

1 **Biallelic expansion of an intronic repeat in *RFC1* is a common cause of late-**  
2 **onset ataxia**

3  
4 Andrea Cortese<sup>1\*</sup>, Roberto Simone<sup>2</sup>, Roisin Sullivan<sup>1§</sup>, Jana Vandrovцова<sup>1§</sup>, Huma Tariq<sup>1</sup>,  
5 Yau Way Yan<sup>1</sup>, Jack Humphrey<sup>1</sup>, Zane Jaunmuktane<sup>2</sup>, Prasanth Sivakumar<sup>1</sup>, James Polke<sup>3</sup>,  
6 Muhammad Ilyas<sup>4</sup>, Eloise Tribollet<sup>1</sup>, Pedro J. Tomaselli<sup>5</sup>, Grazia Devigili<sup>6</sup>, Iliaria Callegari<sup>7</sup>,  
7 Maurizio Versino<sup>7,8</sup>, Vincenzo Salpietro<sup>1</sup>, Stephanie Efthymiou<sup>1</sup>, Diego Kaski<sup>1</sup>, Nick W.  
8 Wood<sup>1</sup>, Nadja S. Andrade<sup>9</sup>, Elena Buglo<sup>10</sup>, Adriana Rebelo<sup>10</sup>, Alexander M. Rossor<sup>1</sup>,  
9 Adolfo Bronstein<sup>2</sup>, Pietro Fratta<sup>1</sup>, Wilson J. Marques<sup>5</sup>, Stephan Züchner<sup>10</sup>, Mary M.  
10 Reilly<sup>1#</sup>, and Henry Houlden<sup>1\*,#</sup>

11  
12 <sup>1</sup>Department of Neuromuscular Disease, UCL Institute of Neurology and The National Hospital for  
13 Neurology, London, UK.

14 <sup>2</sup>Department of Clinical and Movement Neurosciences, UCL Institute of Neurology and The National  
15 Hospital for Neurology, London, UK.

16 <sup>3</sup>Neurogenetics Laboratory, UCL Institute of Neurology and The National Hospital for Neurology,  
17 London, UK.

18 <sup>4</sup>Department of Biotechnology, Islamabad University, Islamabad and Punjab University, Lahore,  
19 Pakistan.

20 <sup>5</sup>Department of Neurology, School of Medicine of Ribeirão Preto, University of São Paulo,  
21 Ribeirão Preto, Brazil.

22 <sup>6</sup>UO Neurologia I, Fondazione IRCCS Istituto Neurologico “Carlo Besta”, Milano, Italy.

23 <sup>7</sup>IRCCS Mondino Foundation, Pavia, Italy.

24 <sup>8</sup>Department of Brain and Behavioral Sciences, University of Pavia, Pavia, Italy.

25 <sup>9</sup>Department of Psychiatry and Behavioural Sciences, Center for Therapeutic Innovation,  
26 University of Miami Miller School of Medicine, Miami, Florida, USA.

27 <sup>10</sup>Dr. John T. Macdonald Foundation Department of Human Genetics and John P. Hussman  
28 Institute for Human Genomics, University of Miami Miller School of Medicine, Miami, Florida,  
29 USA.

30 <sup>§</sup>These authors contributed equally to this work.

31 <sup>#</sup>These authors jointly directed this project.

32 \*E-mail: [andrea.cortese@ucl.ac.uk](mailto:andrea.cortese@ucl.ac.uk), [h.houlden@ucl.ac.uk](mailto:h.houlden@ucl.ac.uk)

33

34 Late-onset ataxia is common, often idiopathic, and can result from cerebellar,  
35 proprioceptive or vestibular impairment, when in combination also termed cerebellar  
36 ataxia, neuropathy, vestibular areflexia syndrome (CANVAS). We used non-parametric  
37 linkage analysis and genome sequencing to identify a biallelic intronic AAGGG  
38 repeat expansion in *RFC1* as the cause of familial CANVAS and a frequent cause of  
39 late-onset ataxia, particularly if sensory neuronopathy and bilateral vestibular  
40 areflexia coexisted. The expansion, which occurs in the polyA tail of an AluSx3  
41 element and differs in terms of both size and nucleotide sequence from the reference  
42 (AAAAG)<sub>11</sub> allele, does not affect *RFC1* expression in patient peripheral and brain  
43 tissue, suggesting no overt loss-of-function. These data, along with an expansion  
44 carrier frequency of 0.7% in Europeans, implies that biallelic AAGGG expansion in  
45 *RFC1* is a frequent cause of late-onset ataxia.

46 Late-onset ataxia, postural imbalance and falls are a frequent reason for neurological  
47 consultation. Physiologically, motor coordination is achieved under visual control thanks  
48 to the cerebellar integration of proprioceptive information conveyed by large-fibre  
49 sensory neurons and vestibular inputs. Failure of any or a combination of these systems  
50 can result in ataxia<sup>1-6</sup>. Both acquired and genetic causes are known, but a large proportion  
51 remains idiopathic.

52 Previous studies suggest that there is a spectrum of clinical signs, from pure  
53 idiopathic late-onset cerebellar degeneration (ILOCA) through to the combined  
54 degeneration of the cerebellum and its vestibular and sensory afferents, which has been  
55 named cerebellar atrophy, neuropathy and vestibular areflexia syndrome (CANVAS)  
56 (**Fig. 1a**)<sup>7</sup>. CANVAS is an adult-onset slowly progressive neurological disorder  
57 characterized by imbalance, sensory neuropathy (neuronopathy), bilateral  
58 vestibulopathy<sup>8</sup>, chronic cough, and occasionally autonomic dysfunction<sup>9</sup>. Typically,  
59 sensory action potentials and somatosensory potentials are absent throughout, brain MRI  
60 shows cerebellar atrophy, and vestibular testing is consistent with impaired vestibular  
61 function bilaterally<sup>9-17</sup>. Late-onset ataxia and CANVAS are usually sporadic, but  
62 occasionally occur in siblings, raising the possibility of recessive transmission. However,  
63 initial attempts to identify the underlying genetic defect by whole-exome sequencing  
64 were unsuccessful.

65 Using non-parametric linkage analysis and whole-genome sequencing, we  
66 identified a recessive intronic AAGGG repeat expansion in the replication factor C  
67 subunit 1 gene (*RFC1*) as a cause of familial CANVAS. The expansion occurs in the polyA  
68 tail of an AluSx3 element and differs in terms of both size and nucleotide sequence from  
69 the reference (AAAAG)<sub>11</sub> allele. Screening of additional sporadic cases with late-onset  
70 ataxia confirmed the presence of the mutated AAGGG repeat expansion in 22% of them,  
71 and in higher percentages if sensory neuronopathy and/or bilateral vestibular areflexia  
72 coexisted, suggesting that it represents a frequent and underrecognized cause of late-  
73 onset ataxia.

74

## 75 RESULTS

76 **Genetic study.** We genotyped 29 individuals (23 affected and 6 unaffected) from 11  
77 families (**Fig. 1b**). The majority of the families consisted of affected sibships, except for  
78 two first-degree cousins from non-consanguineous families (Fam 5b-2 and Fam 6b-1).  
79 None of the families had convincing evidence of vertical disease transmission.

80 Assuming a recessive mode of inheritance, non-parametric linkage analysis  
81 identified a single peak at position 4q14 with cumulative maximum HLOD of 5.8 (**Fig.**  
82 **2a**). Haplotype analysis defined a 1.7-Mb region between markers rs6814637 and  
83 rs10008483 (chr4:38977921-40712231) where, within single families, affected siblings  
84 shared the same maternal and paternal alleles as opposed to unaffected brothers and  
85 sisters, who had at most one of them (**Fig. 2b**). The region contains 21 known HGNC  
86 genes (**Supplementary Table 1**). Homozygosity mapping in consanguineous family Fam  
87 7 showed that the previously identified 1.7-Mb region is encompassed in a larger run of

88 homozygosity of 12 Mb shared by the affected siblings (**Supplementary Fig. 1**). Of  
89 interest, inside the 1.7-Mb region, four SNPs (rs2066790, rs11096992, rs17584703 and  
90 rs6844176, bold highlighted) mapping inside a region encompassing all exons of *RFC1*  
91 and the last exon of WD repeat domain 19 (*WDR19*) were shared by all affected  
92 individuals from different families except for individual Fam 5b-2, raising the possibility  
93 of a founder haplotype (**Fig. 2c,d**).

94 Whole-exome sequencing was previously performed in seven individuals (Fam 1-  
95 1, Fam 1-2, Fam 1-3, Fam 3-1, Fam 3-2, Fam 4-2, Fam 4-3) from three unrelated families  
96 (Fam1, Fam3, Fam4), but did not identify recurrent non-synonymous variants within the  
97 coding regions of the genes encompassed in the 1.7-Mb region (data not shown). We next  
98 performed whole-genome sequencing (WGS) in an additional six affected individuals  
99 (Fam 2-1, Fam 8-2, Fam 8-3, Fam 6a-2, Fam 5a-2, Fam 7-1), one unaffected subject (Fam 8-  
100 1) from four unrelated families, and one sporadic case (s9). Analysis for non-synonymous  
101 variants and copy number variants did not reveal changes recurring in the affected  
102 families. By visually inspecting the aligned paired reads inside the 1.7-Mb region, we  
103 noted in all CANVAS patients a reduced read depth in a region encompassing a simple  
104 tandem (AAAAG)<sub>11</sub> repeat at position chr4:39350045-39350103 (**Fig. 3a**). Inside the  
105 microsatellite region, the reference (AAAAG)<sub>11</sub> repeat was replaced in patients by a  
106 variable number of AAGGG repeated units, which were detected on the reads mapped to  
107 either side of the short tandem repeat. However, none of the reads could span across the  
108 microsatellite region from one side to the other, suggesting the presence of a biallelic  
109 expansion of the AAGGG repeat unit (**Fig. 3b**). WGS from an unaffected sibling (Fam 8-1)  
110 showed an equal distribution of interrupted reads containing the mutated AAGGG  
111 repeated unit change as well as reads containing the AAAAG repeat.

112 We then performed repeat-primed PCR (RPPCR) with primers targeting the  
113 mutant AAGGG pentanucleotide unit and confirmed the presence of an AAGGG repeat  
114 expansion in all affected members from 11 families, as well as in unaffected carriers (**Fig.**  
115 **3c**). Flanking PCR using standard conditions failed to amplify the region in all patients,  
116 suggesting the presence of a large expansion on both alleles, as opposed to their  
117 unaffected siblings for whom at least one allele could be amplified by PCR (data not  
118 shown).

119 We next screened a cohort of 150 patients diagnosed with sporadic late onset ataxia  
120 and identified an additional 33 (22%) sporadic cases carrying the recessive AAGGG  
121 repeat expansion, as defined by a positive RPPCR for AAGGG repeat unit and the  
122 absence of PCR amplifiable products by standard flanking PCR. The percentage of  
123 positive cases raised to 63% (32/51) if considering cases with late-onset cerebellar ataxia  
124 and sensory neuropathy and to 92% (11/12) in cases with full CANVAS syndrome.  
125 Taking advantage of two informative SNPs rs11096992 and rs2066790, by PCR and direct  
126 sequencing we observed that all additional sporadic cases but individual s23 shared the  
127 same haplotype as familial CANVAS cases.

128 By long-range PCR, we were able to amplify and confirm by Sanger sequencing in  
129 all patients the presence of the AAGGG expansion (**Fig. 3d**). However long-range PCR

130 did not allow sizing of the repeat expansion as PCR is error-prone and contraction of  
 131 repeated regions during PCR cycling have been previously demonstrated<sup>18</sup>. Therefore,  
 132 Southern blots were conducted in 34 cases and confirmed the presence of biallelic large  
 133 expansions in all of them. Biallelic expansions could be visualized as two distinct bands  
 134 in subjects carrying expansions of different sizes, or one thick band if the expanded  
 135 alleles had a similar size (**Supplementary Fig. 2**). Four unaffected siblings from four  
 136 families were also included, and they all carried one expanded and one normal allele.  
 137 Although the expansion size could vary across different families, ranging from around  
 138 400 to 2,000 repeats, in the majority of cases approximately 1,000 repeats were observed.  
 139 Repeat size was relatively stable in siblings within single families. There was no  
 140 association between age at onset and the number of AAGGG repeat units on either the  
 141 smaller or larger allele ( $n = 34$ ;  $r = -0.006$ ,  $P = 0.97$  and  $r = -0.04$ ,  $P = 0.81$ , respectively).  
 142

143 **Polymorphic conformations and allelic distribution of the short tandem repeat locus in**  
 144 **the normal population.** Recessive AAGGG expansion, as defined by the combination of  
 145 positive RPPCR targeting the AAGGG repeat and the absence of a PCR amplifiable  
 146 product on flanking PCR, were not observed in 304 healthy controls screened. RPPCR  
 147 analysis targeting the AAGGG repeat showed that 0.7% (4 out of 608 chromosomes  
 148 tested) carried an AAGGG expansion in heterozygous state. Southern blot analysis was  
 149 performed in all of them and confirmed the presence of an expanded allele in all of them.  
 150 The chr4:39350045-39350103 locus, where the expansion resides, was shown to be highly  
 151 polymorphic in the normal population and, besides the rare AAGGG expansion allele  
 152 (AAGGG)<sub>exp</sub>, three other conformations were observed: (AAAAG)<sub>11</sub>, (AAAAG)<sub>exp</sub>, and  
 153 (AAAGG)<sub>exp</sub> (**Fig. 4a**). The (AAAGG)<sub>exp</sub> often showed interruptions and nucleotide  
 154 changes of the expanded sequence. By a combinatory approach of flanking PCR, RPPCR  
 155 targeting one of the three possible nucleotide sequences, as well as Southern blot and  
 156 Sanger sequencing in selected cases, we observed an allelic distribution of 75.5% ( $n = 459$ )  
 157 for the (AAAAG)<sub>11</sub> allele, 13.0% ( $n = 79$ ) for the (AAAAG)<sub>exp</sub> allele, 7.9% ( $n = 48$ ) for the  
 158 (AAAGG)<sub>exp</sub> allele, and, as per above, 0.7% ( $n = 4$ ) for the (AAGGG)<sub>exp</sub> allele (**Fig. 4b**).  
 159 Average size of (AAAAG)<sub>exp</sub> ranged from 15 to 200 repeats (mean  $72 \pm 43$ ), and  
 160 (AAAGG)<sub>exp</sub> ranged from 40 to 1,000 (mean  $173 \pm 232$ ) (**Fig. 4c**).

161 Eight healthy subjects had biallelic repeat expansions of a distinct repeated unit:  
 162 (AAAAG)<sub>exp</sub>/(AAGGG)<sub>exp</sub> in one case, (AAAGG)<sub>exp</sub>/(AAGGG)<sub>exp</sub> in one case, and  
 163 (AAAAG)<sub>exp</sub>/(AAAGG)<sub>exp</sub> in six cases. 22 cases likely had two expansions of the repeated  
 164 AAAAG unit and nine of the repeated AAGGG unit, as defined by a positive RPPCR for  
 165 the target repeat and two distinct bands on the Southern blot, although we cannot  
 166 exclude that one of the two alleles may be characterized by a distinct nucleotide  
 167 sequence, which was not considered in the present study. Indeed, 9 additional subjects  
 168 had no PCR amplifiable product on flanking PCR and were negative for RPPCR targeting  
 169 the AAAAG, AAAGG, or AAGGG repeated units, suggesting the potential existence of  
 170 other possible allelic conformations in 3% ( $n = 18$ ) of tested chromosomes. Southern blot  
 171 could not be performed because of insufficient amount of DNA in these cases.

172 The haplotype associated in most patients with the AAGGG repeat expansion has  
 173 an allelic carrier frequency in the 1000 Genomes Project control population of 18%. Based  
 174 on rs11096992 and rs2066790 genotyping, the disease-associated haplotype rs2066790  
 175 (AA), rs11096992 (AA) was absent in recessive state from healthy individuals who carried  
 176 two (AAAAG)<sub>11</sub> alleles, two (AAAAG)<sub>exp</sub> alleles or a compound (AAAAG)<sub>11</sub>/(AAAAG)<sub>exp</sub>  
 177 genotype, but was observed in three out of nine carriers of two (AAAGG)<sub>exp</sub> alleles and  
 178 one healthy subject with (AAGGG)<sub>exp</sub>/(AAAGG)<sub>exp</sub> alleles, suggesting its possible  
 179 association with both (AAGGG)<sub>exp</sub> and (AAAGG)<sub>exp</sub> configurations of the repeated unit,  
 180 but not (AAAAG)<sub>11</sub> or (AAAAG)<sub>exp</sub>.

181

182 **Clinical features of patients carrying the recessive AAGGG repeat expansion.** The  
 183 clinical features of 56 cases carrying the recessive intronic AAGGG repeat expansion,  
 184 including 23 familial and 33 sporadic cases, are summarized in **Table 1** and detailed in  
 185 **Supplementary Table 2**. All cases were of European ancestry. Apart from a higher  
 186 frequency of vestibular areflexia in familial CANVAS, clinical features were otherwise  
 187 similar in familial and sporadic cases; hence data are presented together. Mean age of  
 188 onset was 54 ± 9 (35-73) years, and mean disease duration at examination was 11 ± 7 (1-  
 189 30) years. The most common complaint at disease onset was unsteadiness, which was  
 190 reported by 84% of patients, and frequently described as being worse in the dark. 37% of  
 191 patients complained of chronic cough, which in some cases could precede by decades the  
 192 onset of the walking difficulties. Neurologic examination invariably showed signs in  
 193 keeping with a large fibre sensory neuropathy, 80% of patients had signs of cerebellar  
 194 involvement, and overall 54% had evidence of bilateral vestibular areflexia. 23% of  
 195 patients had concurrent autonomic nervous system involvement, particularity affecting  
 196 micturition and defecation. Nerve conduction studies confirmed the presence of a non-  
 197 length-dependent sensory neuropathy in all cases tested, as opposed to an entirely  
 198 normal motor conduction study in most patients. Cerebellar atrophy was identified in 35  
 199 (83%) of 42 cases who underwent an MRI or CT scan.

200

201 **Neuropathological examination.** Pathological examination was conducted in a patient  
 202 with CANVAS who carried the biallelic AAGGG repeat expansion and compared with a  
 203 patient with genetically confirmed Friedreich's ataxia, one patient with spinocerebellar  
 204 ataxia 17 (SCA17) and one case with *C9orf72*-related frontotemporal dementia (FTD), as  
 205 well as control brains (**Fig. 5**). The patient with CANVAS showed severe, widespread  
 206 depletion of Purkinje cells with associated prominent Bergmann gliosis, while cell density  
 207 in the granule cell layer was well preserved. Loss of Purkinje cells was also observed in  
 208 Friedreich's ataxia, SCA17 and, to a much lesser extent, in *C9orf72*-related FTD, but not in  
 209 control brain. Similar to Friedreich's ataxia and control brain, and as opposed to SCA17  
 210 and a *C9orf72*-related FTD, which were tested as positive controls, immunostaining for  
 211 p62 showed no pathological cytoplasmic or intranuclear inclusions in the cerebellar  
 212 cortex of the patient with CANVAS. Examination of the brain, in addition to prominent

213 cerebellar atrophy, revealed age-related changes in the form of neurofibrillary tangle tau  
214 pathology and amyloid- $\beta$  pathology (**Supplementary Fig. 3**).

215 Eight nerve biopsies and 10 muscle biopsies were also available for assessment  
216 from patients carrying the homozygous AAGGG repeat expansion. In all nerve biopsies,  
217 there was prominent widespread depletion of myelinated fibres, and the muscle biopsies  
218 confirmed chronic denervation with re-innervation (**Supplementary Fig. 4**).

219 Fluorescence *in situ* hybridization using sense (AAGGG)<sub>5</sub> and anti-sense (TTCCC)<sub>5</sub>  
220 repeat-specific oligonucleotides was performed on vermis post-mortem tissue from one  
221 CANVAS patient and disease and healthy controls. As opposed to SH-SY5Y cells  
222 transfected with pcDNA3.1/CT-GFP TOPO vector containing either (TTCCC)<sub>94</sub> or  
223 (AAGGG)<sub>54</sub>, in which intranuclear and cytoplasmic inclusion were clearly detectable, we  
224 did not observe the presence of endogenous RNA foci in any of the samples examined  
225 (**Supplementary Fig. 5**).

226  
227 **RNA sequencing.** We performed whole transcriptome analysis in order to assess the  
228 presence of changes in *RFC1* expression, as well as *cis* and *trans* effects at more distant  
229 genomic regions. RNA-seq data showed that *RFC1* mRNA was unchanged in CANVAS  
230 ( $n = 4$ ) and control ( $n = 4$ ) fibroblasts ( $P = 0.42$ ) and in CANVAS ( $n = 2$ ) and control ( $n = 3$ )  
231 lymphoblasts ( $P = 0.45$ ). We also performed RNA-seq from frontal cortex and cerebellar  
232 vermis from autopsied brains from one CANVAS patient, Friedreich's ataxia cases ( $n = 3$ )  
233 and controls without evidence of neurological disease ( $n = 3$ ). In the single CANVAS  
234 patient, *RFC1* appears to be unchanged in both cortex and cerebellum compared to the  
235 other samples (**Fig. 6a**). However, frataxin gene (*FXN*) was clearly down regulated in  
236 Friedreich's ataxia frontal cortex and cerebellum compared to controls (cerebellum  $P =$   
237  $0.007$ ;  $\log_2$  fold change =  $-1.2$ ; frontal cortex  $P = 0.0003$ ;  $\log_2$  fold change =  $-1.3$ ) (**Fig. 6a**).  
238 The single CANVAS sample resembled the controls for *FXN* expression.

239 There were no differentially expressed genes between patient and control  
240 fibroblasts, whereas 132 differentially expressed genes were identified between patient  
241 and control lymphoblasts. Gene Ontology analysis showed enrichment for immune  
242 terms, whose relevance to the disease will warrant further work. Notably, only eight  
243 differentially expressed genes were located on chromosome 4 and were all separated by  
244 at least 25 Mb from the locus of the repeat expansion. Analysis of differentially expressed  
245 genes in frontal cortex and vermis was not possible due to the limited numbers of  
246 CANVAS samples ( $n = 1$ ).

247 Splicing analysis was performed in lymphoblasts. We identified 145 exons in 108  
248 genes that had evidence of differential exon usage in CANVAS patients compared to  
249 healthy controls. Motif analysis for the alternatively spliced exons showed enrichment of  
250 motifs targeted by SRSF proteins, and in particular of SRSF3. *RFC1* did not show aberrant  
251 splicing of its coding exons in mature mRNA. Also, no reads containing the AAGGG or  
252 TTCCC repeated unit mapping to intron 2 of *RFC1* pre-mRNA transcript were detected,  
253 and no anti-sense or non-coding transcript was observed at the *RFC1* locus in any of the  
254 tissues examined. Gene Ontology analysis of alternatively spliced genes found

255 enrichment for focal adhesion and non-specific cellular response terms. Lists of  
256 differentially expressed genes and differentially expressed exons in lymphoblasts, their  
257 normalized count values in brain samples, and motif analysis for the alternatively spliced  
258 exons are provided in **Supplementary Data**.

259  
260 ***RFC1* expression in patients' tissues.** Quantitative reverse transcriptase PCR was  
261 performed using two sets of primers (**Fig. 6b**) and, concordantly with RNA-seq data, did  
262 not show any significant decrease of *RFC1* mRNA (RefSeq NM\_002913) level in patients'  
263 fibroblasts ( $n = 5$ ), lymphoblasts ( $n = 2$ ), muscle ( $n = 6$ ), frontal cortex and cerebellar  
264 vermis ( $n = 1$ ) compared to healthy controls or Friedreich's ataxia cases (**Fig. 6c**). Exon 2  
265 and 3 were correctly spliced in the mature *RFC1* mRNA as shown by RNA-seq, qRT-PCR  
266 and sequencing. However, assessment of pre-mRNA expression by qRT-PCR showed a  
267 consistent increase of intron 2 retention (IR) in patients' lymphoblasts ( $n = 2$ ), muscle ( $n =$   
268  $6$ ) ( $P = 0.0077$ ), cerebellar and frontal cortex ( $n = 1$ ) compared to healthy controls  
269 (**Supplementary Fig. 6**). The low level of *RFC1* expression in fibroblasts prevented the  
270 assessment of pre-mRNA processing.

271 Western blot showed that RFC1 protein (Uniprot P35251-1) was not decreased in  
272 patients' fibroblasts ( $n = 5$ ), lymphoblasts ( $n = 4$ ) or brain ( $n = 1$ ) compared to healthy  
273 controls or Friedreich's ataxia cases (**Fig. 6d** and **Supplementary Fig. 7**). Assessment of  
274 RFC1 protein expression in muscle could not be performed due to limited tissue  
275 availability.

276 Since RFC1 plays a key role in DNA damage recognition and recruitment of DNA  
277 repair enzymes, we assessed whether patient-derived fibroblasts have an impaired  
278 response to DNA damage. Patients' fibroblasts did not show an increased susceptibility  
279 to DNA damage, and their treatment with double-stranded DNA break-inducing agents,  
280 UV and methyl methanesulfonate, triggered a grossly normal response to DNA damage  
281 (**Supplementary Fig. 8**).

282  
283

## 284 DISCUSSION

285 We identified a recessive repeat expansion in intron 2 of *RFC1* as a cause of CANVAS  
286 and late-onset ataxia. Twenty-three cases from 11 families and 33 sporadic cases carried  
287 the biallelic AAGGG repeat expansion. Notably, out of 150 cases from a single centre  
288 diagnosed with late-onset ataxia, 22% tested positive for the biallelic AAGGG repeat  
289 expansion, and the percentage was higher if only patients with sensory neuropathy  
290 and cerebellar involvement (62%), CANVAS disease (92%) and familial CANVAS disease  
291 (100%) were considered, highlighting that a higher diagnostic can be achieved in cases  
292 with well-defined clinical features and positive family history. Not since the discovery  
293 two decades ago of the most common genes causing ataxia<sup>19-22</sup> and Charcot-Marie-Tooth  
294 (CMT) disease<sup>23-26</sup> has a novel gene explained percentages above 10% of genetically  
295 undetermined cases<sup>27,28</sup>.

296 We determined that the allelic carrier frequency of the AAGGG repeat expansion  
297 in healthy controls was 0.7%, which is similar to the allelic carrier frequency of the GAA  
298 expansion in *FXN* ranging from 0.9 to 1.6%, and which in the biallelic state causes the  
299 most common recessive ataxia, Friedreich's ataxia. Together, these data suggest that the  
300 recessive AAGGG expansion in *RFC1* may represent a frequent cause of late-onset ataxia  
301 in the general population, with an estimated prevalence at birth of the recessive trait of  
302 ~1/20,000.

303 The expansion resides at the 3'-end of a deep intronic AluSx3 element, and it  
304 increases the polyA-tail size from 11 to over 400 repeated units, but also alters its  
305 sequence. Of interest, expansions in terminal and mid A stretches of Alu elements have  
306 been previously identified to cause Friedreich's ataxia<sup>19</sup>, SCA37 (ref. 29), more recently  
307 benign adult familial myoclonic epilepsy (BAFME)<sup>30</sup> and now CANVAS and late-onset  
308 ataxia. Together, these observations suggest that variations and expansion of these highly  
309 polymorphic regions of Alu elements represent a common mechanism underlying  
310 different inherited neurological disorders. Notably, both SCA37 and BAFME are  
311 characterized by expansion of a mutated repeated unit, ATTTC and TTTCA,  
312 respectively<sup>29,30</sup>. In this study, as well as in BAFME and SCA37, the presence in the  
313 normal population of large expansions of the reference repeated unit suggests that the  
314 nucleotide change rather than the size of the expansion may be the driving pathogenic  
315 mechanism

316 Alu elements are repetitive elements about 300 bp long and are highly conserved  
317 within primate genomes. The 3'-end of an Alu element has a longer A-rich region that  
318 plays a critical role in its amplification mechanism<sup>31</sup>. Active elements degrade rapidly on  
319 an evolutionary time scale by A-tail shortening or heterogeneous base interruptions  
320 accumulating in the A-tail, such as G insertions. We hypothesize that the mutation of the  
321 AAGGG repeated unit occurred as part of the inactivation process by G interruption of  
322 the polyA tail of the retrotransposon AluSx3. As known, repetitive DNA motifs,  
323 particularly G-rich regions, can form secondary or tertiary nucleotide structures such as  
324 hairpins, parallel and antiparallel G-quadruplexes and, if transcribed, DNA-RNA hybrids  
325 also known as R loops. These structures have been shown to increase the exposure of  
326 single-stranded DNA to damaging environmental agents and can initiate repeat  
327 expansion and perpetrate genomic instability across meiotic and mitotic divisions or after  
328 DNA damage<sup>32</sup>.

329 Since the same ancestral haplotype is shared by the majority of familial and  
330 positive cases as well as some healthy carriers of two (AAAGG)<sub>exp</sub> alleles, we speculate  
331 that nucleotide change AAAAG to AAAGG or AAGGG may represent an ancestral  
332 founder event, which was followed by the pathologic expansion of the repeated unit,  
333 whose size seems to correlate positively with its GC content. However, the identification  
334 of two patients (Fam 5b-2 and s23) with a recessive AAGGG repeat expansion who share  
335 only one allele of the common haplotype implies that repeat expansions of the mutated  
336 AAGGG unit can occur also on a different genetic background. Interestingly, Fam 5b-2

337 was also found to carry the largest repeat expansion (10 kb or 2,000 repeats) among the  
338 cohort of patients tested.

339 In the majority of the patients, the expansion encompassed 1,000 repeats, but as  
340 low as 400 AAGGG repeats were shown to be sufficient to cause disease. The size of  
341 expanded alleles was relatively stable in siblings within single families, but no parent of  
342 the affected patients was available to assess whether this also applies across generations.  
343 We did not observe a correlation between age of onset of the neuropathy and size of the  
344 repeat expansion, although the disease course was very slowly progressive and initial  
345 symptoms might have been neglected in some patients but reported by others.

346 So far, approximately 40 neurological or neuromuscular genetic disorders have  
347 been associated with nucleotide repeat expansions. Two of them are known to be  
348 inherited in a recessive mode, namely Friedreich's ataxia and myoclonic epilepsy type 1,  
349 and both are associated with loss-of-function of the repeat-hosting gene<sup>33-35</sup>.

350 A remarkable aspect of the recessive expansion described here is that our data do  
351 not suggest a direct mechanism of loss of function for *RFC1*. We did not observe a  
352 reduced level of *RFC1* expression at either transcript or protein level in CANVAS  
353 patients, although as a known loss-of-function control, we were able to detect a  
354 significant reduction of *FXN* transcript in post-mortem brain from patients with  
355 Friedreich's ataxia. Also, RNA-seq data did not show a clear effect on the expression of  
356 neighboring or distant genes. We cannot exclude that the repeat expansion may cause  
357 more subtle tissue-specific alterations of *RFC1* transcript and protein or alter the  
358 structural organization of the chromatin.

359 *RFC1* encodes the large subunit of replication factor C, a five subunit DNA  
360 polymerase accessory protein. It loads PCNA onto DNA and activates DNA polymerases  
361 delta and epsilon to promote the coordinated synthesis of both strands during replication  
362 or after DNA damage<sup>36-38</sup>. It is interesting to note that mutations in many of the genes  
363 involved in DNA repair have been already associated with degenerative neurological  
364 disorders<sup>39</sup>, including ataxia-telangiectasia, xeroderma pigmentosum, Cockayne  
365 syndrome, and ataxia oculomotor apraxia 1 and 2. Interestingly, ataxia and neuropathy  
366 are common clinical features to all of them, suggesting a particular susceptibility of  
367 cerebellum and peripheral nerves to DNA damage. However, our preliminary study did  
368 not show an impaired response to DNA damage in patient-derived fibroblasts.

369 In fact, late-onset Mendelian disorders represent a unique interpretative challenge,  
370 as risk variants may exert subtle effects, rather than a clear loss of function of the mutated  
371 gene, that are compatible with normal developmental until adult or old age<sup>40</sup>. In this  
372 regard, although unusual in the context of a recessive mode of inheritance, other  
373 mechanisms, including the production of toxic RNA containing the expanded repeat, or  
374 the translation of a repeat-encoded polypeptide, should be considered<sup>41</sup>. We did not  
375 observe in patient's brain the presence of RNA foci of either the sense or anti-sense  
376 repeated unit. However, we were able to detect a consistent increase across different  
377 tissues of the retention of intron 2 in *RFC1* pre-mRNA. Retention of the repeat-hosting  
378 intron was recently identified as a common event associated with other disease-causing

379 GC-rich intronic expansions, such as in myotonic dystrophy type 2 and *C9orf72*-ALS/FTD  
 380 but not AT-rich repeat expansions such as in Friedreich's ataxia<sup>42</sup>. Intron retention and  
 381 abnormal pre-mRNA processing bear potential effects on nuclear retention and  
 382 nucleocytoplasmic transport of the pre-mRNA, which, if efficiently exported to the  
 383 cytoplasm, would be accessible to the translational machinery.

384 Notwithstanding the enormous progress in Mendelian gene identification during  
 385 the last decade, up to 40% of patients with ataxia and inherited neuropathy remain  
 386 genetically undiagnosed, and the percentage can rise up to 80-90% in particular subtypes,  
 387 such as late-onset ataxia<sup>2,5,43</sup> and hereditary sensory neuropathies<sup>27,28</sup>. Our paper, together  
 388 with other studies from recent years<sup>30,44-46</sup>, provides evidence that the combined use of  
 389 whole-genome sequencing and classical genetic investigations such as linkage analysis  
 390 can provide a powerful tool to unravel a part of the missing heritability hidden in non-  
 391 coding regions of the human genome

392

393

#### 394 **ACKNOWLEDGMENTS**

395 A.C. is funded by the inherited neuropathy consortium, which is a part of the NIH Rare Diseases  
 396 Clinical Research Network (RDCRN) (U54NS065712) and Wellcome Trust (204841/Z/16/Z).  
 397 A.M.R. is funded by a Wellcome Trust Postdoctoral Fellowship for Clinicians (110043/Z/15/Z).  
 398 H.H. is also supported by Rosetrees Trust, Ataxia UK, The MSA Trust, Brain Research UK,  
 399 MDUK, The Muscular Dystrophy Association (MDA), Higher Education Commission (HEC) of  
 400 Pakistan and The Wellcome Trust (Synaptopathies Strategic Award). The INC (U54NS065712) is a  
 401 part of the NCATS Rare Diseases Clinical Research Network (RDCRN). RDCRN is an initiative of  
 402 the Office of Rare Diseases Research (ORDR), NCATS, funded through a collaboration between  
 403 NCATS and the NINDS. S.Z. thanks the National Institute of Health (4R01NS075764) for its  
 404 support. This research was also supported by the National Institute for Health Research University  
 405 College London Hospitals Biomedical Research Centre (BRC). Neuromuscular and brain tissue  
 406 samples were obtained from University College London Hospitals NHS Foundation Trust as part of  
 407 the UK Brain Archive Information Network (BRAIN UK), which is funded by the Medical  
 408 Research Council and Brain Tumour Research and the NIH funded NeuroBioBank. We also thank  
 409 Francesca Launchbury from UCL IQPath laboratory for technical assistance in histology slide  
 410 preparation.

411

#### 412 **AUTHOR CONTRIBUTIONS**

413 A.C. designed the study, collected clinical data, performed the genetic analysis that led to the  
 414 discovery of the AAGGG repeat expansions, analyzed the data, drafted the manuscript together with  
 415 contributions from J.V., R. Simone, R. Sullivan, and J.H. R. Simone, N.S.A., E.T., E.B., A.R., Y.W.Y.,  
 416 and M.I. performed the investigation on *RFC1* expression. J.V. performed the computational genetic  
 417 analysis. R. Sullivan and H.T. collected and analyzed the genetic data in healthy controls. P.J.T.,  
 418 W.J.M., A.B., G.D., I.C., M.V., D.K., V.S., S.E., N.W.W. and A.M.R. contributed with collection of  
 419 clinical data and patients' samples. J.H., P.S. and P.F. performed the RNA-seq analysis. Z.J. performed  
 420 the pathological investigation. R. Simone, A.M.R., P.F., and J.P. contributed to the design of the study.  
 421 S.Z. contributed to the design of the study and analyzed the data. H.H. and M.M.R. designed the study,  
 422 collected patients' clinical data and biological samples and analyzed the data. All authors revised the  
 423 manuscript.

424

425

#### 426 **COMPETING INTERESTS STATEMENT**

427 The authors declare no competing interests.

428 **REFERENCES**

429

430

431 1. Harding AE. "Idiopathic" late onset cerebellar ataxia. A clinical and genetic study of 36 cases.  
432 J Neurol Sci. 1981 Aug;51(2):259–71.

433 2. Muzaimi MB, Thomas J, Palmer-Smith S, Rosser L, Harper PS, Wiles CM, et al. Population  
434 based study of late onset cerebellar ataxia in south east Wales. J Neurol Neurosurg Psychiatry.  
435 2004 Aug;75(8):1129–34.

436 3. Sghirlanzoni A, Pareyson D, Lauria G. Sensory neuron diseases. Lancet Neurol. 2005  
437 Jun;4(6):349–61.

438 4. Strupp M, Feil K, Dieterich M, Brandt T. Bilateral vestibulopathy. Handb Clin Neurol.  
439 2016;137:235–40.

440 5. Abele M, Bürk K, Schöls L, Schwartz S, Besenthal I, Dichgans J, et al. The aetiology of  
441 sporadic adult-onset ataxia. Brain J Neurol. 2002 May;125(Pt 5):961–8.

442 6. Kirchner H, Kremmyda O, Hüfner K, Stephan T, Zingler V, Brandt T, et al. Clinical,  
443 electrophysiological, and MRI findings in patients with cerebellar ataxia and a bilaterally  
444 pathological head-impulse test. Ann N Y Acad Sci. 2011 Sep;1233:127–38.

445 7. Migliaccio AA, Halmagyi GM, McGarvie LA, Cremer PD. Cerebellar ataxia with bilateral  
446 vestibulopathy: description of a syndrome and its characteristic clinical sign. Brain J Neurol.  
447 2004 Feb;127(Pt 2):280–93.

448 8. Szmulewicz DJ, Roberts L, McLean CA, MacDougall HG, Halmagyi GM, Storey E. Proposed  
449 diagnostic criteria for cerebellar ataxia with neuropathy and vestibular areflexia syndrome  
450 (CANVAS). Neurol Clin Pract. 2016 Feb;6(1):61–8.

451 9. Wu TY, Taylor JM, Kilfoyle DH, Smith AD, McGuinness BJ, Simpson MP, et al. Autonomic  
452 dysfunction is a major feature of cerebellar ataxia, neuropathy, vestibular areflexia  
453 "CANVAS" syndrome. Brain J Neurol. 2014 Oct;137(Pt 10):2649–56.

454 10. Szmulewicz DJ, Merchant SN, Halmagyi GM. Cerebellar ataxia with neuropathy and bilateral  
455 vestibular areflexia syndrome: a histopathologic case report. Otol Neurotol Off Publ Am Otol  
456 Soc Am Neurotol Soc Eur Acad Otol Neurotol. 2011 Oct;32(8):e63-65.

457 11. Szmulewicz DJ, McLean CA, Rodriguez ML, Chancellor AM, Mossman S, Lamont D, et al.  
458 Dorsal root ganglionopathy is responsible for the sensory impairment in CANVAS.  
459 Neurology. 2014 Apr 22;82(16):1410–5.

460 12. Cazzato D, Bella ED, Dacci P, Mariotti C, Lauria G. Cerebellar ataxia, neuropathy, and  
461 vestibular areflexia syndrome: a slowly progressive disorder with stereotypical presentation. J  
462 Neurol. 2016 Feb;263(2):245–9.

463 13. Rust H, Peters N, Allum JHJ, Wagner B, Honegger F, Baumann T. VEMPs in a patient with  
464 cerebellar ataxia, neuropathy and vestibular areflexia (CANVAS). J Neurol Sci. 2017 Jul  
465 15;378:9–11.

- 466 14. Pelosi L, Leadbetter R, Mulroy E, Chancellor AM, Mossman S, Roxburgh R. Peripheral nerve  
467 ultrasound in cerebellar ataxia neuropathy vestibular areflexia syndrome (CANVAS). *Muscle*  
468 *Nerve*. 2017 Jul;56(1):160–2.
- 469 15. Pelosi L, Mulroy E, Leadbetter R, Kilfoyle D, Chancellor AM, Mossman S, et al. Peripheral  
470 nerves are pathologically small in cerebellar ataxia neuropathy vestibular areflexia syndrome:  
471 a controlled ultrasound study. *Eur J Neurol*. 2018 Apr;25(4):659–65.
- 472 16. Taki M, Nakamura T, Matsuura H, Hasegawa T, Sakaguchi H, Morita K, et al. Cerebellar  
473 ataxia with neuropathy and vestibular areflexia syndrome (CANVAS). *Auris Nasus Larynx*.  
474 2018 Aug;45(4):866–70.
- 475 17. Infante J, García A, Serrano-Cárdenas KM, González-Aguado R, Gazulla J, de Lucas EM, et  
476 al. Cerebellar ataxia, neuropathy, vestibular areflexia syndrome (CANVAS) with chronic  
477 cough and preserved muscle stretch reflexes: evidence for selective sparing of afferent Ia  
478 fibres. *J Neurol*. 2018 Jun;265(6):1454–62.
- 479 18. Hommelsheim CM, Frantzeskakis L, Huang M, Ülker B. PCR amplification of repetitive  
480 DNA: a limitation to genome editing technologies and many other applications. *Sci Rep*. 2014  
481 May 23;4:5052.
- 482 19. Campuzano V, Montermini L, Moltò MD, Pianese L, Cossée M, Cavalcanti F, et al.  
483 Friedreich's ataxia: autosomal recessive disease caused by an intronic GAA triplet repeat  
484 expansion. *Science*. 1996 Mar 8;271(5254):1423–7.
- 485 20. Orr HT, Chung MY, Banfi S, Kwiatkowski TJ, Servadio A, Beaudet AL, et al. Expansion of  
486 an unstable trinucleotide CAG repeat in spinocerebellar ataxia type 1. *Nat Genet*. 1993  
487 Jul;4(3):221–6.
- 488 21. Pulst SM, Nechiporuk A, Nechiporuk T, Gispert S, Chen XN, Lopes-Cendes I, et al. Moderate  
489 expansion of a normally biallelic trinucleotide repeat in spinocerebellar ataxia type 2. *Nat*  
490 *Genet*. 1996 Nov;14(3):269–76.
- 491 22. Kawaguchi Y, Okamoto T, Taniwaki M, Aizawa M, Inoue M, Katayama S, et al. CAG  
492 expansions in a novel gene for Machado-Joseph disease at chromosome 14q32.1. *Nat Genet*.  
493 1994 Nov;8(3):221–8.
- 494 23. Lupski JR, de Oca-Luna RM, Slaugenhaupt S, Pentao L, Guzzetta V, Trask BJ, et al. DNA  
495 duplication associated with Charcot-Marie-Tooth disease type 1A. *Cell*. 1991 Jul  
496 26;66(2):219–32.
- 497 24. Hayasaka K, Himoro M, Sato W, Takada G, Uyemura K, Shimizu N, et al. Charcot-Marie-  
498 Tooth neuropathy type 1B is associated with mutations of the myelin P0 gene. *Nat Genet*.  
499 1993 Sep;5(1):31–4.
- 500 25. Bergoffen J, Scherer SS, Wang S, Scott MO, Bone LJ, Paul DL, et al. Connexin mutations in  
501 X-linked Charcot-Marie-Tooth disease. *Science*. 1993 Dec 24;262(5142):2039–42.
- 502 26. Züchner S, Mersiyanova IV, Muglia M, Bissar-Tadmouri N, Rochelle J, Dadali EL, et al.  
503 Mutations in the mitochondrial GTPase mitofusin 2 cause Charcot-Marie-Tooth neuropathy  
504 type 2A. *Nat Genet*. 2004 May;36(5):449–51.

- 505 27. Fridman V, Bundy B, Reilly MM, Pareyson D, Bacon C, Burns J, et al. CMT subtypes and  
506 disease burden in patients enrolled in the Inherited Neuropathies Consortium natural history  
507 study: a cross-sectional analysis. *J Neurol Neurosurg Psychiatry*. 2015 Aug;86(8):873–8.
- 508 28. Murphy SM, Laura M, Fawcett K, Pandraud A, Liu Y-T, Davidson GL, et al. Charcot-Marie-  
509 Tooth disease: frequency of genetic subtypes and guidelines for genetic testing. *J Neurol*  
510 *Neurosurg Psychiatry*. 2012 Jul;83(7):706–10.
- 511 29. Seixas AI, Loureiro JR, Costa C, Ordóñez-Ugalde A, Marcelino H, Oliveira CL, et al. A  
512 Pentanucleotide ATTTTC Repeat Insertion in the Non-coding Region of DAB1, Mapping to  
513 SCA37, Causes Spinocerebellar Ataxia. *Am J Hum Genet*. 2017 Jul 6;101(1):87–103.
- 514 30. Ishiura H, Doi K, Mitsui J, Yoshimura J, Matsukawa MK, Fujiyama A, et al.  
515 Expansions of intronic TTTCA and TTTTA repeats in benign adult familial myoclonic  
516 epilepsy. *Nat Genet*. 2018 Apr;50(4):581–90.
- 517 31. Deininger P. Alu elements: know the SINEs. *Genome Biol*. 2011 Dec 28;12(12):236.
- 518 32. Haeusler AR, Donnelly CJ, Rothstein JD. The expanding biology of the C9orf72 nucleotide  
519 repeat expansion in neurodegenerative disease. *Nat Rev Neurosci*. 2016;17(6):383–95.
- 520 33. Dürr A, Cossee M, Agid Y, Campuzano V, Mignard C, Penet C, et al. Clinical and genetic  
521 abnormalities in patients with Friedreich’s ataxia. *N Engl J Med*. 1996 Oct 17;335(16):1169–  
522 75.
- 523 34. Lazaropoulos M, Dong Y, Clark E, Greeley NR, Seyer LA, Brigatti KW, et al. Frataxin levels  
524 in peripheral tissue in Friedreich ataxia. *Ann Clin Transl Neurol*. 2015 Aug;2(8):831–42.
- 525 35. Paulson H. Repeat expansion diseases. *Handb Clin Neurol*. 2018;147:105–23.
- 526 36. Majka J, Burgers PMJ. The PCNA-RFC families of DNA clamps and clamp loaders. *Prog*  
527 *Nucleic Acid Res Mol Biol*. 2004;78:227–60.
- 528 37. Tomida J, Masuda Y, Hiroaki H, Ishikawa T, Song I, Tsurimoto T, et al. DNA damage-  
529 induced ubiquitylation of RFC2 subunit of replication factor C complex. *J Biol Chem*. 2008  
530 Apr 4;283(14):9071–9.
- 531 38. Overmeer RM, Gourdin AM, Giglia-Mari A, Kool H, Houtsmuller AB, Siegal G, et al.  
532 Replication factor C recruits DNA polymerase delta to sites of nucleotide excision repair but is  
533 not required for PCNA recruitment. *Mol Cell Biol*. 2010 Oct;30(20):4828–39.
- 534 39. McKinnon PJ. Maintaining genome stability in the nervous system. *Nat Neurosci*. 2013  
535 Nov;16(11):1523–9.
- 536 40. Higuchi Y, Hashiguchi A, Yuan J, Yoshimura A, Mitsui J, Ishiura H, et al. Mutations in MME  
537 cause an autosomal-recessive Charcot-Marie-Tooth disease type 2. *Ann Neurol*. 2016  
538 Apr;79(4):659–72.
- 539 41. La Spada AR, Taylor JP. Repeat expansion disease: progress and puzzles in disease  
540 pathogenesis. *Nat Rev Genet*. 2010 Apr;11(4):247–58.

- 541 42. Sznajder ŁJ, Thomas JD, Carrell EM, Reid T, McFarland KN, Cleary JD, et al. Intron  
542 retention induced by microsatellite expansions as a disease biomarker. *Proc Natl Acad Sci U S*  
543 *A*. 2018 17;115(16):4234–9.
- 544 43. Gebus O, Montaut S, Monga B, Wirth T, Cheraud C, Alves Do Rego C, et al. Deciphering the  
545 causes of sporadic late-onset cerebellar ataxias: a prospective study with implications for  
546 diagnostic work. *J Neurol*. 2017 Jun;264(6):1118–26.
- 547 44. DeJesus-Hernandez M, Mackenzie IR, Boeve BF, Boxer AL, Baker M, Rutherford NJ, et al.  
548 Expanded GGGGCC hexanucleotide repeat in noncoding region of C9ORF72 causes  
549 chromosome 9p-linked FTD and ALS. *Neuron*. 2011 Oct 20;72(2):245–56.
- 550 45. Renton AE, Majounie E, Waite A, Simón-Sánchez J, Rollinson S, Gibbs JR, et al. A  
551 hexanucleotide repeat expansion in C9ORF72 is the cause of chromosome 9p21-linked ALS-  
552 FTD. *Neuron*. 2011 Oct 20;72(2):257–68.
- 553 46. Aneichyk T, Hendriks WT, Yadav R, Shin D, Gao D, Vaine CA, et al. Dissecting the Causal  
554 Mechanism of X-Linked Dystonia-Parkinsonism by Integrating Genome and Transcriptome  
555 Assembly. *Cell*. 2018 Feb 22;172(5):897-909.e21.
- 556
- 557

558 **FIGURES LEGENDS**

559

560 **Fig. 1 | Clinical spectrum and pedigrees of late-onset ataxia.** **a**, Clinical spectrum of  
 561 idiopathic late-onset ataxia from isolated cerebellar, vestibular and sensory variants to  
 562 full-blown CANVAS. ILOCA, idiopathic late-onset cerebellar ataxia; CANVAS, cerebellar  
 563 ataxia, neuropathy, vestibular areflexia syndrome. **b**, Pedigrees of CANVAS families.  
 564 Squares indicate males and circles females. Diagonal lines are used for deceased  
 565 individuals. CANVAS patients are indicated with filled symbols. Black dots indicate  
 566 genotyped individuals. Red dots indicate patients enrolled for whole-genome sequencing  
 567 study.

568

569 **Fig. 2 | Identification of CANVAS locus.** **a**, Non-parametric multipoint linkage analysis  
 570 identifies a unique locus associated with the disease in chromosomal region 4p14 with  
 571 maximal HLOD score of 5.8. **b**, Schematic representation of shared haplotypes within  
 572 single families. Light blue bars indicate a genomic region shared by affected siblings in a  
 573 family and for which unaffected siblings are discordant. Two red dashed lines define a  
 574 1.7-Mb region common to the different families. Single-nucleotide polymorphisms  
 575 defining the haplotypes are represented on the top line. **c**, Fine-mapping inside the 1.7-  
 576 Mb region identifies a recessive haplotype shared by all distinct families (green  
 577 highlighted), except for individual Fam 5b-2, who likely shares only one allele (light  
 578 green highlighted). **d**, Schematic representation of the candidate 1.7-Mb region  
 579 encompassing all 24 exons and flanking regions of *RFC1* and the last exon and flanking  
 580 intron of *WDR19*.

581

582 **Fig. 3 | A recessive expansion of a mutated AAGGG repeated unit in intron 2 of *RFC1***  
 583 **causes CANVAS and late-onset ataxia in familial and sporadic cases.** **a**, A reduced read  
 584 depth of whole genome sequencing is observed in CANVAS patients ( $n = 6$ ) in a region  
 585 corresponding to a short tandem AAAAG repeat in intron 2 of *RFC1*. STR, short tandem  
 586 repeat. **b**, Visualization on IGV of reads aligned to the short repeat and flanking region  
 587 shows in patients ( $n = 6$ ) the presence of a mutated AAGGG repeat unit (representative  
 588 image). Reads from both sides are interrupted and are unable to cover the entire length of  
 589 the microsatellite region. Note that, as per IGV default setting, AAGGG repeated units  
 590 which do not map to the (AAAAG)<sub>11</sub> reference sequence are soft-clipped and do not  
 591 contribute to the coverage of the STR in **a**, which is virtually absent. However,  $\geq 20$  reads  
 592 containing the AAGGG repeated unit could be observed in each patient if soft-clipped  
 593 reads are shown. **c**, Repeat-primed PCR (RPPCR) targeting the mutated AAGGG  
 594 repeated unit. FAM-labelled PCR products are separated on an ABI3730 DNA Analyzer.  
 595 Electropherograms are visualized on GENEMAPPER at 2,000 relative fluorescence units.  
 596 Representative plots from a patient carrying the AAGGG repeat expansion and one non-  
 597 carrier are shown. RPPCR experiments were repeated independently twice with similar  
 598 results. **d**, Sanger sequencing of long-range PCR reactions confirms in patients the  
 599 AAAAG to AAGGG nucleotide change of the repeated unit.

600

601 **Fig. 4 | Polymorphic configurations of the repeat expansion locus and allelic**  
 602 **distribution in healthy controls. a**, Schematic representation of the repeat expansion  
 603 locus in intron 2 of *RFC1* and its main allelic variants. **b**, Estimated allelic frequencies in  
 604 608 chromosomes from 304 healthy controls. **c**, Average size and standard deviation of  
 605  $(AAAAG)_{exp}$  ( $n = 24$ ) and  $(AAAGG)_{exp}$  ( $n = 30$ ) expansions in healthy controls and  
 606  $(AAAGG)_{exp}$  ( $n = 72$ ) in controls and CANVAS patients

607

608 **Fig. 5 | Pathology of cerebellar degeneration in a patient with CANVAS carrying the**  
 609 **recessive AAGGG repeat expansion. a-j**, Haematoxylin and eosin (H&E) stained  
 610 sections (**a-e**) and sections immunostained for p62 (**f-j**). In a control brain (**a**), age-  
 611 matched for the patient with CANVAS syndrome, there is well preserved density of  
 612 Purkinje cells (yellow arrow) and also granule cell layer is densely populated with small  
 613 neurocytes (green asterisk). In CANVAS syndrome (**b**), there is severe, widespread  
 614 depletion of Purkinje cells with associated prominent Bergmann gliosis (blue arrow),  
 615 while cell density in the granule cell layer is well preserved. In a patient with genetically  
 616 confirmed Friedreich's ataxia (**c**), there is patchy depletion of Purkinje cells associated  
 617 with Bergmann gliosis and unremarkable appearance of the granule cell layer. In a  
 618 patient with genetically confirmed spinocerebellar ataxia 17 (SCA17) (**d**), there is  
 619 widespread Purkinje cell loss with only occasional Purkinje cells remaining; also, in this  
 620 patient, granule cell layer is densely populated with small neurocytes. In a patient with  
 621 frontotemporal dementia due to *C9orf72* expansion (**e**), the Purkinje cell loss is patchy and  
 622 granule cell layer is unremarkable. Immunostaining for p62 shows no pathological  
 623 cytoplasmic or intranuclear inclusions in the cerebellar cortex in the control patient (**f**),  
 624 the patient with CANVAS syndrome (**g**) and also in the patient with Friedreich's ataxia  
 625 (**h**). In the SCA17 patient, there are scattered discrete intranuclear p62 immunoreactive  
 626 inclusions in the small neurones within granule cell layer (**i**; high-power view of a  
 627 representative intranuclear inclusion is demonstrated in the inset within **i**). In the patient  
 628 with *C9orf72* expansion, there are frequent characteristic perinuclear p62 positive  
 629 inclusions in the granule cell layer (**j** and high-power view of a representative inclusion is  
 630 shown in the inset within **j**). Scale bar: 100  $\mu\text{m}$  in **a-e**, 30  $\mu\text{m}$  in **f-j**, and 5  $\mu\text{m}$  in insets in **i**  
 631 and **j**. Stainings were carried out once on patients' samples with appropriate controls  
 632 according to standard practice and histopathology procedures in an ISO15189 accredited  
 633 laboratory.

634

635 **Fig. 6 | *RFC1* expression is not affected by the AAGGG repeat expansion. a**, Plots  
 636 showing expression levels of *RFC1* and *FXN* in controls ( $n = 3$ ), patients with Friedreich's  
 637 ataxia ( $n = 2$ ) and one CANVAS patient ( $n = 1$ ) in post-mortem cerebellum and frontal  
 638 cortex. **b**, Mapping on *RFC1* transcript 1 of the primers used for assessment by qRT-PCR  
 639 of *RFC1* mRNA (cF1-cR1 and cF2-cR2) and pre-mRNA (cF1/iR1) expression. Blue arrows  
 640 indicate primers mapping to exonic and intronic regions of canonical *RFC1* transcript.  
 641 Primers spanning across exonic junctions are connected by dotted lines. A red triangle

642 indicates the site of the AAGGG repeat expansion. **c**, Expression levels of the canonical  
643 coding *RFC1* mRNA as measured by qRT-PCR using two separate set of primers cF1-cR1  
644 and cF2-cR2 in control ( $n = 3$ ) and CANVAS ( $n = 2$ ) lymphoblasts, control ( $n = 5$ ) and  
645 CANVAS ( $n = 5$ ) fibroblasts, control ( $n = 3$ ), Friedreich's ataxia ( $n = 3$ ) and CANVAS ( $n =$   
646 1) cerebellum and frontal cortex, and control ( $n = 5$ ) and CANVAS muscles ( $n = 6$ ). **d**,  
647 *RFC1*-encoded protein levels as measured by Western blotting using the polyclonal  
648 antibody (GTX129291) and normalized to  $\beta$ -actin in control ( $n = 5$ ) and CANVAS ( $n = 5$ )  
649 fibroblasts, control ( $n = 3$ ) and CANVAS ( $n = 3$ ) lymphoblasts, control ( $n = 3$ ), Friedreich's  
650 ataxia ( $n = 3$ ) and CANVAS ( $n = 1$ ) post-mortem cerebellum and frontal cortex. Bar graphs  
651 show mean  $\pm$  s.d. and data distribution (black dots). Two-tailed *t*-test was performed to  
652 compare *RFC1* transcript and encoded protein expression in patients versus healthy or  
653 disease controls. All experiments were repeated independently twice with similar results.  
654 CANVAS, cerebellar ataxia, neuropathy, vestibular areflexia syndrome; CBM,  
655 cerebellum; Ctrl, control; FBs, fibroblasts; FCX, frontal cortex; FRDA, Friedreich's ataxia;  
656 *FXN*, frataxin; LBLs, lymphoblasts; *RFC1*, replication factor C subunit 1.  
657  
658  
659

## 660 TABLES

661

662 **Table 1 | Clinical features of patients with familial or sporadic late-onset ataxia**  
663 **carrying the recessive AAGGG repeat expansion in *RFC1***

664

	<b>Familial cases (<i>n</i> = 23)</b>	<b>Sporadic cases (<i>n</i> = 33)</b>	<b>All cases (<i>n</i> = 56)</b>	<b><i>P</i>-value</b>
<b>Male</b>	12 (52%)	11 (52%)	27 (48%)	NS
<b>Age of onset</b>	53 ± 8	54 ± 10	54 ± 9	NS
<b>Disease duration at examination</b>	13 ± 9	10 ± 6	11 ± 7	NS
<b>Sensory neuropathy</b>	23 (100%)	33 (100%)	56 (100%)	NS
<b>Cerebellar syndrome</b>	18 (78%)	27 (82%)	45 (80%)	NS
<b>Bilateral vestibular impairment</b>	17 (74%)	13 (39%)	30 (53%)	0.01
<b>Dysautonomia</b>	4 (17%)	9 (27%)	13 (23%)	NS
<b>Cough</b>	7 (30%)	14 (42%)	21 (37%)	NS
<b>SAPs upper limbs</b>				NS
<b>Reduced</b>	6/21 (29%)	4/31 (13%)	10/46 (22%)	
<b>Absent</b>	15/21 (71%)	27/31 (87%)	36/46 (78%)	
<b>SAPs lower limbs</b>				NS
<b>Reduced</b>	2/21 (10%)	1/31 (3%)	3/52 (6%)	
<b>Absent</b>	19/21 (90%)	30/31 (97%)	49/52 (94%)	
<b>Normal motor conduction</b>	19/21 (90%)	26/31 (84%)	45/52 (87%)	NS
<b>Cerebellar atrophy at CT/MRI scan</b>	14/17 (82%)	21/25 (84%)	35/42 (83%)	NS
<b>Full-blown CANVAS syndrome</b>	15 (65%)	11 (33%)	26 (46%)	0.02

665 cMAP, compound motor action potential; CT, computed tomography; MRI, magnetic resonance  
666 imaging; NS, not significant; SAP, sensory action potential.

667 **METHODS**

668 **Patients.** For the initial linkage study, we enrolled 29 individuals (23 affected and 6  
669 unaffected) from 11 families with a clinical diagnosis of CANVAS across four centres:  
670 National Hospital for Neurology and Neurosurgery (London, UK), C. Mondino National  
671 Neurological Institute (Pavia, Italy), C. Besta Neurological Institute and Department of  
672 Neurology, School of Medicine (Ribeirão Preto, Brazil).

673 An additional 150 patients with sporadic CANVAS or late onset ataxia (onset after  
674 35 years of age) were identified from the neurogenetic database of the National Hospital  
675 for Neurology and Neurosurgery (London, UK). For the experimental procedures,  
676 patients' samples are generally referred to as CANVAS, and no distinction between  
677 samples from patients with full-blown CANVAS or other more limited variants of late-  
678 onset ataxia is made. A skin biopsy was performed in five (Fam 1-3, Fam 2-2, Fam 5a-2,  
679 Fam 5b-2, Fam 6b-1) genetically confirmed subjects and six age- and gender-matched  
680 controls. Fibroblast cultures were maintained according to standard procedures<sup>47</sup>.  
681 Epstein-Barr virus-transformed lymphoblast cultures from four patients (Fam 6-1, Fam 8-  
682 1, Fam 8-2, Fam 11-2) were generated and maintained. Epstein-Barr virus-transformed  
683 lymphoblast cultures from three age- and gender-matched healthy controls were  
684 provided by the European Collection of Authenticated Cell Cultures (ECACC) (Salisbury,  
685 UK)

686 Paraffin-embedded and snap-frozen cerebellar (vermis) and frontal cortex from  
687 post-mortem brain from one sporadic CANVAS patient carrying the biallelic AAGGG  
688 repeat expansion (s16), three patients with genetically confirmed Friedreich's ataxia, one  
689 patient with genetically confirmed spinocerebellar ataxia 17, one patient with genetically  
690 confirmed *C9orf72*-related FTD and three neurologically healthy controls were obtained  
691 from the Queen Square Brain Bank for Neurological Disorders (London, UK).

692 Eight nerve biopsies and 10 muscle biopsies were obtained from patients carrying  
693 the homozygous AAGGG repeat expansion and healthy controls for pathological  
694 examination. Muscle biopsy tissue from six patients (Fam 6b-1, s1, s2, s18, s19, s22) and  
695 five controls was also used for qRT-PCR.

696 The study was approved by the UCL Institute of Neurology Institutional Review  
697 Board, and all subjects gave written informed consent to participate. The study has  
698 complied with all relevant ethical regulations.

699  
700 **SNP genotyping and linkage analysis.** Genotype calls were generated by the UCL  
701 genomics genotyping facility using InfiniumCoreExome arrays (Illumina). Raw data were  
702 processed and QC'ed using GenomeStudio (Illumina). All individual passed the 99% call  
703 rate threshold and were included in the subsequent analysis using PLINK 1.9 software<sup>48</sup>.  
704 Uninformative markers or markers with missing genotypes > 10% were removed, and the  
705 resulting dataset was further pruned to remove markers in high linkage equilibrium.  
706 Finally, the dataset was thinned to include 1-cM spaced markers covering all autosomes.  
707 In total, 3,476 markers were included. For fine-mapping analyses, all available  
708 informative markers were included.

709 Parametric linkage analysis was performed using MERLIN<sup>49</sup> assuming a highly  
710 penetrant recessive model of inheritance and disease allele frequency less than 1:10,000.  
711 MERLIN software was also used to obtain the most likely haplotypes in the candidate  
712 region. All genotyped individuals were included for haplotype analysis.

713 Single nucleotide polymorphisms rs11096992 and rs2066790 were genotyped in  
714 sporadic CANVAS patients and unaffected individuals by PCR followed by Sanger  
715 sequencing. Primers sequences, concentrations and PCR thermocycling conditions are  
716 provided in **Supplementary Table 3**.

717  
718 **Whole genome sequencing.** Whole genome sequencing was performed by deCODE  
719 genetics, Inc. Paired-end sequencing reads (100 bp) were generated using HiSeq4000  
720 (Illumina) and aligned to GRCH37 using Burrows-Wheeler Aligner<sup>50</sup>. The mean coverage  
721 per sample was 35x. Variants were called according to the GATK UnifiedGenotyper<sup>51</sup>  
722 workflow and annotated using ANNOVAR<sup>52</sup>. Variants were prioritized based on  
723 segregation, minor allele frequency (<0.0001 in the 1000 Genomes Project<sup>53</sup>, NHLBI GO  
724 Exome Sequencing project (Exome Variant Server, NHLBI GO Exome Sequencing Project  
725 (ESP), Seattle, WA (URL: <http://evs.gs.washington.edu/EVS/>) (September 2017), or  
726 gnomAD<sup>54</sup>, evolutionary conservation and in-silico prediction of pathogenicity for coding  
727 variants. Copy number analysis was performed using LUMPY<sup>55</sup> with default parameters.  
728 The candidate region on chromosome 4 was also visually inspected for any copy number  
729 or structural variants using IGV<sup>56</sup>.

730  
731 **Repeat-primed PCR.** Repeat-primed PCR was performed in order to provide qualitative  
732 assessment of the presence of an expanded AAGGG repeat as well expansions of the  
733 reference AAAAG allele or the AAAGG variant. The repeat-primed PCR was designed  
734 such that the reverse primers bind at different points within the repeat expansion to  
735 produce multiple amplicons of incremental size. 25 to 27 nucleotides flanking the repeat  
736 were added in order to increase binding affinity of the reverse primer to the polymorphic  
737 (A/AA/-) 3' end of the microsatellite and flanking region and give preferential  
738 amplification of the larger PCR product, thus allowing sizing of the expansion in some  
739 cases. Primers sequences, concentrations and PCR thermocycling conditions are provided  
740 in **Supplementary Table 3**.

741 Reverse primers were used in equimolar concentrations. Fragment length analysis  
742 was performed on an ABI 3730xl genetic analyzer (Applied Biosystems), and data were  
743 analyzed using GeneMapper software. Expansions with a characteristic "saw-tooth"  
744 pattern were identified and put forward for Southern blotting where sufficient DNA  
745 allowed.

746  
747 **Southern blotting.** Five µg of gDNA was digested for 3 h with EcoRI (10U) prior to  
748 electrophoresis. DNA was transferred to positively charged nylon membrane (Roche  
749 Applied Science) by capillary blotting and was crosslinked by exposure by ultraviolet  
750 light. Digoxigenin (DIG)-labelled probes were prepared by PCR amplification of a

751 genomic fragment cloned into a pGEM®-T Easy Vector using PCR DIG Probe Synthesis  
752 Kit (Roche Applied Science). Primer pairs used for cloning of gDNA fragment and PCR  
753 amplification of digoxigenin-labelled probe and PCR conditions are shown in  
754 **Supplementary Table 3**. Filter hybridization was undertaken as recommended in the  
755 DIG Application Manual (Roche Applied Science) except for the supplementation of DIG  
756 Easy Hyb buffer with 100 mg/ml denatured fragmented salmon sperm DNA. After  
757 prehybridization at 46 °C for 3 h, hybridization was allowed to proceed at 46 °C  
758 overnight. A total of 600 µl of PCR products containing the labelled oligonucleotide probe  
759 was used in 30 ml of hybridization solution. Membranes were washed initially in 23  
760 standard sodium citrate (SSC) and 0.1% sodium dodecyl sulfate (SDS), while the oven  
761 was being ramped from 48 °C to 65 °C and then washed three times in fresh solution at  
762 65 °C for 15 min. Detection of the hybridized probe DNA was carried out as  
763 recommended in the DIG Application Manual with CSPD ready-to-use (Roche Applied  
764 Science) as a chemiluminescent substrate. Signals were visualized on Fluorescent  
765 Detection Film (Roche Applied Science) after 1 h. All samples were electrophoresed  
766 against DIG-labelled DNA molecular-weight markers II and III (Roche Applied Science).  
767 Pentanucleotide repeat number was estimated after subtraction of the wild-type allele  
768 fragment size (5,037 bp). Sizes of the detected bands were recorded for each individual  
769 and number of expanded repeated unit was estimated using the formula repeated  
770 pentanucleotides unit = (size of the expanded band in bp – 5,000 bp)/5.

771  
772 **Neuropathological examination.** The formalin fixed cerebellar tissue was embedded in  
773 paraffin wax, from which 5-µm thick sections were cut for routine haematoxylin and  
774 eosin staining and immunohistochemistry. The sections were immunostained for p62  
775 (Abcam, ab56416, 1:500), TDP43 (Novus Biologicals, 2E2-D3, 1:500),  $\alpha$ -synuclein (Abcam,  
776 4D6, 1:1,000), phospho-Tau (AT-8, Innogenetics, 1:100) and anti  $\beta$ A4 (DAKO 6F3D, 1:50).  
777 Immunostaining, together with appropriate controls, was performed on a Roche Ventana  
778 Discovery automated staining platform following the manufacturer's guidelines, using  
779 biotinylated secondary antibodies and streptavidin-conjugated horseradish peroxidase  
780 and diaminobenzidine as the chromogen. Assessment of neuronal density in the  
781 cerebellar cortex was performed semi-quantitatively. Nerve and muscle biopsy  
782 specimens were performed and analysed according to standard procedures<sup>57,58</sup>. In brief,  
783 all nerve biopsies were examined after processing for paraffin histology (immunostaining  
784 for neurofilaments was performed with SMI31 antibody (Sternberger, 1:5,000) and in  
785 resin blocks (semithin resin sections were stained with methylene blue azure – basic  
786 fuchsin). The muscle biopsies were examined with routine histochemical stains after  
787 freezing in isopentane cooled in liquid nitrogen.

788  
789 **qRT-PCR.** Total RNA was extracted from fibroblasts, lymphoblasts and brain regions  
790 using 1 ml of Qiazol (Qiagen) and 200 µl chloroform. Aqueous phase was loaded and  
791 purified on columns using the RNeasy Lipid Tissue Mini kit (Qiagen) and treated with  
792 RNase-free DNase I (Qiagen). cDNA was synthesized using 500 ng of total RNA for all

793 samples, with a Superscript III first strand cDNA synthesis kit (Invitrogen) and an  
794 equimolar mixture of oligo dT and random hexamer primers. Real-time qRT-PCR was  
795 carried out using Power SYBR Green Master Mix (Applied Biosystems) and measured  
796 using a QuantStudio 7 Flex Real-Time PCR platform (Applied Biosystems).  
797 Glyceraldehyde 3-phosphate dehydrogenase (*GAPDH*) was used as housekeeping gene  
798 to normalize across different samples. Amplified transcripts were quantified using the  
799 comparative Ct method and presented as normalized fold expression change ( $2^{-\Delta\Delta Ct}$ ).  
800 Oligonucleotide sequences and thermocycling conditions are provided in  
801 **Supplementary Table 3.**

802  
803 **Western blotting.** Cells and tissues were lysed in radioimmunoprecipitation assay  
804 (RIPA) buffer supplemented with complete EDTA-free protease inhibitor cocktail  
805 (Roche). Brain lysates were homogenized on ice using a tissue ruptor with disposable  
806 probes (Qiagen). Protein lysate concentrations were measured by the BCA protein assay  
807 (Bio-Rad). After adding 5  $\mu$ l of sample buffer (Bio-Rad) and 2  $\mu$ l of NuPAGE reducing  
808 agent (Invitrogen) and boiling at 95 °C for 5 min, 15-30  $\mu$ g proteins for each sample were  
809 separated on 4-12% SDS-polyacrylamide gel (Bio-Rad) in MES buffer and transferred  
810 onto nitrocellulose membranes (GE-Healthcare) using a Turbo Transfer Pack (Bio-Rad).  
811 After blocking in 5% milk, immunoblotting was performed incubating over night at 4 °C  
812 with the following primary antibodies: anti-RFC1 (GTX129291, GeneTex 1:1,000), anti- $\beta$ -  
813 actin (A2228, Sigma, 1:2,000). Secondary antibodies were as follows: IRDye-800CW or  
814 IRDye-680CW conjugated goat anti-rabbit, donkey anti-mouse, IgG (Li-COR Bioscience).  
815 Signals of RFC1 bands were normalized to those of the corresponding  $\beta$ -actin bands as  
816 internal controls. Signals were digitally acquired by using an Odyssey Fc infrared scanner  
817 (Li-COR Bioscience) and quantified using Image Studio software (Li-COR Bioscience).

818  
819 **RNA-sequencing.** Reads were aligned to the hg38 human genome build using STAR  
820 (2.4.2a)<sup>59</sup>. BAM files were sorted, and duplicate reads flagged using NovoSort (1.03.09)  
821 (Novocraft). The aligned reads overlapping human exons (Ensembl 82) were counted  
822 using HTSeq (0.1)<sup>60</sup>. For each gene and each sample, the fragments per kilobase of exon  
823 per million mapped reads (FPKM) was calculated. Any gene with a mean FPKM across  
824 all samples in a dataset < 1 was discarded from further analysis. Differential gene  
825 expression was assessed with DESeq2 (1.8.2)<sup>61</sup> and differential splicing was assessed with  
826 DEXSeq<sup>62</sup>, running on R (3.3.2) (R project for statistical computing). The thresholds for  
827 significance for differential expression and splicing were set at a Benjamini-Hochberg  
828 false discovery rate of 10%. Quality control reports were collated using MultiQC<sup>63</sup>. Gene  
829 Ontology enrichment testing was done using g:Profiler<sup>64</sup> with both GO and KEGG  
830 ontologies, with minimum term size of 5 genes and all *P*-values Bonferroni corrected for  
831 multiple testing. Motif analysis was conducted on 49 alternatively spliced exons in  
832 lymphoblasts identified by unambiguous sequences with known strand using RBPmap<sup>65</sup>.  
833 Prediction of non-coding RNAs sequences in intron 2 of *RFC1* was tested by Rfam<sup>66</sup>.

834

835 **Statistical analyses.** Clinical variables were compared between familial and sporadic  
836 cases with two-tailed Student's *t* test (continuous variables) and Chi<sup>2</sup> (categorical  
837 variables). Correlation between repeat expansion size and age of onset of neuropathy was  
838 calculated using Pearson's correlation coefficient. FPKM of *FXN* and *RFC1* was compared  
839 using the two-tailed Student's *t* test. The relative expression of *RFC1* transcript 1 versus  
840 *GAPDH* as measured by qRT-PCR was compared with two-tailed Student's *t* test.  
841 Statistical analysis of the results of the western blot analysis was performed with two-  
842 tailed Student's *t* test after confirmation of equality of variances. *P* values of < 0.05 were  
843 considered to be significant.

844  
845 **Cloning of *RFC1* repeat expansion locus.** The *RFC1* locus containing the AAGGG repeat  
846 expansion was amplified by long-range PCR from genomic DNA from a CANVAS  
847 patient carrying the homozygous AAGGG repeat expansion and a healthy control  
848 carrying two (AAAAG)<sub>11</sub> alleles. PCR products were cloned into the pcDNA3.1/TOPO  
849 vector (Invitrogen) according to manufacturer's instructions. Primers and thermocycling  
850 conditions are provided in **Supplementary Table 3**. The size of the insert was determined  
851 by digestion with BstXI. Integrity of repeats and their orientation was confirmed by DNA  
852 sequencing (Eurofins Genomics), which revealed uninterrupted 94x (CCCTT) and 54x  
853 (AAGGG) repeats in mutant clones, as well as 11x (CTTTT) and 11x (AAAAG) repeat  
854 sequences in wild-type clone. Once confirmed, the four clones used for experimental  
855 procedures were amplified using a maxi-prep plasmid purification system.

856  
857 **RNA *in situ* hybridization.** Paraffin-embedded formalin-fixed post-mortem vermis  
858 sections from a CANVAS case, two healthy and two cerebellar degeneration age-  
859 matched controls were deparaffinized in xylene twice for 10 min, then rehydrated in  
860 100%, 90% and 70% ethanol, then in phosphate-buffered saline (PBS). About 10<sup>5</sup> SH-  
861 SY5Y cells were seeded on coverslips in 24-well plates and transfected using  
862 lipofectamine 3000 (Invitrogen) with plasmids expressing wild-type sense (TTTTTC)<sub>11</sub>,  
863 wild-type anti-sense (AAAAG)<sub>11</sub>, mutant sense (TTCCC)<sub>94</sub> or mutant anti-sense  
864 (AAGGG)<sub>54</sub> repeat sequences and were analyzed after 24 h. Cells were fixed in 4%  
865 methanol-free paraformaldehyde (Pierce) for 10 min at room temperature, dehydrated  
866 in a graded series of alcohols, air dried and rehydrated in PBS, permeabilized for 10  
867 min in 0.1% Triton X100 in PBS, briefly washed in 2× SSC and incubated for 30 min in  
868 pre-hybridisation solution (40% formamide, 2× SSC, 1 mg/ml tRNA, 1 mg/ml salmon  
869 sperm DNA, 0.2 %BSA, 10 % dextran sulphate, and 2 mM ribonucleoside vanadyl  
870 complex) at 57 °C. Hybridization solution (40% formamide, 2× SSC, 1 mg/ml tRNA,  
871 1 mg/ml salmon sperm DNA, 0.2% BSA, 10% dextran sulphate, 2 mM ribonucleoside  
872 vanadyl complex, 0.2 ng/μl (AAGGG)<sub>5</sub> or (CCCTT)<sub>5</sub> LNA probe, 5' TYE563-labeled  
873 (Exiqon), was heated at 95 °C for 10 min prior to incubation with sections for 1 h at  
874 57 °C. Cells were washed for 30 min at 57 °C with high-stringency buffer (2x SSC, 0.2%  
875 Triton X100, 40% formamide) and then for 20 min each, in 0.2x SSC buffer. Nuclei were  
876 stained by DAPI. Coverslips were then dehydrated in 70% then 100% ethanol and

877 mounted onto slides in Vectashield mounting medium. Images were acquired using an  
878 LSM710 confocal microscope (Zeiss) using a plan-apochromat 63x oil immersion  
879 objective.

880

881 **Response to DNA damage.** Fibroblasts were grown in 10-cm dishes in Dulbecco's  
882 modified Eagle's medium supplemented with 10% fetal bovine serum. Asynchronous  
883 cell cultures were grown to approximately 80% confluency and treated with UV,  
884 methyl methanesulfonate or untreated. For UV irradiation, cells were washed with  
885 PBS, and exposed to 30 or 120 J/m<sup>2</sup> UV light (254 nm) using a Stratalinker UV  
886 crosslinker®. For genotoxin treatment, methyl methanesulfonate (Sigma-Aldrich) was  
887 added to the culture media to give a final concentration of 1 mM, and cells were  
888 exposed for 8 h. After UV irradiation or genotoxin treatment, cells were allowed to  
889 recover for 24 h prior to analysis.

890 Cells were homogenized in RIPA Buffer containing 50 mM Tris pH 7.4, 150 mM  
891 NaCl, 1% Triton X-100, 0.5% Na deoxycholate, 0.1% SDS, 1 mM EDTA, and protease  
892 inhibitor. Samples were sonicated and centrifuged before protein levels were  
893 quantified using a BCA assay (Thermo Fisher Scientific Pierce). For western blot  
894 analysis, protein (5 µg) was size separated by SDS-PAGE, transferred to nitrocellulose  
895 membranes, and subjected to standard immunoblotting procedures using the  
896 following antibodies: γH2AX (Abcam; 1:1,000), β-Actin (Sigma-Aldrich; 1:1,000).  
897 γH2AX has been extensively used as a marker for DNA double strand breaks (DSBs).  
898 It is one of the initial markers of DSB being common to all DNA repair pathways.  
899 Secondary HRP-conjugated antibodies were purchased from PorteinTech and used at a  
900 1:2,000 concentration. Antibody staining was detected by ECL (Thermo Fisher  
901 Scientific Pierce) and visualized by X-ray film.

902 Cell viability was assessed using CellTiter-Glo® Luminescent Cell Viability  
903 Assay following manufacturers protocol. For cell-viability assessment, 20,000 cells/well  
904 were seeded in 96-well plates prior to treatment and treated as previously described.

905

906 **Life Sciences Reporting Summary.** Further information on experimental design is  
907 available in the **Life Sciences Reporting Summary**.

908

909 **Data availability.** The genotyping microarray data and sequence data obtained by  
910 whole-genome sequencing and RNA sequencing are available on request from the  
911 corresponding authors (A.C., H.H.). They are not publicly available because some of the  
912 study participants did not give full consent for releasing data publicly. Since whole-  
913 genome sequence data are protected by the Personal Information Protection Law,  
914 availability of these data is under the regulation by the institutional review board. The  
915 data obtained RNA sequencing have been deposited on SRA under accession number  
916 SUB5043763.

917

918

919 **METHODS-ONLY REFERENCES**

920

921

922 47. Manole A, Jaunmuktane Z, Hargreaves I, Ludtmann MHR, Salpietro V, Bello OD, et al.  
 923 Clinical, pathological and functional characterization of riboflavin-responsive neuropathy.  
 924 *Brain J Neurol.* 2017 01;140(11):2820–37.

925 48. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a tool  
 926 set for whole-genome association and population-based linkage analyses. *Am J Hum Genet.*  
 927 2007 Sep;81(3):559–75.

928 49. Abecasis GR, Cherny SS, Cookson WO, Cardon LR. Merlin--rapid analysis of dense genetic  
 929 maps using sparse gene flow trees. *Nat Genet.* 2002 Jan;30(1):97–101.

930 50. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
 931 *Bioinforma Oxf Engl.* 2009 Jul 15;25(14):1754–60.

932 51. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome  
 933 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing  
 934 data. *Genome Res.* 2010 Sep;20(9):1297–303.

935 52. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from  
 936 high-throughput sequencing data. *Nucleic Acids Res.* 2010 Sep;38(16):e164.

937 53. 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang  
 938 HM, et al. A global reference for human genetic variation. *Nature.* 2015 Oct 1;526(7571):68–  
 939 74.

940 54. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of  
 941 protein-coding genetic variation in 60,706 humans. *Nature.* 2016 18;536(7616):285–91.

942 55. Layer RM, Chiang C, Quinlan AR, Hall IM. LUMPY: a probabilistic framework for structural  
 943 variant discovery. *Genome Biol.* 2014 Jun 26;15(6):R84.

944 56. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al.  
 945 Integrative genomics viewer. *Nat Biotechnol.* 2011 Jan;29(1):24–6.

946 57. Weis J, Brandner S, Lammens M, Sommer C, Vallat J-M. Processing of nerve biopsies: a  
 947 practical guide for neuropathologists. *Clin Neuropathol.* 2012 Feb;31(1):7–23.

948 58. Dubowitz V, Sewry C, Oldfors A. *Muscle Biopsy—A Practical Approach*, 4th edn. Elsevier  
 949 Limited, Philadelphia. 4th ed. Philadelphia: Elsevier Limited; 2013.

950 59. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast  
 951 universal RNA-seq aligner. *Bioinforma Oxf Engl.* 2013 Jan 1;29(1):15–21.

952 60. Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-throughput  
 953 sequencing data. *Bioinforma Oxf Engl.* 2015 Jan 15;31(2):166–9.

954 61. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-  
 955 seq data with DESeq2. *Genome Biol.* 2014;15(12):550.

- 956 62. Anders S, Reyes A, Huber W. Detecting differential usage of exons from RNA-seq data.  
957 Genome Res. 2012 Oct;22(10):2008–17.
- 958 63. Ewels P, Magnusson M, Lundin S, Källner M. MultiQC: summarize analysis results for  
959 multiple tools and samples in a single report. Bioinforma Oxf Engl. 2016 01;32(19):3047–8.
- 960 64. Reimand J, Arak T, Adler P, Kolberg L, Reisberg S, Peterson H, et al. g:Profiler-a web server  
961 for functional interpretation of gene lists (2016 update). Nucleic Acids Res. 2016  
962 08;44(W1):W83-89.
- 963 65. Paz I, Kosti I, Ares M, Cline M, Mandel-Gutfreund Y. RBPmap: a web server for mapping  
964 binding sites of RNA-binding proteins. Nucleic Acids Res. 2014 Jul;42(Web Server  
965 issue):W361-367.
- 966 66. Griffiths-Jones S, Bateman A, Marshall M, Khanna A, Eddy SR. Rfam: an RNA family  
967 database. Nucleic Acids Res. 2003 Jan 1;31(1):439–41.
- 968 67. Podhorecka M, Skladanowski A, Bozko P. H2AX Phosphorylation: Its Role in DNA Damage  
969 Response and Cancer Therapy. J Nucleic Acids. 2010 Aug 3;2010.
- 970 68. Sharma A, Singh K, Almasan A. Histone H2AX phosphorylation: a marker for DNA damage.  
971 Methods Mol Biol Clifton NJ. 2012;920:613–26.

972

973

974 **Editorial summary:**

975

976 Biallelic expansion of an intronic AAGGG repeat in *RFC1* is identified here as a common  
977 cause of late-onset ataxia. This expansion occurs in the polyA tail of an AluSx3 element  
978 and is observed at a carrier frequency of 0.7% in populations of European ancestry.