased.gds")
17 gds_file <- seqOpen("path/to/chr21_phased.gds")
18
19 # create a map file for chr21 from VCF with phased haplotype data
20 map_chr21 <- data.frame(chr=seqGetData(gds_file, "chromosome"), snp_id=
  seqGetData(gds_file, "annotation/id"), pos=seqGetData(gds_file, "
  position"), alleles=seqGetData(gds_file, "allele"),
21 allele_m1=t(seqGetData(gds_file, "genotype")
  [1,,]), allele_m2=t(seqGetData(gds_file, "genotype")[2,,]))
22 map_chr21 <- map_chr21[order(map_chr21[, "pos"]),]
23
24 # reformat phased haplotypes
25 sample_id <- seqGetData(gds_file, "sample.id")
26 allele_m1 <- map_chr21[,grep("allele_m1", colnames(map_chr21))]
27 allele_m2 <- map_chr21[,grep("allele_m2", colnames(map_chr21))]
28 map_chr21 <- map_chr21[,1:3]
29 phased_gt <- data.frame()
30 for (i in 1:ncol(allele_m1)) {
31   map_chr21 <- cbind.data.frame(map_chr21, data.frame(paste(allele_m1[,i
  ], allele_m2[,i], sep="|")))
32 }
33 colnames(map_chr21)[4:ncol(map_chr21)] <- sample_id
34 seqClose(gds_file)
35
36
37 *****
38 ***** get founder haplotype endpoints within a cluster *****
39 *****
40
41 # NOTE: this code needs to run seperately for each cluster of interest
42
43 # manually extract the sample IDs within a cluster
44 samples_cluster <- c("sample1", "sample2", "sample3", "sample4")
45
46 # SOD1 locus - hg19 coordinates
47 gene_interval <- c(21, 33031935, 33041243)
```

```

48
49 # subset IBD segments by chr 21 and samples within the cluster of
      interest
50 ibd_cluster <- ibd_tribes[ibd_tribes[,"iid1"] %in% samples_cluster & ibd
      _tribes[,"iid2"] %in% samples_cluster,]
51 ibd_cluster_chr <- ibd_cluster[ibd_cluster[,"chr"] == gene_interval[1],]
52
53 # function to extract overlap between two intervals
54 getOverlap <- function(region.1, region.2) {
55   a <- max(region.1[1], region.2[1])
56   b <- min(region.1[2], region.2[2])
57   return(c(a,b))
58 }
59
60 # get IBD segments overlapping SOD1 for all samples within the cluster
      of interest
61 ibd_overlap <- rep(0, nrow(ibd_cluster_chr))
62 for(i in 1:nrow(ibd_cluster_chr)){
63   my_overlap <- getOverlap(ibd_cluster_chr[i,c("start.position.bp", "end.
      position.bp")], c(gene_interval[2], gene_interval[3]))
64   if(my_overlap[2] > my_overlap[1])
65     ibd_overlap[i] <- 1
66 }
67 ibd_interval <- ibd_cluster_chr[ibd_overlap == 1,]
68
69 # get genomic region common to all IBD samples over SOD1
70 common_region <- c(21, max(ibd_interval[,"start.position.bp"]), min(ibd_
      interval[,"end.position.bp"]))
71
72 # subset haplotype data for samples of interest over common region
73 haplotype_c1 <- map_chr21[map_chr21[,"pos"] >= common_region[2] & map_
      chr21[,"pos"] <= common_region[3],c("chr", "snp_id", "pos", "alleles"
      , samples_cluster)]
74
75 # for each subsequent cluster, generate a dataframe called haplotype_c2,
      haplotype_c3, ..., haplotype_cn, where n is the number of clusters
76
77
78 *****
79 ***** extract shared haplotype over common interval *****
80 *****
81
82 # NOTE: this code needs to run seperately for each cluster of interest
83
84 # trim founder haplotypes to overlapping interval between all clusters -
      have done manually here by looking at haplotype_c1, haplotype_c2,
      etc
85 interval_overlap <- c(32792145, 33274026)
86
87 # function that extracts the most common haplotype shared between
      samples within a cluster over an interval
88 get_haplotypes <- function(haplotypes_interval) {
89   haplotypes_test <- haplotypes_interval[,1:4]
90   haplotypes_test_v2 <- data.frame()
91   for (i in 5:ncol(haplotypes_interval)) {
92     hap_i <- do.call('rbind', strsplit(as.character(haplotypes_interval[,
      i]), "|"))
93     hap_i_df <- data.frame(hap_i[,1], hap_i[,3])
94     colnames(hap_i_df) <- rep(colnames(haplotypes_interval)[i], 2)
95     haplotypes_test <- cbind.data.frame(haplotypes_test, hap_i_df)

```

```

96   haplotypes_test_v2 <- rbind.data.frame(haplotypes_test_v2, data.
    frame(sample_id=colnames(haplotypes_interval)[i],hap=c(paste(hap_i_df
      [,1], collapse=""), paste(hap_i_df[,2], collapse="))))
97 }
98 tab_hap <- table(haplotypes_test_v2[,2])
99
100 # what is the frequency of most common haplotype
101 cat(paste0("Haplotype count = ", max(tab_hap),"\n"))
102
103 # samples with common haplotype
104 cat("\nSamples with haplotype:\n")
105 print(as.character(haplotypes_test_v2[haplotypes_test_v2[,2] == names(
    tab_hap[tab_hap == max(tab_hap)]],1)))
106
107 # samples without common haplotype - genotype/phasing errors, ends
    might differ slightly
108 cat("\nSamples without haplotype:\n")
109 print(as.character(unique(haplotypes_test_v2[,1])[!(unique(haplotypes_
    test_v2[,1]) %in% as.character(haplotypes_test_v2[haplotypes_test_v2
      [,2] == names(tab_hap[tab_hap == max(tab_hap)]],1)))]))
110
111 # extract haplotype over internal in data frame format
112 col_same <- NULL
113 for (i in 5:ncol(haplotypes_test)) {
114   if (paste0(haplotypes_test[,i],collapse="") == names(tab_hap[tab_hap
      == max(tab_hap)]))
115     col_same <- c(col_same, i)
116 }
117
118 return(list(haplotypes_test[,c(1:4,col_same)], haplotypes_test,
    haplotypes_test_v2, tab_hap))
119 }
120
121 # extract shared haplotype for a cluster
122 hap_cluster_1 <- haplotype_c1[haplotype_c1[,"pos"] >= interval_overlap
    [1] & haplotype_c1[,"pos"] <= interval_overlap[2],]
123 haplotypes_cluster_1 <- get_haplotypes(hap_cluster_1)
124 haplotype_C1 <- haplotypes_cluster_1[[1]][,c(1:5)]
125
126 # for each subsequent cluster, generate haplotypes_cluster_2, haplotype_
    C2 etc.
127
128
129 *****
130 ***** filter and format haplotypes *****
131 *****
132
133 # combine the shared haplotypes for each cluster into a single dataframe
134 haplotypes_all <- cbind.data.frame(haplotype_C1, haplotype_C4[,6],
    haplotype_C2[,6], haplotype_C3[,6], haplotype_C5[,6])
135 for (i in 6:10) haplotypes_all[,i] <- as.numeric(as.character(haplotypes
    _all[,i]))
136
137 # remove SNPs that are the same in all haplotypes
138 keep_snp <- NULL
139 for (i in 1:nrow(haplotypes_all)) {
140   if (!all(haplotypes_all[i,6] == haplotypes_all[i,6:10]))
141     keep_snp <- c(keep_snp, i)
142 }
143 haplotypes_pruned <- haplotypes_all[keep_snp,]

```

```

144
145 # manually add 3x SOD1 variants
146 sod1_mut <- data.frame(chr=c(21,21,21), snp_id=c("rs121912439",
147           rs121912441", "rs1476760624"), pos=c(33039633,33039672,33040872),
148           alleles=c("A,G","T,C","T,G"), C1=c(0,1,0), C4=c(
149           (0,1,0), C2=c(1,0,0), C3=c(0,0,1), C5=c(0,0,1))
150 for (i in 1:3) {
151   haplotypes_pruned <- rbind.data.frame(haplotypes_pruned[haplotypes_
152     pruned[, "pos"] < sod1_mut[i, "pos"], ],
153     sod1_mut[i, ],
154     haplotypes_pruned[haplotypes_
155     pruned[, "pos"] > sod1_mut[i, "pos"], ])
156 }
157 dim(haplotypes_pruned)
158
159 # add strand and mutation status to haplotype dataframe
160 ref_allele <- do.call('rbind', strsplit(as.character(haplotypes_pruned[, "
161   alleles"]), ", "))[,1]
162 alt_allele <- do.call('rbind', strsplit(as.character(haplotypes_pruned[, "
163   alleles"]), ", "))[,2]
164 sod1_mutation <- rep(".", nrow(haplotypes_pruned))
165 sod1_mutation[haplotypes_pruned[, "pos"] == 33039633] <- "p.E101G"
166 sod1_mutation[haplotypes_pruned[, "pos"] == 33039672] <- "p.I114T"
167 sod1_mutation[haplotypes_pruned[, "pos"] == 33040872] <- "p.V149G"
168 strand <- rep("+", nrow(haplotypes_pruned))
169 haplotypes_save <- data.frame(haplotypes_pruned[,1:3], strand, ref_
170   allele, alt_allele, sod1_mutation, haplotypes_pruned[,c(5:9)])
171 for (i in 4:7) haplotypes_save[,i] <- as.character(haplotypes_save[,i])
172
173 # change haplotypes to AGCT annotation rather than 0/1
174 for (i in 8:12) {
175   risk_hap <- NULL
176   for (j in 1:nrow(haplotypes_save)) {
177     risk_hap <- c(risk_hap, ifelse(haplotypes_save[j,i] == 0, haplotypes
178       _save[j, "ref_allele"], haplotypes_save[j, "alt_allele"]))
179   }
180   haplotypes_save[,i] <- risk_hap
181 }
182 colnames(haplotypes_save) <- c("Chromosome", "SNP ID", "Position (hg19)"
183   , "Strand", "Ref", "Alt", "SOD1 Mutation",
184     "C1 - p.I114T", "C4 - p.I114T", "C2 - p.
185     V149G", "C3 - p.E101G", "C5 - p.E101G")
186
187 #*****
188 #***** add 1000K EUR + gnomAD allele frequencies *****
189 #*****
190
191 #*** 1000K ***
192
193 # read in phased 1000 genomes chr 21 VCF file - downloaded from: ftp://
194   ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/
195 vcf.fn <- "path/to/ALL.chr21.phase3_shapeit2_mvncall_integrated_v5a
196   .20130502.genotype.vcf.gz"
197 seqVCF2GDS(vcf.fn, "path/to/ALL.chr21.phase3_shapeit2_mvncall_integrated
198   _v5a.20130502.genotype.gds")
199 gds_file <- seqOpen("path/to/ALL.chr21.phase3_shapeit2_mvncall_
200   integrated_v5a.20130502.genotype.gds")
201 seqSetFilter(gds_file)
202
203

```

```

190 # subset by variants in haplotypes
191 snp_id <- seqGetData(gds_file, "annotation/id")
192 variant_id <- seqGetData(gds_file, "variant.id")
193 seqSetFilter(gds_file, variant.id=variant_id[snp_id %in% haplotypes_save
194   [, "SNP ID"]])
195
196 # subset allele freq
197 ceu_freq <- seqGetData(gds_file, "annotation/info/EUR_AF")
198
199 # get alleles
200 alleles_1k <- seqGetData(gds_file, "allele")
201 ref_1k <- do.call('rbind', strsplit(alleles_1k, ","))[,1]
202 alt_1k <- do.call('rbind', strsplit(alleles_1k, ","))[,2]
203
204 # get alleles
205 map_1k <- data.frame(chr=seqGetData(gds_file, "chromosome"), snp_id=
206   seqGetData(gds_file, "annotation/id"), pos=seqGetData(gds_file, "
207   position"),
208   REF=ref_1k, ALT=alt_1k, ceu_freq[[2]])
209 seqClose(gds_file)
210
211 *** gnomad ***
212
213 # read in allele frequencies manually collated for haplotype SNPs from
214 # webpage
215 gnomad <- fread("path/to/haplotype_coords_gnomAD_anno.csv", data.table=
216   FALSE)
217
218 # merge with 1000 genomes map file - 3 SOD1 mutations missing from 1k
219 # data, 2 missing from gnomad
220 gnomad_1k <- merge(map_1k, gnomad, by.x=c("snp_id", "REF", "ALT"), by.y=c(
221   "snp", "Reference", "Alternate"))
222 gnomad_1k <- gnomad_1k[order(gnomad_1k[, "pos.x"]),]
223
224 # look at difference in allele freq between 1k and gnomad
225 hist(abs(gnomad_1k[, "ceu_freq..2.."] - gnomad_1k[, "AF"]), xlim=c(0,1),
226   xlab="diff", main="")
227 # looks great!
228
229 # create new haplotype file with updated allele frequencies
230 haplotypes_save$freq_gnomad <- 0
231 haplotypes_save$freq_ceu <- 0
232 for (i in 1:nrow(haplotypes_save)) {
233   if (haplotypes_save[i, "SNP ID"] %in% gnomad_1k[, "snp_id"]) {
234     haplotypes_save[i, "freq_ceu"] <- gnomad_1k[gnomad_1k[, "snp_id"] %in%
235     haplotypes_save[i, "SNP ID"], "ceu_freq..2.."]
236     haplotypes_save[i, "freq_gnomad"] <- gnomad_1k[gnomad_1k[, "snp_id"] %
237     in% haplotypes_save[i, "SNP ID"], "AF"]
238   }
239 }
240 haplotypes_save_v2 <- haplotypes_save[, c(1:3, 13, 14, 4:12)]
241 colnames(haplotypes_save_v2)[4:5] <- c("Reference Allele Freq - 1000
242   Genomes EUR", "Reference Allele Freq - gnomAD non-neuro non-Finnish
243   EUR")
244
245 # change SOD1 allele freq
246 haplotypes_save_v2[haplotypes_save_v2[, "SNP ID"] %in% c("rs121912439", "
247   rs121912441", "rs1476760624"), "Reference Allele Freq - 1000 Genomes
248   EUR"] <- NA
249 haplotypes_save_v2[haplotypes_save_v2[, "SNP ID"] %in% c("rs121912439", "

```



```
    rs121912441", "rs1476760624"), "Reference Allele Freq - gnomAD non-  
    neuro non-Finnish EUR"] <- NA  
236 haplotypes_save_v2[haplotypes_save_v2[, "SNP ID"] %in% "rs121912441", "  
    Reference Allele Freq - gnomAD non-neuro non-Finnish EUR"] <-  
    0.0000111664  
237  
238 # look at haplotypes  
239 haplotypes_save_v2
```