1  **Genome-wide association study identifies genetic factors that modify age at onset in**

2  **Machado-Joseph disease**

3  Fulya Akçimen[1,2], Sandra Martins[3,4], Calwing Liao[1,2], Cynthia V. Bourassa[2,5], Hélène Catoire[2,5],

4  Garth A. Nicholson[6], Olaf Riess[7], Mafalda Raposo[8], Marcondes C. França Jr.[9], João

5  Vasconcelos[10], Manuela Lima[8], Iscia Lopes-Cendes[11,12], Maria Luiza Saraiva-Pereira[13,14], Laura

6  B. Jardim[13,15], Jorge Sequeiros[4,16,17], Patrick A. Dion[2,5], Guy A. Rouleau[1,2,5*]

7  [1]Department of Human Genetics, McGill University, Montréal, QC, Canada;

8  [2]Montreal Neurological Institute and Hospital, McGill University, Montréal, QC, Canada;

9  [3]i3S – Instituto de Investigação e Inovação em Saúde, Universidade do Porto, Portugal;

10  [4]IPATIMUP – Institute of Molecular Pathology and Immunology of the University of Porto, Portugal;

11  [5]Department of Neurology and Neurosurgery, McGill University, Montréal, QC, Canada;

12  [6]University of Sydney, Department of Medicine, Concord Hospital, Australia;

13  [7]Institute of Medical Genetics and Applied Genomics, University of Tuebingen, Tuebingen, Germany;

14  [8]Faculdade de Ciências e Tecnologia, Universidade dos Açores e Instituto de Biologia Molecular e Celular

15  (IBMC), Instituto de Investigação e Inovação em Saúde (i3S), Universidade do Porto, Portugal;

16  [9]Department of Neurology, Faculty of Medical Sciences, UNICAMP, Campinas, SP, Brazil;

17  [10]School of Medical Sciences, Department of Medical Genetics and Genomic Medicine, University of

18  Campinas (UNICAMP), Campinas, SP, Brazil;

19  [11]The Brazilian Institute of Neuroscience and Neurotechnology (BRAINN), Campinas, SP, Brazil;

20  [12]Departamento de Neurologia, Hospital do Divino Espírito Santo, Ponta Delgada, Portugal;

21  [13]Medical Genetics Service, Hospital de Clínicas de Porto Alegre (HCPA), Porto Alegre, Brazil;

22  [14]Depto. de Bioquímica – ICBS, Universidade Federal do Rio Grande do Sul (UFRGS);

23  [15]Depto de Medicina Interna, Universidade Federal do Rio Grande do Sul (UFRGS), Porto Alegre, Brazil;

24  [16]Institute for Molecular and Cell Biology (IBMC), Universidade do Porto, Porto, Portugal;

25  [17]Instituto de Ciências Biomédicas Abel Salazar (ICBAS), Universidade do Porto, Portugal.

26  **\*Corresponding author:** Guy A. Rouleau

27  Address: Montreal Neurological Institute and Hospital,

28  3801 University Street, Room 636,

29  Montréal, Québec, Canada H3A 2B4

30  Email: guy.rouleau@mcgill.ca

31    Phone: +1 (514) 398-6644

32

33

**Abstract**

Machado-Joseph disease (MJD/SCA3) is the most common form of dominantly inherited ataxia worldwide. The disorder is caused by an expanded CAG repeat in the *ATXN3* gene. Past studies have revealed that the length of the expansion partly explains the disease age at onset (AO) variability of MJD, which is confirmed in this study. Using a total of 786 MJD patients from five different geographical origins, a genome-wide association study (GWAS) was conducted to identify additional AO modifying factors that could explain some of the residual AO variability. We identified nine suggestively associated loci ($P < 1 \times 10^{-5}$). These loci were enriched for genes involved in vesicle transport, olfactory signaling, and synaptic pathways. Furthermore, associations between AO and the *TRIM29* and *RAG* genes suggests that DNA repair mechanisms might be implicated in MJD pathogenesis. Our study demonstrates the existence of several additional genetic factors, along with CAG expansion, that may lead to a better understanding of the genotype-phenotype correlation in MJD.

**Keywords**

## Introduction

Machado-Joseph disease, also known as spinocerebellar ataxia type 3 (MJD/SCA3), is an autosomal dominant neurodegenerative disorder that is characterized by progressive cerebellar ataxia and pyramidal signs, which can be associated with a complex clinical picture and includes extrapyramidal signs or amyotrophy [1, 2]. MJD is caused by an abnormal CAG trinucleotide repeat expansion in exon 10 of the ataxin-3 gene (*ATXN3*), located at 14q32.1. Deleterious expansions consensually contain 61 to 87 CAG repeats, whereas wild type alleles range from 12 to 44 [2].

As with other diseases caused by repeat expansions, such as Huntington's disease (HD) and other spinocerebellar ataxias, there is an inverse correlation between expanded repeat size and the age at which pathogenesis leads to disease onset [3]. Depending on the cohort structure, the size of the repeat expansion explains 55 to 70%  of the age at onset (AO) variability in MJD, suggesting the existence of additional modifying factors [3,4]. Although several genetic factors have been proposed as modifiers, such as CAG repeat size of normal *ATXN3* (SCA3), *HTT* (HD), *ATXN2* (SCA2) and *ATN1* (DRPLA) alleles, *APOE* status, and expression level of *HSP40* [4,5,6], these were not replicated by subsequent studies [7, 8]. Since CAG tract profile and allelic frequencies of the potential modifier loci can have unique characteristics in different populations, large collaborative studies are required to identify genetic modifiers in MJD, as well as replicate the findings of such studies [8].

Previously, Genetic Modifiers of Huntington's Disease (GeM-HD) Consortium carried out a GWA approach of HD individuals to reveal genetic modifiers of AO in HD [9,10]. A total of eleven [9] and fourteen loci [10] were found to be associated with residual age at HD onset. In the present

71 study, we performed the first GWAS to identify some possible genetic modifiers of AO in MJD.

72 First, we assessed the relationship between AO and size of the expanded ($CAG_{exp}$) and normal

73 ($CAG_{nor}$) alleles, biological sex and geographical origin. Next, we determined a residual AO for

74 each subject, which is the difference between the measured AO and the predicted/estimated AO

75 from expanded CAG repeat size alone. Using the residuals as a quantitative phenotype for a

76 GWAS, we looked for genetic factors that modulate AO in MJD.

**Methods**

**Study subjects**

79 A total of 786 MJD patients from five distinct geographical origins (Portugal, Brazil, North

80 America, Germany and Australia) were included in the present study. The overall average age at

81 onset (standard deviation) was 38 ($\pm$ 1.82) years, with a 1:1 male to female ratio. All subjects

82 provided informed consent, and the study was approved by the respective institutional review

83 boards. Detailed cohort demographics are shown in Supplementary Table 1.

**Assessment of the *ATXN3* CAG repeat length**

85 A singleplex polymerase chain reaction was performed to determine the length of the $CAG_{exp}$ and

86 $CAG_{nor}$ alleles at exon 10 of *ATXN3* [11]. The final volume for each assay was 10 µL: 7.5 ng of

87 gDNA, 0.2 µM of each primer, 5 µL of Taq PCR Master Mix Kit Qiagen®, 1 µL of Q-Solution

88 from Qiagen® and $H_2O$. Fragment length analysis was done using ABIPrism 3730xl sequencer

89 (Applied Biosystems®, McGill University and Genome Québec Innovation Centre) and

90 GeneMapper software [12]. A stepwise regression model was performed to assess the correlation

91 between AO and $CAG_{exp}$ size, as well as gender, origin, $CAG_{nor}$ size, and interaction between these

92  variables. Residual AO was calculated for each subject by subtracting individual's expected AO

93  based upon $CAG_{exp}$ size from actual AO, to be used as the primary phenotype for following genetic

94  approach.

**Genotyping, quality control and imputation**

96  Samples were genotyped using the Global Screening Array v.1.0 from Illumina (636,139 markers).

97  Sample-based (missingness, relatedness, sex, and multidimensional scaling analysis) and SNP-

98  based quality assessments (missingness, Hardy-Weinberg equilibrium, and minor allele

99  frequency) were conducted using PLINK version 1.9 [13]. In sample level QC, samples were

100  excluded with one or more of the following: high missingness (missingness rate > 0.05), close

101  relationship (pi-hat value > 0.2), discrepancy between genetically-inferred sex and reported sex,

102  population outliers (deviation $\geq$ 4 SD from the population mean in multidimensional scaling

103  analysis). All SNPs were checked for marker genotyping call rate (> 98%), minor allele frequency

104  (MAF) > 0.05, and HWE (p-value threshold = $1.0 \times 10^{-5}$).

105  Phasing and imputation were performed using SHAPEIT [14] and PBWT [15] pipelines,

106  implemented on the Sanger Imputation Service [16]. Haplotype Reference Consortium (HRC)

107  reference panel r1.1 containing 64,940 human haplotypes at 40,405,505 genetic markers were used

108  as the reference panel. Imputed variants with an allele count of 30 (MAF > 0.02), an imputation

109  quality score above 0.3 and an HWE p-value of $> 1.0 \times 10^{-5}$ were included for subsequent analysis.

**Genome-wide association analysis**

111  A genome-wide linear mixed model based association analysis was conducted using GCTA

112  version 1.91.7 [17]. Residual AO was modelled as a function of minor allele count of the test SNP,

113   sex, and the first three principal components based on the scree plot (Supplementary Figure 1).

114   The --mlma-loco option, which takes into account the difference in allele frequency between

115   populations, was used to control for population structure. QQ plots and Manhattan plots were

116   generated in FUMA v.1.3.4 [18]. Regional association plots were generated using LocusZoom

117   [19] (Supplementary Figure 3).

118   **Functional annotation of SNPs**

119   Genomic risk loci were defined using SNP2GENE function implemented in FUMA. Independent

120   suggestive SNPs ($P < 1 \times 10^{-5}$) with a threshold of $r^2 < 0.6$ were selected within a 250 kb window.

121   The UK Biobank release 2 European population consisting of randomly selected 10,000 subjects

122   was used as the reference population panel. The ANNOVAR [20] categories and combined

123   annotation-dependent depletion (CADD) [21] scores were obtained from FUMA for functional

124   annotation. Functionally annotated variants were mapped to genes based on genomic position

125   using FUMA positional mapping tool.

126   **Pathway analysis**

127   To identify known biological pathways and gene sets at the associated loci, an enrichment

128   approach   was   applied   using   public   datasets   containing   Gene   Ontology   (GO,

129   http://geneontology.org),   the   Kyoto   Encyclopaedia   of   Genes   and   Genomes   (KEGG,

130   https://www.genome.jp/kegg)   and   Reactome   (https://reactome.org)   pathways.   The   primary

131   enrichment analysis was performed using the i-GSEA4GWAS v2 [22]. It uses a candidate list of

132   a genome-wide set of genes mapped within the SNP loci and ranks them based on the strength of

133   their association with the phenotype. Genes were mapped within 20 kb up or downstream of the

134   SNPs with a $P < 0.05$. Gene and pathway sets meeting a false discovery rate *(FDR)*-corrected *q-*

135    value $< 0.05$ were regarded as significantly associated with high confidence, and $q$-value $< 0.25$

136    was regarded to be possibly associated with the phenotype of interest. We performed a secondary

137    gene-based association test using the Versatile Gene-based Association Study (VEGAS) [23]

138    algorithm that controls the number of SNPs in each gene and the linkage disequilibrium (LD)

139    between these SNPs using the HapMap European population. As a third algorithm to identify

140    enriched pathways, we used Pathway Scoring Algorithm (PASCAL) [24], which controls for

141    potential bias from gene size, SNP density, as well as LD. ClueGO [25] and CluePedia [26] plug-

142    ins in Cytoscape were employed to visualize identified pathways and their clustering.

143    **Results**

144    **The inverse correlation between CAG$_{exp}$ and age at onset**

145    In the first phase of the study, the expanded *ATXN3*-CAG repeat lengths of 786 MJD patients were

146    assessed. The mean (SD) CAG$_{exp}$ size were Australia: 68.2 ($\pm$3.3), Brazil: 74.3 (3.9), Germany:

147    72.9 ($\pm$3.6), North America: 73 ($\pm$4.3) and Portugal: 72 ($\pm$4.0). Next, the relationship between AO

148    and CAG$_{exp}$ size, CAG$_{nor}$ size, sex and ethnicity was examined (Supplementary Table 1). The

149    previously observed negative correlation between *ATXN3* CAG$_{exp}$ size and AO [3] was confirmed

150    (Pearson's correlation coefficient $R^2 = 0.62$) (Figure 1). The CAG$_{nor}$ size (P $= 0.39$), sex (P $= 0.02$)

151    and geographic origin (P [Brazil] $= 0.38$, P [Germany] $= 0.38$, P [North America] $= 0.33$, P

152    [Portugal] $= 0.29$) were not significant and their addition had little contribution to the model ($\Delta R^2$

153    $= 0.0072$). Residual AO for each sample was calculated and used as a quantitative phenotype to

154    identify the modifiers of AO. The distribution of residual AO was close a theoretical normal

155    distribution (Figure 1).

156    **Genome-wide association study**

157 After post-imputation quality assessments, a total of 700 individuals with genotyping information

158 for 6,716,580 variants remained for GWAS. The resulting Manhattan plots and quantile-quantile

159 (QQ) plots are shown in Figure 2. The genomic inflation factor was close to one ($\lambda = 0.98$),

160 indicating the p-values were not inflated. No association signal was identified meeting genome-

161 wide significance ($P < 5 \times 10^{-8}$, the genome-wide Bonferroni-corrected significance threshold);

162 however, genome-wide suggestive associations ($P < 1 \times 10^{-5}$) with 204 variants across 9 loci were

163 identified (Supplementary Table 3). The most significantly associated SNP at each locus are shown

164 in Table 1. Positional gene mapping aligned SNPs to 17 genes by their genomic location. Fourteen

165 of the 204 variants had a Combined Annotation Dependent Depletion (CADD)-PHRED score

166 higher than the suggested threshold for deleterious SNPs (12.37), arguing the given loci have a

167 functional role [27].

168 **Interaction analysis between $CAG_{exp}$ and SNP genotype**

169 To assess a possible interaction between $CAG_{exp}$ size and the variants identified, each of the nine

170 variants was added to the initial linear regression, modelling AO as a function of $CAG_{exp}$ size,

171 SNP, sex, the first three principal components, $CAG_{nor}$ size, and interaction of SNP:$CAG_{exp}$.

172 Association of each independent SNP with AO revealed nominally significant p-values. Among

173 the nine variants, only rs585809 (mapped to *TRIM29*) had a significant interaction with $CAG_{exp}$ (P

174 $= 0.01$), suggesting that rs585809 might modulate AO through this epistatic interaction on $CAG_{exp}$.

175 **Association of HD-AO modifier variants in MJD**

176 Association of previously identified HD-AO modifier loci in MJD were assessed. Among the 25

177 HD-AO modifier variants in 17 loci, a total of 18 variants (MAF > 0.02) in 12 loci were tested in

178 this study (Supplementary Table 4). None of these HD-AO modifiers reached the genome-wide

179    suggestive threshold. However, two variants rs144287831 ($P = 0.02$, effect size $= -0.98$) and

180    rs1799977 ($P = 0.02$, effect size $= -0.98$) in the *MLH1* locus were found to be nominally associated

181    with a later AO in MJD.

182    **Pathway and gene-set enrichment analysis**

183    A gene-set enrichment and pathway analysis was conducted using i-GSEA4GWAS. Various

184    approaches and algorithms are currently in use to conduct similar analyses. To be able to make

185    better comparisons with other studies that may use different approaches, we performed a secondary

186    gene-set enrichment and pathway analysis using the VEGAS2 and PASCAL software

187    (Supplementary Tables 5-7). We also used these results for replication purposes in our own study.

188    A total of 13 overrepresented pathways were found, after FDR-multiple testing correction (q-value

189    $< 0.05$) in the primary GSEA analysis and replicated using at least one of the secondary gene-set

190    enrichment algorithms (Table 2). Overall, the most significantly enriched gene-sets and pathways

191    were vesicle transport, olfactory signaling, and synaptic pathways. Visualization and clustering of

192    pathways are shown in Figure 3.

193

194    **Discussion**

195    Using five cohorts from different geographical origins, we performed the first GWAS to examine

196    the presence of genetic factors that could modify AO in MJD. We identified a total of nine loci

197    that were potentially associated with either an earlier or later AO. Concomitantly, we confirmed

198    the previously observed negative correlation between $CAG_{exp}$ and AO [3]. It was shown previously

199    that normal *ATXN3* allele ($CAG_{nor}$) had a significant influence on AO of MJD [28]; however,

200    several studies did not replicate this effect [6,8]. Indeed, we did not observe an association between

201    $CAG_{nor}$ and AO. However, it had little contribution to our model, with a minor difference in the

202    correlation coefficient ($\Delta R^2 = 0.0012$).

203    In our GWAS, the strongest signal is for the rs11529293 variant ($P = 3.30 \times 10^{-6}$) within the

204    *C11orf72* and *RAG* loci at 11p12. Within this locus, two *RAG* genes, recombination-activating

205    genes *RAG1* and *RAG2*, were shown to be implicated in DNA damage response and DNA repair

206    machineries [29,30]. The rs585809 variant, which was mapped to the *TRIM29* gene, was found to

207    interact with $CAG_{exp}$, suggesting that it might have an effect on AO through this interaction. Both

208    *RAG* and *TRIM29* loci were identified as AO-hastening modifiers. *TRIM29* encodes for tripartite

209    motif protein 29, which is implicated in mismatch repair and double strand breaks pathways

210    [31,32]. TRIM29 is involved both upstream and downstream of these pathways, in the regulation

211    of DNA repair proteins into chromatin by mediating the interaction between them. One of these

212    DNA repair proteins is MLH1, which is implicated in mismatch repair complex [32]. Previously,

213    the *MLH1* locus was identified as an AO modifier in another neurodegenerative disease caused by

214    CAG repeat expansion, Huntington's disease [9,10,33]. Additionally, in a genome-wide genetic

215    screening study, MLH1-knock out was shown to modify the somatic expansion of the CAG repeat

216    and slow the pathogenic process in HD mouse model [34]. Overall, the association of *TRIM29* and

217    *RAG* loci suggests that DNA repair mechanisms may be implicated in the alteration of AO of MJD,

218    as well as HD, and may have a role in the pathogenesis of other CAG repeat diseases. Interestingly,

219    in a previous study, we found variants in three transcription-coupled repair genes (ERCC6, RPA,

220    and CDK7) associated with different CAG instability patterns in MJD [35].

221    We identified gene-sets enriched in olfactory signaling, vesicle transport, and synaptic pathways.

222    Olfactory dysfunction is one of the main non-motor symptoms that was already described in

223    patients with MJD [36,37]. In a previous study, transplantation of olfactory ensheathing cells,

224    which are specialized glial cells of the primary olfactory system, were found to improve motor

225    function in an MJD mice model, and were suggested as a novel potential strategy for MJD

226    treatment [38]. Vesicle transport and synaptic pathways were also implicated in MJD, as well as

227    in other neurodegenerative diseases [39,40]. An interruption of synaptic transmission caused by

228    an expanded polyglutamine repeat and mutant ataxin-3 aggregates were shown in *Drosophila* and

229    *Caenorhabditis elegans* models of MJD. Therefore, the interaction between synaptic vesicles and

230    mutant aggregates supports the role of synaptic vesicle transport in the pathogenesis of MJD

231    [41,42]. Overall, we suggest that these gene-sets and pathways might construct a larger molecular

232    network that could modulate the AO in MJD.

233    In summary, our study identified nine genetic loci that may modify the AO of MJD. Identification

234    of *TRIM29* and *RAG* genetic variants, as well as our gene-set enrichment analyses, implicated

235    DNA repair, olfactory signaling, synaptic, and vesicle transport pathways in the pathogenesis of

236    MJD. Although we used different cohorts from five distinct geographical ethnicities, a replication

237    study in similar or additional populations would add valuable evidence to support our findings.

238

239     **Description of Supplemental Data**

240     Supplemental Data include three figures and seven tables

241     **Declaration of Interests**

242     The authors declare no conflict of interest.

243     **Acknowledgements**

255

256

257 **Table 1.** Suggestive loci associated with residual age at onset in MJD. Chr: chromosome, MAF: minor allele frequency, 1KGP: 1000
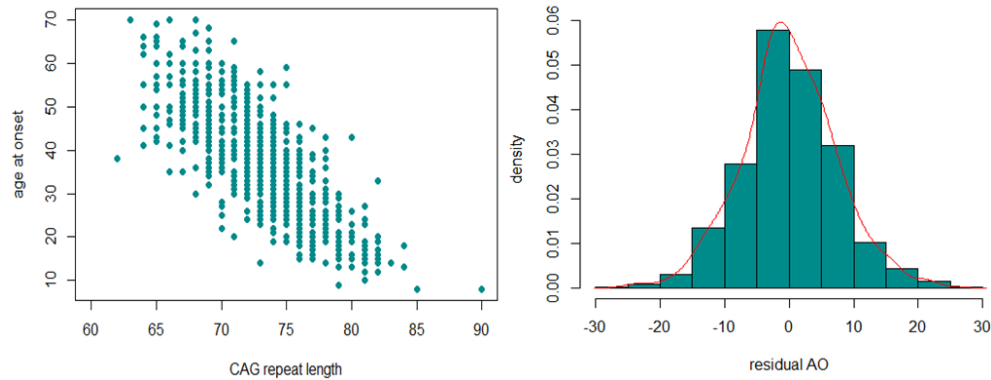258 Genomes Project

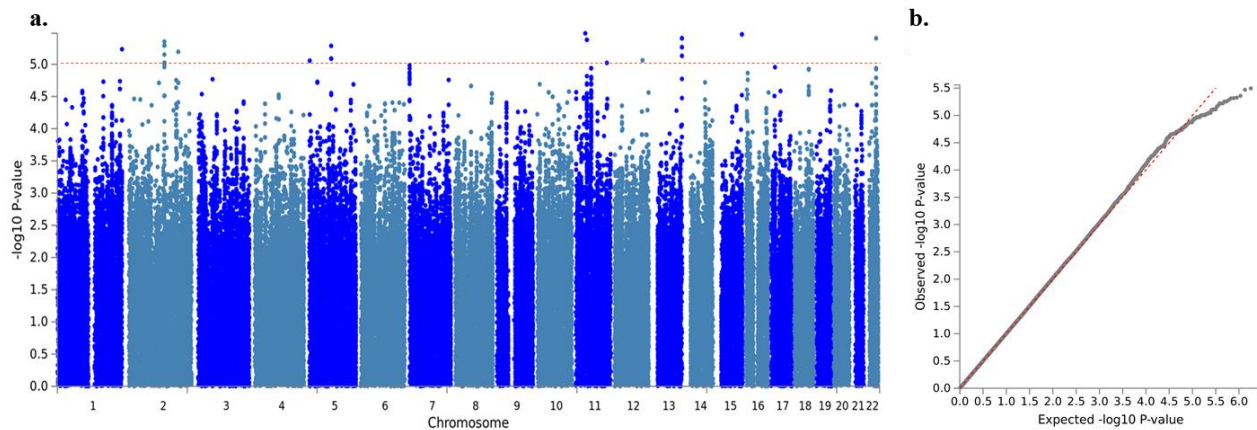| SNP | Chr | Position (GRCh37) | Nearest gene | Minor allele | Major allele | MJD MAF | 1KGP MAF | b (SNP effect) | P-value |
|---|---|---|---|---|---|---|---|---|---|
| rs62171220 | 2 | 137802855 | *THSD7B* | G | C | 0.13 | 0.11 | 2.71 | $4.45 \times 10^{-6}$ |
| rs2067390 | 2 | 191209028 | *HIBCH, INPP1* | A | T | 0.04 | 0.06 | 4.74 | $6.39 \times 10^{-6}$ |
| rs144891322 | 5 | 85135387 | *RPL5P17,* | C | T | 0.02 | 0.007 | 6.10 | $5.18 \times 10^{-6}$ |
| rs11529293 | 11 | 36855388 | *C11orf74,  RAG1, RAG2* | T | C | 0.14 | 0.26 | -2.71 | $3.30 \times 10^{-6}$ |
| rs7480166 | 11 | 42984753 | *HNRNPKP3* | A | G | 0.40 | 0.40 | -1.86 | $4.17 \times 10^{-6}$ |
| rs585809 | 11 | 119949979 | *TRIM29* | T | C | 0.06 | 0.17 | -3.76 | $9.50 \times 10^{-6}$ |
| rs72660056 | 13 | 113507543 | *ATP11A* | A | G | 0.08 | 0.05 | -3.29 | $3.94 \times 10^{-6}$ |
| rs11857349 | 15 | 99924857 | *TTC23,  SYNM, LRRC28* | G | A | 0.04 | 0.02 | -4.58 | $3.43 \times 10^{-6}$ |
| rs8141510 | 22 | 42821185 | *NFAM1,  CYP2D6, NAGA, NDUFA6* | C | T | 0.43 | 0.49 | 1.83 | $3.94 \times 10^{-6}$ |

259

260

**Table 2.** Pathways significant after multiple-correction ($q < 5 \times 10^{-2}$) in the primary GSEA analysis and replicated using at least one of the secondary gene-set enrichment algorithms. NA means that the pathway was not enriched by at least two significant genes in VEGAS.

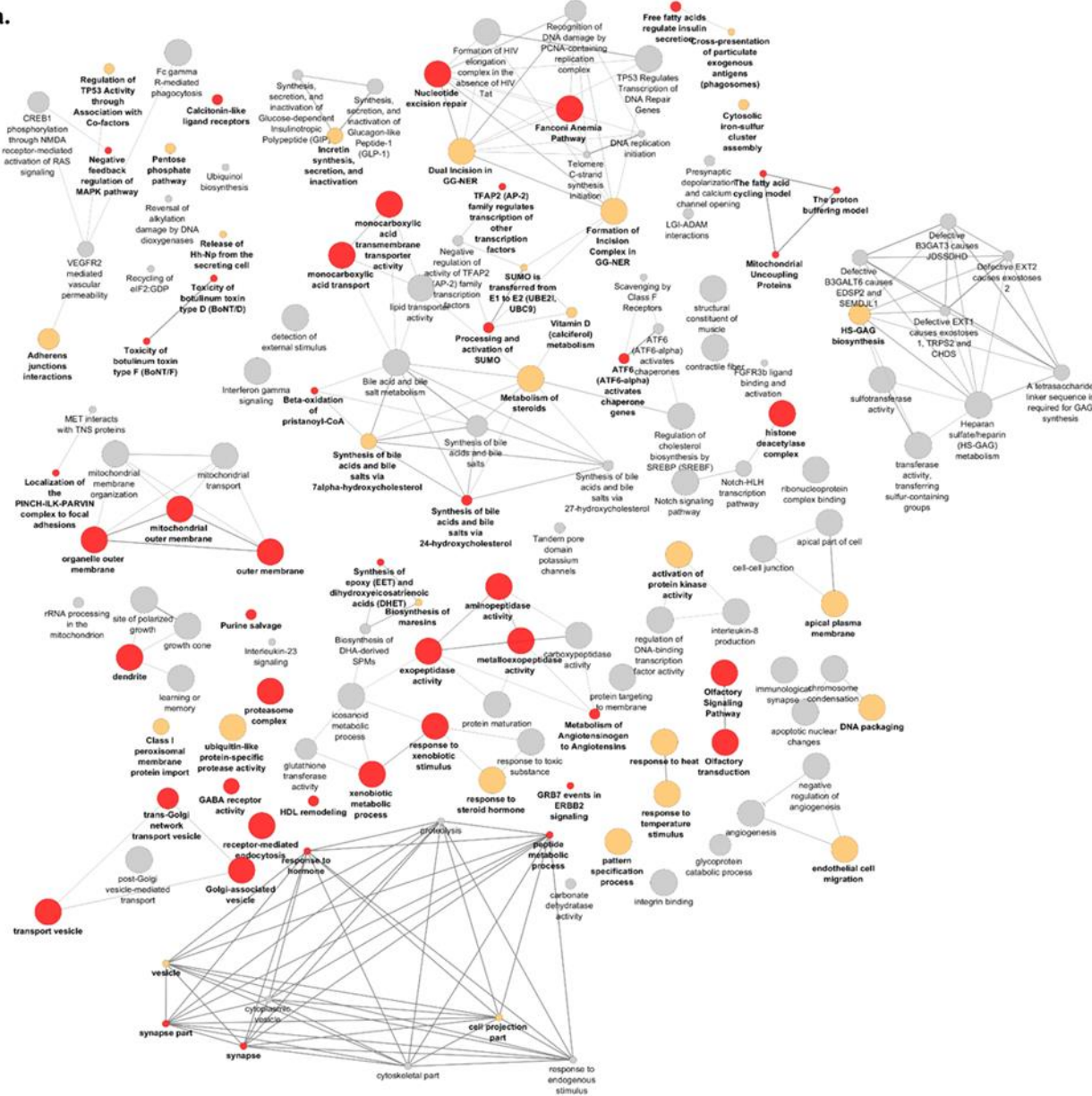| Pathway | Description | p-value (GSEA) | q-value (GSEA) | p-value (VEGAS) | permuted p-value (VEGAS) | p-value (PASCAL) |
|---|---|---|---|---|---|---|
| GO:0030133 | transport vesicle | $< 1.0 \times 10^{-3}$ | $8.20 \times 10^{-3}$ | $6.15 \times 10^{-40}$ | $4.46 \times 10^{-1}$ | $6.70 \times 10^{-3}$ |
| KEGG:04740 | olfactory transduction | $< 1.0 \times 10^{-3}$ | $8.30 \times 10^{-3}$ | NA | NA | $3.89 \times 10^{-4}$ |
| R-HSA:381753 | olfactory signaling pathway | $< 1.0 \times 10^{-3}$ | $8.80 \times 10^{-3}$ | $1.10 \times 10^{-27}$ | $7.71 \times 10^{-1}$ | $2.51 \times 10^{-4}$ |
| GO:0044456 | synapse part | $< 1.0 \times 10^{-3}$ | $9.30 \times 10^{-3}$ | $1.25 \times 10^{-182}$ | $< 1.0 \times 10^{-6}$ | $< 1.0 \times 10^{-7}$ |
| R-HSA:74217 | purine salvage | $< 1.0 \times 10^{-3}$ | $1.06 \times 10^{-2}$ | $1.06 \times 10^{-2}$ | $2.15 \times 10^{-1}$ | $6.48 \times 10^{-3}$ |
| GO:0045202 | synapse | $< 1.0 \times 10^{-3}$ | $1.15 \times 10^{-2}$ | $1.15 \times 10^{-2}$ | $< 1.0 \times 10^{-6}$ | $< 1.0 \times 10^{-7}$ |
| GO:0004177 | aminopeptidase activity | $< 1.0 \times 10^{-3}$ | $1.50 \times 10^{-2}$ | $1.50 \times 10^{-2}$ | $3.41 \times 10^{-1}$ | $1.24 \times 10^{-2}$ |
| GO:0008238 | exopeptidase activity | $< 1.0 \times 10^{-3}$ | $1.80 \times 10^{-2}$ | $1.80 \times 10^{-2}$ | $2.80 \times 10^{-2}$ | $8.31 \times 10^{-3}$ |
| GO:0006898 | receptor mediated endocytosis | $< 1.0 \times 10^{-3}$ | $2.25 \times 10^{-2}$ | $2.25 \times 10^{-2}$ | $2.03 \times 10^{-1}$ | $6.64 \times 10^{-3}$ |
| GO:0016917 | GABA receptor activity | $< 1.0 \times 10^{-3}$ | $2.26 \times 10^{-2}$ | $2.26 \times 10^{-2}$ | $1.30 \times 10^{-4}$ | $2.30 \times 10^{-5}$ |
| GO:0030140 | trans Golgi network transport vesicle | $< 1.0 \times 10^{-3}$ | $2.36 \times 10^{-2}$ | $2.36 \times 10^{-2}$ | $2.80 \times 10^{-2}$ | $1.28 \times 10^{-1}$ |
| GO:0009725 | response to hormone stimulus | $< 1.0 \times 10^{-3}$ | $2.73 \times 10^{-2}$ | $2.73 \times 10^{-2}$ | $1.32 \times 10^{-1}$ | $1.30 \times 10^{-4}$ |
| GO:0030425 | Dendrite | $< 1.0 \times 10^{-3}$ | $3.86 \times 10^{-2}$ | $3.86 \times 10^{-2}$ | $< 1.0 \times 10^{-6}$ | $< 1.0 \times 10^{-7}$ |

264

265 **Figure 1.** The inverse correlation between $CAG_{exp}$ and AO (left) and the distribution of residual AO (right) observed in our MJD cohort.
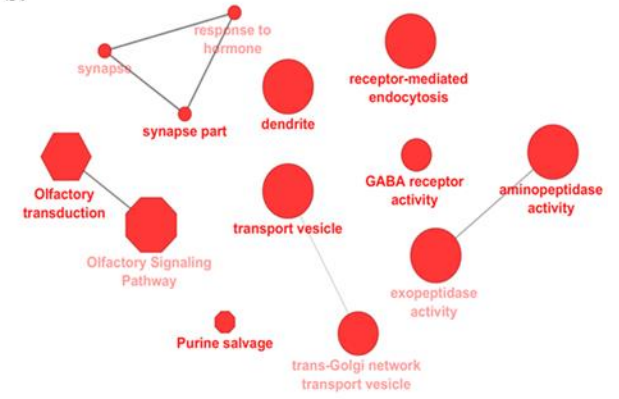


266

267 **Figure 2.** Manhattan plot (a) and QQ plot (b) of the GWAS for residual AO of MJD. Imputed using the HRC panel, 6,716,580 variants

268 that passed QC are included in the plot. The x-axis shows the physical position along the genome. The y-axis shows the $-\log_{10}$(p-value)

269 for association. The red line indicates the level of genome-wide suggestive association ($P = 1 \times 10^{-5}$).

271 **Figure 3.** Visualization of the gene-sets and pathways enriched in primary GSEA analysis (a) and replicated in VEGAS and PASCAL

272 (b). The size of the nodes corresponds to the number of the genes associated with a term. The significance is represented by the color of

273 the nodes ($P < 0.05$, $0.05 < P < 0.1$ and $P > 0.1$ are represented by red, yellow and gray, respectively).

274 **References**

275 1. Twist, E.C. *et al.* Machado Joseph disease maps to the same region of chromosome 14 as the
276 spinocerebellar ataxia type 3 locus. *Journal of Medical Genetics* **32**, 25-31 (1995).

277 2. Bettencourt, C. & Lima, M. Machado-Joseph Disease: from first descriptions to new
278 perspectives. *Orphanet J Rare Dis* **6**, 35 (2011).

279 3. Maciel, P. *et al.* Correlation between CAG repeat length and clinical features in Machado-Joseph
280 disease. *Am J Hum Genet* **57**, 54-61 (1995).

281 4. de Mattos, E.P., Kolbe Musskopf, M., Bielefeldt Leotti, V., Saraiva-Pereira, M.L. & Jardim, L.B.
282 Genetic risk factors for modulation of age at onset in Machado-Joseph disease/spinocerebellar
283 ataxia type 3: a systematic review and meta-analysis. *Journal of Neurology, Neurosurgery &amp;*
284 *Psychiatry* **90**, 203-210 (2019).

285 5. Zijlstra, M.P. *et al.* Levels of DNAJB family members (HSP40) correlate with disease onset in
286 patients with spinocerebellar ataxia type 3. *European Journal of Neuroscience* **32**, 760-770
287 (2010).

288 6. Tezenas du Montcel, S. *et al.* Modulation of the age at onset in spinocerebellar ataxia by CAG
289 tracts in various genes. *Brain* **137**, 2444-2455 (2014).

290 7. Chen, Z. *et al.* (CAG)n loci as genetic modifiers of age-at-onset in patients with Machado-Joseph
291 disease from mainland China. *Brain* **139**, e41-e41 (2016).

292 8. Raposo, M., Ramos, A., Bettencourt, C. & Lima, M. Replicating studies of genetic modifiers in
293 spinocerebellar ataxia type 3: can homogeneous cohorts aid? *Brain* **138**, e398-e398 (2015).

294 9. Genetic Modifiers of Huntington's Disease, C. Identification of Genetic Factors that Modify
295 Clinical Onset of Huntington's Disease. *Cell* **162**, 516-26 (2015).

296 10. Lee, J.-M. *et al.* CAG Repeat Not Polyglutamine Length Determines Timing of Huntington's
297 Disease Onset. *Cell* **178**, 887-900.e14 (2019).

298 11. Martins, S., Calafell, F., Wong, V.C., Sequeiros, J. & Amorim, A. A multistep mutation mechanism
299 drives the evolution of the CAG repeat at MJD/SCA3 locus. *Eur J Hum Genet* **14**, 932-40 (2006).

300 12. Chatterji, S. & Pachter, L. Reference based annotation with GeneMapper. *Genome biology* **7**,
301 R29-R29 (2006).

302 13. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets.
303 *Gigascience* **4**, 7 (2015).

304 14. Delaneau, O., Marchini, J. & Zagury, J.-F. A linear complexity phasing method for thousands of
305 genomes. *Nature Methods* **9**, 179 (2011).

306 15. Durbin, R. Efficient haplotype matching and storage using the positional Burrows-Wheeler
307 transform (PBWT). *Bioinformatics (Oxford, England)* **30**, 1266-1272 (2014).

308 16. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nature*
309 *genetics* **48**, 1279-1283 (2016).

310 17. Yang, J., Lee, S.H., Goddard, M.E. & Visscher, P.M. GCTA: a tool for genome-wide complex trait
311 analysis. *American journal of human genetics* **88**, 76-82 (2011).

312 18. Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping and
313 annotation of genetic associations with FUMA. *Nature Communications* **8**, 1826 (2017).

314 19. Pruim, R.J. *et al.* LocusZoom: regional visualization of genome-wide association scan results.
315 *Bioinformatics (Oxford, England)* **26**, 2336-2337 (2010).

316 20. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from
317 high-throughput sequencing data. *Nucleic Acids Res* **38**, e164 (2010).

318 21. Rentzsch, P., Witten, D., Cooper, G.M., Shendure, J. & Kircher, M. CADD: predicting the
319 deleteriousness of variants throughout the human genome. *Nucleic Acids Research* **47**, D886-
320 D894 (2018).

321  22.  Zhang, K., Cui, S., Chang, S., Zhang, L. & Wang, J. i-GSEA4GWAS: a web server for identification
322      of pathways/gene sets associated with traits by applying an improved gene set enrichment
323      analysis to genome-wide association study. *Nucleic Acids Research* **38**, W90-W95 (2010).
324  23.  Mishra, A. & Macgregor, S. VEGAS2: Software for More Flexible Gene-Based Testing. *Twin*
325      *Research and Human Genetics* **18**, 86-91 (2014).
326  24.  Lamparter, D., Marbach, D., Rueedi, R., Kutalik, Z. & Bergmann, S. Fast and Rigorous
327      Computation of Gene and Pathway Scores from SNP-Based Summary Statistics. *PLoS Comput*
328      *Biol* **12**, e1004714 (2016).
329  25.  Bindea, G. *et al.* ClueGO: a Cytoscape plug-in to decipher functionally grouped gene ontology
330      and pathway annotation networks. *Bioinformatics (Oxford, England)* **25**, 1091-1093 (2009).
331  26.  Bindea, G., Galon, J. & Mlecnik, B. CluePedia Cytoscape plugin: pathway insights using integrated
332      experimental and in silico data. *Bioinformatics (Oxford, England)* **29**, 661-663 (2013).
333  27.  Amendola, L.M. *et al.* Actionable exomic incidental findings in 6503 participants: challenges of
334      variant classification. *Genome research* **25**, 305-315 (2015).
335  28.  Franca, M.C., Jr. *et al.* Normal ATXN3 Allele but Not CHIP Polymorphisms Modulates Age at
336      Onset in Machado-Joseph Disease. *Front Neurol* **3**, 164 (2012).
337  29.  Lescale, C. & Deriano, L. The RAG recombinase: Beyond breaking. *Mechanisms of Ageing and*
338      *Development* **165**, 3-9 (2017).
339  30.  Bahjat, M. & Guikema, J.E.J. The Complex Interplay between DNA Injury and Repair in
340      Enzymatically Induced Mutagenesis and DNA Damage in B Lymphocytes. *Int J Mol Sci* **18**(2017).
341  31.  Wikiniyadhanee, R., Lerksuthirat, T., Stitchantrakul, W., Chitphuk, S. & Dejsuphong, D. AB064.
342      TRIM29: a novel gene involved in DNA repair mechanisms. *Annals of Translational Medicine* **5**,
343      AB064-AB064 (2017).
344  32.  Masuda, Y. *et al.* TRIM29 regulates the assembly of DNA repair proteins into damaged
345      chromatin. *Nat Commun* **6**, 7299 (2015).
346  33.  Lee, J.-M. *et al.* A modifier of Huntington's disease onset at the MLH1 locus. *Human Molecular*
347      *Genetics* **26**, 3859-3867 (2017).
348  34.  Pinto, R.M. *et al.* Mismatch repair genes Mlh1 and Mlh3 modify CAG instability in Huntington's
349      disease mice: genome-wide and candidate approaches. *PLoS Genet* **9**, e1003930 (2013).
350  35.  Martins, S. *et al.* Modifiers of (CAG)(n) instability in Machado-Joseph disease (MJD/SCA3)
351      transmissions: an association study with DNA replication, repair and recombination genes. *Hum*
352      *Genet* **133**, 1311-8 (2014).
353  36.  Braga-Neto, P. *et al.* Clinical correlates of olfactory dysfunction in spinocerebellar ataxia type 3.
354      *Parkinsonism Relat Disord* **17**, 353-6 (2011).
355  37.  Pedroso, J.L. *et al.* Nonmotor and extracerebellar features in Machado-Joseph disease: a review.
356      *Mov Disord* **28**, 1200-8 (2013).
357  38.  Hsieh, J. *et al.* Human Olfactory Ensheathing Cell Transplantation Improves Motor Function in a
358      Mouse Model of Type 3 Spinocerebellar Ataxia. *Cell Transplant* **26**, 1611-1621 (2017).
359  39.  Wiatr, K. *et al.* Altered Levels of Proteins and Phosphoproteins, in the Absence of Early Causative
360      Transcriptional Changes, Shape the Molecular Pathogenesis in the Brain of Young
361      Presymptomatic Ki91 SCA3/MJD Mouse. *Mol Neurobiol* (2019).
362  40.  Gissen, P. & Maher, E.R. Cargos and genes: insights into vesicular transport from inherited
363      human disease. *J Med Genet* **44**, 545-55 (2007).
364  41.  Gunawardena, S. & Goldstein, L.S.B. Polyglutamine Diseases and Transport Problems: Deadly
365      Traffic Jams on Neuronal Highways. *JAMA Neurology* **62**, 46-51 (2005).
366  42.  Khan, L.A. *et al.* Expanded polyglutamines impair synaptic transmission and ubiquitin-
367      proteasome system in Caenorhabditis elegans. *J Neurochem* **98**, 576-87 (2006).

368