

1 **Biallelic Expansion of an intronic Repeat in the *RFC1* Gene is a common cause**
2 **of Late-Onset Ataxia**

3
4 Andrea Cortese (1)*, Roberto Simone (2), Rosin Sullivan (1)[§], Jana Vandrovцова (1)[§], Huma
5 Tariq (1), Yau Way Yan (1), Jack Humphrey (1), Zane Jaunmuktane (2), Prasanth
6 Sivakumar (1), James Polke (3), Muhammad Ilyas (4), Eloise Tribollet (1), Pedro J Tomaselli
7 (5), Grazia Devigili (6), Ilaria Callegari (7), Maurizio Versino (7,8), Vincenzo Salpietro (1),
8 Stephanie Efthymiou (1), Diego Kaski (1), Nick W Wood (1), Nadja S Andrade (9), Elena
9 Buglo (10), Adriana Rebelo (10), Alexander M Rossor (1), Adolfo Bronstein (2), Pietro
10 Fratta (1), Wilson J Marques (5), Stephan Züchner (10), Mary M Reilly (1)[#], Henry Houlden
11 (1)^{*,#}

12
13 (1) Department of Neuromuscular Disease and (2) Department of Clinical and Movement
14 Neurosciences (3) Neurogenetics Laboratory, UCL Institute of Neurology and The National Hospital
15 for Neurology, Queen Square, London, WC1N 3BG, UK.

16 (4) Department of Biotechnology, Islamabad University, Islamabad and Punjab University, Lahore,
17 Pakistan.

18 (5) Department of Neurology, School of Medicine of Ribeirão Preto, University of São Paulo,
19 Ribeirão Preto, Brazil.

20 (6) UO Neurologia I, Fondazione IRCCS Istituto Neurologico “Carlo Besta”, Milano

21 (7) IRCCS Mondino Foundation, Pavia, Italy (8) Department of Brain and Behavioral Sciences,
22 University of Pavia, Pavia, Italy

23 (9) Department of Psychiatry and Behavioural Sciences, Center for Therapeutic Innovation, (10)
24 Dr. John T. Macdonald Foundation Department of Human Genetics and John P. Hussman Institute
25 for Human Genomics, University of Miami Miller School of Medicine, Miami, USA.

26 [§] equal contribution

27 [#] these authors jointly directed this project

28 ^{*} corresponding authors

29
30 Correspondence to

31 Andrea Cortese, MD PhD

32 UCL Institute of Neurology, Department of Neuromuscular Disease,

33 Queen Square, WC1N 3BG

34 London, United Kingdom

35 Email: andrea.cortese@ucl.ac.uk

36
37 Henry Houlden, MD PhD

38 UCL Institute of Neurology, Department of Neuromuscular Disease,

39 Queen Square, WC1N 3BG

40 London, United Kingdom

41 Email: h.houlden@ucl.ac.uk

43 **Abstract**

44 Late-onset ataxia is common, often idiopathic and can result from cerebellar,
45 proprioceptive or vestibular impairment, when in combination also termed cerebellar
46 ataxia, neuropathy, vestibular areflexia syndrome (CANVAS). We used non-parametric
47 linkage analysis and genome sequencing to identify a biallelic intronic AAGGG repeat
48 expansion in the replication factor C subunit-1 (*RFC1*) as the cause of familial CANVAS
49 and a frequent cause of late-onset ataxia, particularly if sensory neuropathy and bilateral
50 vestibular areflexia coexisted. The expansion, which occurs in the polyA tail of an AluSx3
51 element and differs in terms of both size and nucleotide sequence from the reference
52 (AAAAG)₁₁ allele, does not affect *RFC1* expression in patient peripheral and brain tissue
53 suggesting no overt loss-of-function. These data, along with the European expansion
54 carrier frequency of 0.7%, implies that biallelic AAGGG expansion in *RFC1* is a frequent
55 cause of late-onset ataxia.

56 INTRODUCTION

57 Late-onset ataxia, postural imbalance and falls are a frequent reason for neurological
58 consultation. Physiologically, motor coordination is achieved under visual control thanks
59 to the cerebellar integration of proprioceptive information conveyed by large-fibre sensory
60 neurons and vestibular inputs. Failure of any, or a combination of these systems can result
61 in ataxia (1–6). Both acquired and genetic causes are known but a large proportion remain
62 idiopathic.

63
64 Previous studies suggest that, there is a spectrum of clinical signs from pure
65 idiopathic late-onset cerebellar degeneration (ILOCA) through to the combined
66 degeneration of the cerebellum and its vestibular and sensory afferents, which has been
67 named cerebellar atrophy, neuropathy and vestibular areflexia syndrome (CANVAS)
68 (**Figure 1A**) (7). CANVAS is an adult-onset slowly progressive neurological disorder
69 characterized by imbalance, sensory neuropathy (neuronopathy), bilateral
70 vestibulopathy(8), chronic cough, and occasionally autonomic dysfunction. (9). Typically,
71 sensory action potentials and somatosensory potentials are absent throughout, brain MRI
72 shows cerebellar atrophy and vestibular testing is consistent with impaired vestibular
73 function bilaterally (10,9,11–17). Late-onset ataxia and CANVAS are usually sporadic, but
74 occasionally occur in siblings, raising the possibility of recessive transmission. However,
75 initial attempts to identify the underlying genetic defect by whole-exome sequencing were
76 unsuccessful.

77
78 Using non-parametric linkage analysis and whole-genome sequencing we identified
79 a recessive intronic AAGGG repeat expansion in replication factor C subunit 1 (*RFC1*) as a
80 the cause of familial CANVAS. The expansion occurs in the polyA tail of an AluSx3 element
81 and differs in terms of both size and nucleotide sequence from the reference (AAAAG)¹¹
82 allele. Screening of additional sporadic cases with late-onset ataxia confirmed the presence
83 of the mutated AAGGG repeat expansion in 22% of them, and in higher percentages in if
84 sensory neuronopathy and/or bilateral vestibular areflexia coexisted, suggesting that it
85 represents a frequent and certainly underrecognized cause of late-onset ataxia.

86

87

88 RESULTS

89 Genetic study

90 We genotyped 29 individuals, 23 affected and 6 unaffected, from 11 families (**Figure**
91 **1B**). The majority of the families consisted of affected sibships except two first-degree
92 cousins from non-consanguineous families (Fam 5b-2 and Fam 6b-1). None of the families
93 had convincing evidence of vertical disease transmission.

94

95 Assuming a recessive mode of inheritance, non-parametric linkage analysis
96 identified a single peak at position 4q14 with cumulative maximum HLOD of 5.8 (**Figure**
97 **2A**). Haplotype analysis defined a 1.7 MB region between markers rs6814637 and

98 rs10008483 (chr4:38977921-40712231) where within single families affected siblings shared
99 the same maternal and paternal alleles as opposed to unaffected brothers and sisters, who
100 had at most one of them (**Figure 2B**). The region contains 21 known HGNC genes
101 (**Supplementary Table 1**). Homozygosity mapping in consanguineous families showed
102 that the previously identified 1.7 MB region is encompassed in a larger run of
103 homozygosity of 12 MB shared by the affected siblings (**Supplementary Figure 1**). Of
104 interest, inside the 1.7Mb region four SNPs (rs2066790, rs11096992, rs17584703 and
105 rs6844176, bold highlighted), mapping inside a region encompassing all exons of
106 replication factor C subunit 1 (*RFC1*) and the last exon of WD repeat domain 19 (*WDR19*)
107 were shared by all affected individuals from different families except for individual Fam
108 5b-2, raising the possibility of a founder haplotype (**Figure 2C and 2D**).

109
110 Whole-exome sequencing was previously performed in seven individuals (Fam 1-1,
111 Fam 1-2, Fam 1-3, Fam 3-1, Fam 3-2, Fam 4-2, Fam 4-3) from three unrelated families (Fam1,
112 Fam3, Fam4), but did not identify recurrent non-synonymous variants within the coding
113 regions of the genes encompassed in the 1.7 Mb region (data not shown). We next
114 performed whole-genome sequencing (WGS) in an additional 6 affected individuals (Fam
115 2-1, Fam 8-2, Fam 8-3, Fam 6a-2, Fam 5a-2, Fam 7-1), 1 unaffected subject (Fam 8-1) from
116 four unrelated families and one sporadic case (s9). Analysis for non-synonymous variants
117 and copy number variants did not reveal changes recurring in the affected families. By
118 visually inspecting the aligned paired reads inside the 1.7 Mb region we noted in all
119 CANVAS patients a reduced read depth in a region encompassing a simple tandem
120 (AAAAG)₁₁ repeat at position chr4:39350045-39350103 (**Figure 3A**). Inside the
121 microsatellite region, the reference (AAAAG)₁₁ repeat was replaced in patients by a
122 variable number of AAGGG repeated units, which were detected on the reads mapped to
123 either side of the short tandem repeat. However, none of the reads could span across the
124 microsatellite region from one side to the other, suggesting the presence of a biallelic
125 expansion of the AAGGG repeat unit (**Figure 3B**). WGS from an unaffected sibling (Fam 8-
126 1) showed an equal distribution of interrupted reads containing the mutated AAGGG
127 repeated unit change as well as reads containing the AAAAG repeat.

128
129 We then performed repeat-primed PCR (RPPCR) with primers targeting the mutant
130 AAGGG pentanucleotide unit and confirmed the presence of an AAGGG repeat expansion
131 in all affected members from 11 families, as well as in unaffected carriers (**Figure 3C**).
132 Flanking PCR using standard conditions failed to amplify the region in all patients
133 suggesting the presence of a large expansion on both alleles, as opposed to their unaffected
134 siblings for whom at least one allele could be amplified by PCR (data not shown).

135
136 We next screened a cohort of 150 patients diagnosed with sporadic late onset ataxia
137 and we identified additional 33 (22%) sporadic cases carrying the recessive AAGGG repeat
138 expansion, as defined by a positive RPPCR for AAGGG repeat unit and the absence of PCR
139 amplifiable products by standard flanking PCR. The percentage of positive cases raised to

140 63% (32/51) if considering cases with late onset cerebellar ataxia and sensory neuropathy
 141 and to 92% (11/12) in cases with full CANVAS syndrome. Taking advantage of two
 142 informative single-nucleotide polymorphisms rs11096992 and rs2066790, by PCR and
 143 direct sequencing we observed that all additional sporadic cases but individual s23 shared
 144 the same haplotype as familial CANVAS cases.

145

146 By long-range PCR we were able to amplify and confirm by Sanger sequencing in
 147 all patients the presence of the AAGGG expansion (**Figure 3D**). However long-range PCR
 148 did not allow sizing of the repeat expansion as PCR is error-prone and contraction of
 149 repeated regions during PCR cycling have been previously demonstrated (18). Therefore,
 150 Southern blots were conducted in 34 cases and confirmed the presence of biallelic large
 151 expansions in all of them. Biallelic expansions could be visualized as two distinct bands in
 152 subjects carrying expansions of different sizes, or one thick band if the expanded alleles
 153 had a similar size (**Supplementary Figure 2**). Four unaffected siblings from four families
 154 were also included and they all carried one expanded and one normal allele. Although the
 155 expansion size could vary across different families, ranging from around 400 to 2000
 156 repeats, in the majority of cases approximately 1000 repeats were observed. Repeat size
 157 was relatively stable in siblings within single families. There was no association between
 158 age at onset and the number of AAGGG repeat units on either the smaller or larger allele
 159 ($n=34$; $r=-0.006$, $p=0.97$ and $r=-0.04$, $p=0.81$, respectively)

160

161 **Polymorphic conformations and allelic distribution of the short tandem repeat locus in** 162 **the normal population**

163 Recessive AAGGG expansion, as defined by the combination of positive RPPCR
 164 targeting the AAGGG repeat and the absence of a PCR amplifiable product on flanking
 165 PCR, were not observed in 304 healthy controls screened. RPPCR analysis targeting the
 166 AAGGG repeat showed that 0.7% (4 out of 608 chromosomes tested) carried an AAGGG
 167 expansion in heterozygous state. Southern blot analysis was performed in all of them and
 168 confirmed the presence of an expanded allele ranging from 200 to 880 repeats (mean $720 \pm$
 169 360) The chr4:39350045-39350103 locus, where the expansion resides, was shown to be
 170 highly polymorphic in the normal population and, besides the rare AAGGG expansion
 171 allele $(AAGGG)_{exp}$, three other conformations were observed: $(AAAAG)_{11}$, $(AAAAG)_{exp}$,
 172 $(AAAGG)_{exp}$ (**Figure 4A**). The $(AAAGG)_{exp}$ was often interrupted by XX By a combinatory
 173 approach of flanking PCR, RPPCR targeting one of the three possible nucleotide sequences,
 174 as well as Southern blot and Sanger sequencing in selected cases, we observed an allelic
 175 distribution of 75.5% ($n=459$) for the $(AAAAG)_{11}$ allele, 13.0% ($n=79$) for $(AAAAG)_{exp}$
 176 allele, 7.9% ($n=48$) for $(AAAGG)_{exp}$ allele, and, as per above, 0.7% ($n=4$) for the $(AAGGG)_{exp}$
 177 allele (**Figure 4B**). **Average size of $(AAAAG)_{exp}$ ranged from 40 to 400 repeats (mean $160 \pm$**
 178 **72) and $(AAAGG)_{exp}$ ranged from 40 to 880 (mean 236 ± 181) (**Figure 4C**).**

179

180 Eight healthy subjects had biallelic repeat expansions of a distinct repeated unit:
 181 $(AAAAG)_{exp}/(AAGGG)_{exp}$ in one case, $(AAAGG)_{exp}/(AAGGG)_{exp}$ in one case and

182 (AAAAG)_{exp} / (AAAGG)_{exp} in six cases. 22 cases likely had two expansions of the repeated
 183 AAAAG unit and nine of the repeated AAGGG unit, as defined by a positive RPPCR for
 184 the target repeat and two distinct bands on the southern blot, although we cannot exclude
 185 that one of the two alleles may be characterized by a distinct nucleotide sequence, which
 186 was not considered in the present study. Indeed, 9 additional subjects had no PCR
 187 amplifiable product on flanking PCR and were negative for RPPCR targeting the AAAAG,
 188 AAAGG, or AAGGG repeated units suggesting the potential existence of other possible
 189 allelic conformations in 3% (n=18) of tested chromosomes. Southern blot could not be
 190 performed because of insufficient amount of DNA in these cases.

191

192 The haplotype associated in most patients with the AAGGG repeat expansion has
 193 an allelic carrier frequency in 1000genome control population of 18%. Based on rs11096992
 194 and rs2066790 markers genotyping, the disease-associated haplotype rs2066790 (AA),
 195 rs11096992 (AA) was absent in recessive state from healthy individuals who carried two
 196 (AAAAG)₁₁ alleles, two (AAAAG)_{exp} alleles or a compound (AAAAG)₁₁ / (AAAAG)_{exp}
 197 genotype, but was observed in three out of nine carriers of two (AAAGG)_{exp} alleles and one
 198 healthy subject with (AAGGG)_{exp} / (AAAGG)_{exp} alleles, suggesting its possible association
 199 with both (AAGGG)_{exp} and (AAAGG)_{exp} configurations of the repeated unit, but not
 200 (AAAAG)₁₁ or (AAAAG)_{exp}.

201

202 **Clinical features of patients carrying the recessive AAGGG repeat expansion**

203 The clinical features of 56 cases carrying the recessive intronic AAGGG repeat
 204 expansion, including 23 familial and 33 sporadic cases, are summarized in **Table 1** and
 205 detailed in **Supplementary Table 2**. All cases were of Caucasian ancestry. Apart from a
 206 higher frequency of vestibular areflexia in familial CANVAS, clinical features were
 207 otherwise similar in familial and sporadic cases; hence data are presented together. Mean
 208 age of onset was 54 ± 9 (35-73) years and mean disease duration at examination 11 ± 7 (1-
 209 30) years. The most common complain at disease onset was unsteadiness, which was
 210 reported by 84% of patients, and frequently described as being worse in the dark. 37% of
 211 patients complained of chronic cough, which in some cases could precede by decades the
 212 onset of the walking difficulties. The neurologic examination invariably showed signs in
 213 keeping with a large fibre sensory neuropathy, 80% of patients had signs of cerebellar
 214 involvement, and overall 54% had evidence of bilateral vestibular areflexia. 23% of patients
 215 had concurrent autonomic nervous system involvement, particularly affecting micturition
 216 and defecation. Nerve conduction studies confirmed the presence of a non-length-
 217 dependent sensory neuropathy in all cases tested, as opposed to an entirely normal motor
 218 conduction study in most patients. Cerebellar atrophy was identified in 35 (83%) of 42 cases
 219 who underwent an MRI or CT scan.

220

221 **Neuropathological examination**

222 Pathological examination was conducted in a patient with CANVAS who carried
 223 the biallelic AAGGG repeat expansion and compared with a patient with genetically

224 confirmed Friedreich's ataxia, one patient with spinocerebellar ataxia 17 (SCA17) and one
225 case with *C9orf72*-related frontotemporal dementia (FTD), as well as control brains (**Figure**
226 **5**). The patient with CANVAS showed severe, widespread depletion of Purkinje cells with
227 associated prominent Bergmann gliosis, whilst cell density in the granule cell layer was
228 well preserved. Loss of Purkinje cells was also observed in Friedreich's ataxia, SCA17 and,
229 to a much lesser extent, in *C9orf72*-related FTD, but not in control brain. Similar to
230 Friedreich's ataxia and control brain, and as opposed to SCA17 and a *C9orf72*-related FTD
231 which were tested as positive controls, immunostaining for p62 showed no pathological
232 cytoplasmic or intranuclear inclusions in the cerebellar cortex of the patient with CANVAS.
233 Examination of the brain, in addition to prominent cerebellar atrophy, revealed age-related
234 changes in the form of neurofibrillary tangle tau pathology and amyloid- β pathology
235 (**Supplementary Figure 3**).

236

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

241

242 Fluorescence *in situ* hybridization using sense (AAGGG)₅ and anti-sense (TTCCC)₅
243 repeat specific oligonucleotides was performed on vermis post-mortem tissue from one
244 CANVAS patients, disease and healthy controls. As opposed to SH-SY5Y cells transfected
245 with pcDNA3.1/CT-GFP TOPO vector containing either (TTCCC)₉₄ or (AAGGG)₅₄ where
246 intranuclear and cytoplasmic inclusion were clearly detectable, we did not observe the
247 presence of endogenous RNA foci in any of the samples examined (**Supplementary Figure**
248 **5**).

249

250 RNA-sequencing

251 We performed whole transcriptome analysis in order to assess the presence of
252 changes in *RFC1* expression, as well as *in-cis* and *in-trans* effects at more distant genomic
253 regions. RNA-seq data showed that *RFC1* mRNA was unchanged in CANVAS (n=4) and
254 control (n=4) fibroblasts (P = 0.42) and in CANVAS (n=2) and control (n=3) lymphoblasts
255 (P = 0.45). We also performed RNA-seq from frontal cortex and cerebellar vermis from
256 autopsied brains from one CANVAS patient, Friedreich's ataxia cases (n=3) and controls
257 without evidence of neurological disease (n=3). In the single CANVAS patient, *RFC1*
258 appears to be unchanged in both cortex and cerebellum compared to the other samples
259 (**Figure 6A**). However, frataxin gene (*FXN*) was clearly down regulated in Friedreich's
260 ataxia frontal cortex and cerebellum compared to controls (cerebellum P = 0.007; log₂ fold
261 change = -1.2; frontal cortex P = 0.0003; log₂ fold change = -1.3) (**Figure 6A**). The single
262 CANVAS sample resembled the controls for *FXN* expression.

263

264 There were no differentially expressed genes between patient and control
265 fibroblasts, whereas 132 differentially expressed genes were identified between patient and

266 control lymphoblasts. Gene Ontology analysis showed enrichment for immune terms,
267 whose relevance to the disease will warrant further work. Notably, only eight
268 differentially expressed genes were located on chromosome 4 and were all separated by at
269 least 25Mb from the locus of the repeat expansion. Analysis of differentially expressed
270 genes in frontal cortex and vermis was not possible due to the limited numbers of
271 CANVAS samples (n=1).

272

273 Splicing analysis was performed in lymphoblasts. We identified 145 exons in 108
274 genes that had evidence of differential exon usage in CANVAS patients compared to
275 healthy controls. Motif analysis for the alternatively spliced exons showed enrichment of
276 motifs targeted by SRSF proteins, and in particular of SRSF3. *RFC1* did not show aberrant
277 splicing of its coding exons in mature mRNA. Also, no reads containing the AAGGG or
278 TTCCC repeated unit mapping to intron 2 of *RFC1* pre-mRNA transcript were detected
279 and no anti-sense or non-coding transcript was observed at *RFC1* locus in any of the tissues
280 examined. Gene ontology analysis of alternatively spliced genes found enrichment for
281 focal adhesion and non-specific cellular response terms. Lists of differentially expressed
282 genes and differentially expressed exons in lymphoblasts, their normalised count values in
283 brain samples and motif analysis for the alternatively spliced exons are provided in
284 **Supplementary Table 4.**

285

286 ***RFC1* expression in patients' tissue**

287 Quantitative reverse transcriptase PCR was performed using two sets of primers
288 (**Figure 6B**) and, concordantly with RNA-seq data, did not show any significant decrease
289 of *RFC1* mRNA (RefSeq NM_002913) level in patients' fibroblasts (n=5), lymphoblasts
290 (n=2), muscle (n=6), frontal cortex and cerebellar vermis (n=1) compared to healthy controls
291 or Friedreich's ataxia cases (**Figure 6C**). Exon 2 and 3 were correctly spliced in the mature
292 *RFC1* mRNA as shown by RNA-seq, qRT-PCR and sequencing. However, assessment of
293 pre-mRNA expression by qRT-PCR showed a consistent increase of intron 2 retention (IR)
294 in patients' lymphoblasts (n=2), muscle (n=6) (p=0.0077), cerebellar and frontal cortex (n=1)
295 compared to healthy controls (**Supplementary Figure 6**). The low level of *RFC1* expression
296 in fibroblasts prevented the assessment of pre-mRNA processing.

297 Western blot showed that *RFC1* protein (Uniprot P35251-1) was not decreased in
298 patients' fibroblasts (n=5), lymphoblasts (n=4) or brain (n=1) compared to healthy controls
299 or Friedreich's ataxia cases (**Figure 6D and Supplementary Figure 7**). Assessment of *RFC1*
300 protein expression in muscle could not be performed due to limited tissue availability.

301 Since *RFC1* play a key role in DNA damage recognition and recruitment of DNA
302 repair enzymes, we assessed whether patients' derived fibroblasts have an impaired
303 response to DNA damage. Patients' fibroblasts did not show an increased susceptibility to
304 DNA-damage and their treatment with double-stranded DNA break inducing agents, UV
305 and Methyl methanesulfonate, triggered a grossly normal response to DNA-damage
306 (**Supplementary Figure 8**).

307

308

309 **DISCUSSION**

310 We identified a recessive repeat expansion in intron 2 of *RFC1* as the cause of
311 CANVAS and late-onset ataxia. Twenty-three cases from 11 families and 33 sporadic cases
312 carried the biallelic AAGGG repeat expansion. Notably, out of 150 cases from a single
313 centre diagnosed with late-onset ataxia 22% resulted positive for the biallelic AAGGG
314 repeat expansion and the percentage was higher if only patients with sensory
315 neuropathy and cerebellar involvement (62%), CANVAS disease (92%) and familial
316 CANVAS disease (100%) were considered, highlighting that a higher diagnostic can be
317 achieved in cases with well-defined clinical features and positive family history. Not since
318 the discovery two decades ago of the most common genes causing ataxia (19–22) and
319 Charcot-Marie-Tooth (CMT) genes, (23–26) has a novel gene explained percentages above
320 10% of genetically undetermined cases (27,28).

321

322 We determined that the allelic carrier frequency of the AAGGG repeat expansion
323 in healthy controls was 0.7%, which is similar to the allelic carrier frequency of the GAA
324 expansion in *FXN* gene ranging from 0.9 to 1.6%, and which in the biallelic state causes
325 the most common recessive ataxia, Friedreich's ataxia. Together, this data suggests that
326 the recessive AAGGG expansion in *RFC1* may represent a frequent cause of late onset
327 ataxia in the general population, with an estimated prevalence at birth of the recessive
328 trait of ~1/20,000.

329

330 The expansion resides at the 3'-end of a deep intronic AluSx3 element and it
331 increases the polyA-tail size from 11 to over 400 repeated units, but also alters its sequence.
332 Of interest, expansions in terminal and mid A stretches of Alu elements have been
333 previously identified to cause Friedreich's ataxia (19), SCA37 (29), more recently benign
334 adult familial myoclonic epilepsy (BAFME) (30) and now CANVAS and late-onset ataxia.
335 Together, these observations suggest that variations and expansion of this highly
336 polymorphic regions of Alu elements represent a common mechanism underlying
337 different inherited neurological disorders. Notably, both SCA37 and BAFME are
338 characterized by expansion of a mutated repeated unit, ATTTC and TTTC A, respectively
339 (29,30). In this study, as well as in BAFME and SCA37, the presence in normal population
340 of large expansions of the reference repeated unit suggests that the nucleotide change
341 rather than the size of the expansion may be the driving pathogenic mechanism

342

343 Alu elements are repetitive elements about 300 base pairs long highly conserved
344 within primate genomes. The 3'-end of an Alu element has a longer A-rich region that plays
345 a critical role in its amplification mechanism (31). Active elements degrade rapidly on an
346 evolutionary time scale by A-tail shortening or heterogeneous base interruptions
347 accumulating in the A-tail, such as G insertions. We hypothesize that the mutation of the
348 AAGGG repeated unit occurred as part of the inactivation process by G interruption of the
349 polyA tail of the retrotransposon AluSx3. As known, repetitive DNA motives, particularly

350 G-rich regions, can form secondary or tertiary nucleotide structures such as hairpins,
351 parallel and antiparallel G-quadruplexes and, if transcribed, DNA-RNA hybrids also
352 known as R loops. These structures have been shown to increase the exposure of single-
353 stranded DNA to damaging environmental agents and can initiate repeat expansion and
354 perpetrate genomic instability across meiotic and mitotic divisions or after DNA damage
355 (32).

356
357 Since the same ancestral haplotype is shared by the majority of familial and positive
358 cases as well as some healthy carriers of two (AAAGG)_{exp} alleles, we speculate that
359 nucleotide change AAAAG to AAAGG or AAGGG may represent an ancestral founder
360 event, which was followed by the pathologic expansion of the repeated unit, whose size
361 seems to correlate positively with its GC content. However, the identification of two
362 patients (fam 5b-2 and s23) with a recessive AAGGG repeat expansion who share only one
363 allele of the common haplotype implies that repeat expansions of the mutated AAGGG
364 unit can occur also on a different genetic background. Interestingly, fam 5b-2 was also
365 found to carry the largest repeat expansion (10kb or 2,000 repeats) among the cohort of
366 patients tested.

367
368 In the majority of the patients the expansion encompassed 1,000 repeats, but as low
369 as 400 AAGGG repeats were shown to be sufficient to cause disease. The size of expanded
370 alleles was relatively stable in siblings within single families, but no parent of the affected
371 patients was available to assess whether this also applies across generations. We did not
372 observe a correlation between age of onset of the neuropathy and size of the repeat
373 expansion, although the disease course was very slowly progressive and initial symptoms
374 might have been neglected in some patients but reported by others.

375
376 So far, approximately 40 neurological or neuromuscular genetic disorders have been
377 associated with nucleotide repeat expansions. Two of them are known to be inherited in a
378 recessive mode, namely Friedreich's ataxia and myoclonic epilepsy type 1, and are both
379 associated loss-of-function of the repeat hosting gene (33–35).

380 A remarkable aspect of the recessive expansion described here is that our data does
381 not suggest a direct mechanism of loss of function for the *RFC1* gene. We did not observe
382 a reduced level of *RFC1* expression at either transcript or protein level in CANVAS
383 patients, although as a known loss of function control we were able to detect a significant
384 reduction of *FXN* transcript in post-mortem brain from patients with Friedreich's ataxia.
385 Also, RNA-seq data did not show a clear effect on the expression of neighbouring or distant
386 genes. We cannot exclude that the repeat expansion may cause more subtle tissue-specific
387 alterations of *RFC1* transcript and protein or alter the structural organization of the
388 chromatin.

389
390 *RFC1* encodes the large subunit of replication factor C, a five subunit DNA
391 polymerase accessory protein. It loads *PCNA* onto DNA and activates DNA polymerases

392 delta and epsilon to promote the coordinated synthesis of both strands during replication
393 or after DNA damage (36–38). It is interesting to note that mutations in many of the genes
394 involved in DNA repair have been already associated with degenerative neurological
395 disorders, including ataxia-telangiectasia, xeroderma pigmentosum, Cockayne syndrome
396 and ataxia oculomotor apraxia 1 and 2 (39). Interestingly, ataxia and neuropathy are
397 common clinical features to all of them suggesting a particular susceptibility of cerebellum
398 and peripheral nerves to DNA damage. However, our preliminary study did not show an
399 impaired response to DNA damage in patients' derived fibroblasts.

400

401 In fact, late-onset Mendelian disorders represent a unique interpretative challenge,
402 as risk variants may exert subtle effects rather than a clear loss of function of the mutated
403 gene that are compatible with normal developmental until adult or old age (40). To this
404 regard, although unusual in the context of a recessive mode of inheritance, other
405 mechanisms, including the production of toxic RNA containing the expanded repeat, and
406 the translation of a repeat-encoded polypeptide, should be considered (41). We did not
407 observe in patients brain the presence of RNA foci of either the sense or anti-sense repeated
408 unit. However, we were able to detect a consistent increase across different tissues of the
409 retention of intron 2 in *RFC1* pre-mRNA. Retention of the repeat-hosting intron was
410 recently identified as a common event associated with other disease-causing GC-rich
411 intronic expansions, such as in myotonic dystrophy type 2 and *C9orf72*-ALS/FTD but not
412 AT-rich repeat expansions such as in Friedreich's ataxia (42). Intron retention and
413 abnormal pre-mRNA processing bear potential effects on nuclear retention and
414 nucleocytoplasmic transport of the pre-mRNA, which, if efficiently exported to the
415 cytoplasm, would be accessible to the translational machinery.

416

417 Notwithstanding the enormous progress in Mendelian gene identification during
418 the last decade, up to 40% of patients with ataxia and inherited neuropathy remain
419 genetically undiagnosed and the percentage can rise up to 80-90% in particular subtypes,
420 such as late-onset ataxia (2,5,43) and hereditary sensory neuropathies (27,28). Our paper,
421 together with other studies from recent years (30,44–46), provides evidence that the
422 combined use of whole-genome sequencing and classical genetic investigations such as
423 linkage analysis, can provide a powerful tool to unravel a significant part of the missing
424 heritability hidden in non-coding regions of the human genome

425

426

427

428

429

430

431

432

433

434 **ACKNOWLEDGMENTS**

435 AC is funded by the inherited neuropathy consortium, which is a part of the NIH Rare Diseases
436 Clinical Research Network (RDCRN) (U54NS065712) and Wellcome Trust (204841/Z/16/Z).
437 AMR is funded by a Wellcome Trust Postdoctoral Fellowship for Clinicians (110043/Z/15/Z). HH
438 is also supported by Rosetrees Trust, Ataxia UK, The MSA Trust, Brain Research UK, MDUK, The
439 Muscular Dystrophy Association (MDA), Higher Education Commission (HEC) of Pakistan and
440 The Wellcome Trust (Synaptopathies Strategic Award). The INC (U54NS065712) is a part of the
441 NCATS Rare Diseases Clinical Research Network (RDCRN). RDCRN is an initiative of the Office
442 of Rare Diseases Research (ORDR), NCATS, funded through a collaboration between NCATS and
443 the NINDS. Stephan Zuchner is thankful to the National Institute of Health (4R01NS075764) for its
444 support. This research was also supported by the National Institute for Health Research University
445 College London Hospitals Biomedical Research Centre (BRC). Neuromuscular and brain tissue
446 samples were obtained from University College London Hospitals NHS Foundation Trust as part of
447 the UK Brain Archive Information Network (BRAIN UK) which is funded by the Medical
448 Research Council and Brain Tumour Research and the NIH funded NeuroBioBank.. We also thank
449 Francesca Launchbury from UCL IQPath laboratory for the technical assistance in histology slide
450 preparation.

451

452

453 **AUTHORS CONTRIBUTION**

454 AC designed the study, collected clinical data, performed the genetic analysis which led to the discovery
455 of the AAGGG repeat expansions, analysed the data, drafted the manuscript together with contributions
456 from JV, RS, RS, JH. RS, ASN, ET, EB, AR, YW, MI performed the investigation on *RFC1* expression.
457 JV performed the computational genetic analysis; SR and HT collected and analysed the genetic data in
458 healthy controls; PJT, WJM, AB, GD, IC, MV, DK, VS, SE, AMR contributed with collection of clinical
459 data and patients' samples. HJ, SP, PF performed the RNA-seq analysis; ZJ performed the pathological
460 investigation; RS, AMR, PF, JP contributed to the design of the study. SZ contributed to the design of
461 the study and analysed the data. HH, MMR designed the study, collected patients' clinical data and
462 biological samples and analysed the data. All authors revised the manuscript.

463

464

465 **COMPETING INTEREST STATEMENT**

466 The authors declare no competing financial interests.

467

468 **REFERENCES**

469

470

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

473 2. Muzaimi MB, Thomas J, Palmer-Smith S, Rosser L, Harper PS, Wiles CM, et al. Population
474 based study of late onset cerebellar ataxia in south east Wales. J Neurol Neurosurg Psychiatry.
475 2004 Aug;75(8):1129–34.

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

478 4. Strupp M, Feil K, Dieterich M, Brandt T. Bilateral vestibulopathy. Handb Clin Neurol.
479 2016;137:235–40.

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

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

485 7. Migliaccio AA, Halmagyi GM, McGarvie LA, Cremer PD. Cerebellar ataxia with bilateral
486 vestibulopathy: description of a syndrome and its characteristic clinical sign. Brain J Neurol.
487 2004 Feb;127(Pt 2):280–93.

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

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

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

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

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

503 13. Rust H, Peters N, Allum JHJ, Wagner B, Honegger F, Baumann T. VEMPs in a patient with
504 cerebellar ataxia, neuropathy and vestibular areflexia (CANVAS). J Neurol Sci. 2017 Jul
505 15;378:9–11.

- 506 14. Pelosi L, Leadbetter R, Mulroy E, Chancellor AM, Mossman S, Roxburgh R. Peripheral nerve
507 ultrasound in cerebellar ataxia neuropathy vestibular areflexia syndrome (CANVAS). *Muscle*
508 *Nerve*. 2017 Jul;56(1):160–2.
- 509 15. Pelosi L, Mulroy E, Leadbetter R, Kilfoyle D, Chancellor AM, Mossman S, et al. Peripheral
510 nerves are pathologically small in cerebellar ataxia neuropathy vestibular areflexia syndrome:
511 a controlled ultrasound study. *Eur J Neurol*. 2018 Apr;25(4):659–65.
- 512 16. Taki M, Nakamura T, Matsuura H, Hasegawa T, Sakaguchi H, Morita K, et al. Cerebellar
513 ataxia with neuropathy and vestibular areflexia syndrome (CANVAS). *Auris Nasus Larynx*.
514 2018 Aug;45(4):866–70.
- 515 17. Infante J, García A, Serrano-Cárdenas KM, González-Aguado R, Gazulla J, de Lucas EM, et
516 al. Cerebellar ataxia, neuropathy, vestibular areflexia syndrome (CANVAS) with chronic
517 cough and preserved muscle stretch reflexes: evidence for selective sparing of afferent Ia
518 fibres. *J Neurol*. 2018 Jun;265(6):1454–62.
- 519 18. Hommelsheim CM, Frantzeskakis L, Huang M, Ülker B. PCR amplification of repetitive
520 DNA: a limitation to genome editing technologies and many other applications. *Sci Rep*. 2014
521 May 23;4:5052.
- 522 19. Campuzano V, Montermini L, Moltò MD, Pianese L, Cossée M, Cavalcanti F, et al.
523 Friedreich's ataxia: autosomal recessive disease caused by an intronic GAA triplet repeat
524 expansion. *Science*. 1996 Mar 8;271(5254):1423–7.
- 525 20. Orr HT, Chung MY, Banfi S, Kwiatkowski TJ, Servadio A, Beaudet AL, et al. Expansion of
526 an unstable trinucleotide CAG repeat in spinocerebellar ataxia type 1. *Nat Genet*. 1993
527 Jul;4(3):221–6.
- 528 21. Pulst SM, Nechiporuk A, Nechiporuk T, Gispert S, Chen XN, Lopes-Cendes I, et al. Moderate
529 expansion of a normally biallelic trinucleotide repeat in spinocerebellar ataxia type 2. *Nat*
530 *Genet*. 1996 Nov;14(3):269–76.
- 531 22. Kawaguchi Y, Okamoto T, Taniwaki M, Aizawa M, Inoue M, Katayama S, et al. CAG
532 expansions in a novel gene for Machado-Joseph disease at chromosome 14q32.1. *Nat Genet*.
533 1994 Nov;8(3):221–8.
- 534 23. Lupski JR, de Oca-Luna RM, Slaugenhaupt S, Pentao L, Guzzetta V, Trask BJ, et al. DNA
535 duplication associated with Charcot-Marie-Tooth disease type 1A. *Cell*. 1991 Jul
536 26;66(2):219–32.
- 537 24. Hayasaka K, Himoro M, Sato W, Takada G, Uyemura K, Shimizu N, et al. Charcot-Marie-
538 Tooth neuropathy type 1B is associated with mutations of the myelin P0 gene. *Nat Genet*.
539 1993 Sep;5(1):31–4.
- 540 25. Bergoffen J, Scherer SS, Wang S, Scott MO, Bone LJ, Paul DL, et al. Connexin mutations in
541 X-linked Charcot-Marie-Tooth disease. *Science*. 1993 Dec 24;262(5142):2039–42.
- 542 26. Züchner S, Mersiyanova IV, Muglia M, Bissar-Tadmouri N, Rochelle J, Dadali EL, et al.
543 Mutations in the mitochondrial GTPase mitofusin 2 cause Charcot-Marie-Tooth neuropathy
544 type 2A. *Nat Genet*. 2004 May;36(5):449–51.

- 545 27. Fridman V, Bundy B, Reilly MM, Pareyson D, Bacon C, Burns J, et al. CMT subtypes and
546 disease burden in patients enrolled in the Inherited Neuropathies Consortium natural history
547 study: a cross-sectional analysis. *J Neurol Neurosurg Psychiatry*. 2015 Aug;86(8):873–8.
- 548 28. Murphy SM, Laura M, Fawcett K, Pandraud A, Liu Y-T, Davidson GL, et al. Charcot-Marie-
549 Tooth disease: frequency of genetic subtypes and guidelines for genetic testing. *J Neurol*
550 *Neurosurg Psychiatry*. 2012 Jul;83(7):706–10.
- 551 29. Seixas AI, Loureiro JR, Costa C, Ordóñez-Ugalde A, Marcelino H, Oliveira CL, et al. A
552 Pentanucleotide ATTTTC Repeat Insertion in the Non-coding Region of DAB1, Mapping to
553 SCA37, Causes Spinocerebellar Ataxia. *Am J Hum Genet*. 2017 Jul 6;101(1):87–103.
- 554 30. Ishiura H, Doi K, Mitsui J, Yoshimura J, Matsukawa MK, Fujiyama A, et al.
555 Expansions of intronic TTTCA and TTTTA repeats in benign adult familial myoclonic
556 epilepsy. *Nat Genet*. 2018 Apr;50(4):581–90.
- 557 31. Deininger P. Alu elements: know the SINEs. *Genome Biol*. 2011 Dec 28;12(12):236.
- 558 32. Haeusler AR, Donnelly CJ, Rothstein JD. The expanding biology of the C9orf72 nucleotide
559 repeat expansion in neurodegenerative disease. *Nat Rev Neurosci*. 2016;17(6):383–95.
- 560 33. Dürr A, Cossee M, Agid Y, Campuzano V, Mignard C, Penet C, et al. Clinical and genetic
561 abnormalities in patients with Friedreich’s ataxia. *N Engl J Med*. 1996 Oct 17;335(16):1169–
562 75.
- 563 34. Lazaropoulos M, Dong Y, Clark E, Greeley NR, Seyer LA, Brigatti KW, et al. Frataxin levels
564 in peripheral tissue in Friedreich ataxia. *Ann Clin Transl Neurol*. 2015 Aug;2(8):831–42.
- 565 35. Paulson H. Repeat expansion diseases. *Handb Clin Neurol*. 2018;147:105–23.
- 566 36. Majka J, Burgers PMJ. The PCNA-RFC families of DNA clamps and clamp loaders. *Prog*
567 *Nucleic Acid Res Mol Biol*. 2004;78:227–60.
- 568 37. Tomida J, Masuda Y, Hiroaki H, Ishikawa T, Song I, Tsurimoto T, et al. DNA damage-
569 induced ubiquitylation of RFC2 subunit of replication factor C complex. *J Biol Chem*. 2008
570 Apr 4;283(14):9071–9.
- 571 38. Overmeer RM, Gourdin AM, Giglia-Mari A, Kool H, Houtsmuller AB, Siegal G, et al.
572 Replication factor C recruits DNA polymerase delta to sites of nucleotide excision repair but is
573 not required for PCNA recruitment. *Mol Cell Biol*. 2010 Oct;30(20):4828–39.
- 574 39. McKinnon PJ. Maintaining genome stability in the nervous system. *Nat Neurosci*. 2013
575 Nov;16(11):1523–9.
- 576 40. Higuchi Y, Hashiguchi A, Yuan J, Yoshimura A, Mitsui J, Ishiura H, et al. Mutations in MME
577 cause an autosomal-recessive Charcot-Marie-Tooth disease type 2. *Ann Neurol*. 2016
578 Apr;79(4):659–72.
- 579 41. La Spada AR, Taylor JP. Repeat expansion disease: progress and puzzles in disease
580 pathogenesis. *Nat Rev Genet*. 2010 Apr;11(4):247–58.

- 581 42. Sznajder ŁJ, Thomas JD, Carrell EM, Reid T, McFarland KN, Cleary JD, et al. Intron
582 retention induced by microsatellite expansions as a disease biomarker. *Proc Natl Acad Sci U S*
583 *A*. 2018 17;115(16):4234–9.
- 584 43. Gebus O, Montaut S, Monga B, Wirth T, Cheraud C, Alves Do Rego C, et al. Deciphering the
585 causes of sporadic late-onset cerebellar ataxias: a prospective study with implications for
586 diagnostic work. *J Neurol*. 2017 Jun;264(6):1118–26.
- 587 44. DeJesus-Hernandez M, Mackenzie IR, Boeve BF, Boxer AL, Baker M, Rutherford NJ, et al.
588 Expanded GGGGCC hexanucleotide repeat in noncoding region of C9ORF72 causes
589 chromosome 9p-linked FTD and ALS. *Neuron*. 2011 Oct 20;72(2):245–56.
- 590 45. Renton AE, Majounie E, Waite A, Simón-Sánchez J, Rollinson S, Gibbs JR, et al. A
591 hexanucleotide repeat expansion in C9ORF72 is the cause of chromosome 9p21-linked ALS-
592 FTD. *Neuron*. 2011 Oct 20;72(2):257–68.
- 593 46. Aneichyk T, Hendriks WT, Yadav R, Shin D, Gao D, Vaine CA, et al. Dissecting the Causal
594 Mechanism of X-Linked Dystonia-Parkinsonism by Integrating Genome and Transcriptome
595 Assembly. *Cell*. 2018 Feb 22;172(5):897-909.e21.
- 596 47. Manole A, Jaunmuktane Z, Hargreaves I, Ludtmann MHR, Salpietro V, Bello OD, et al.
597 Clinical, pathological and functional characterization of riboflavin-responsive neuropathy.
598 *Brain J Neurol*. 2017 01;140(11):2820–37.
- 599 48. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a tool
600 set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*.
601 2007 Sep;81(3):559–75.
- 602 49. Abecasis GR, Cherny SS, Cookson WO, Cardon LR. Merlin--rapid analysis of dense genetic
603 maps using sparse gene flow trees. *Nat Genet*. 2002 Jan;30(1):97–101.
- 604 50. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
605 *Bioinforma Oxf Engl*. 2009 Jul 15;25(14):1754–60.
- 606 51. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome
607 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing
608 data. *Genome Res*. 2010 Sep;20(9):1297–303.
- 609 52. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from
610 high-throughput sequencing data. *Nucleic Acids Res*. 2010 Sep;38(16):e164.
- 611 53. 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang
612 HM, et al. A global reference for human genetic variation. *Nature*. 2015 Oct 1;526(7571):68–
613 74.
- 614 54. Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of
615 protein-coding genetic variation in 60,706 humans. *Nature*. 2016 18;536(7616):285–91.
- 616 55. Layer RM, Chiang C, Quinlan AR, Hall IM. LUMPY: a probabilistic framework for structural
617 variant discovery. *Genome Biol*. 2014 Jun 26;15(6):R84.
- 618 56. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al.
619 Integrative genomics viewer. *Nat Biotechnol*. 2011 Jan;29(1):24–6.

- 620 57. Weis J, Brandner S, Lammens M, Sommer C, Vallat J-M. Processing of nerve biopsies: a
621 practical guide for neuropathologists. *Clin Neuropathol*. 2012 Feb;31(1):7–23.
- 622 58. Dubowitz V, Sewry C, Oldfors A. *Muscle Biopsy—A Practical Approach*, 4th edn. Elsevier
623 Limited, Philadelphia. 4th ed. Philadelphia: Elsevier Limited; 2013.
- 624 59. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
625 universal RNA-seq aligner. *Bioinforma Oxf Engl*. 2013 Jan 1;29(1):15–21.
- 626 60. Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-throughput
627 sequencing data. *Bioinforma Oxf Engl*. 2015 Jan 15;31(2):166–9.
- 628 61. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-
629 seq data with DESeq2. *Genome Biol*. 2014;15(12):550.
- 630 62. Anders S, Reyes A, Huber W. Detecting differential usage of exons from RNA-seq data.
631 *Genome Res*. 2012 Oct;22(10):2008–17.
- 632 63. Ewels P, Magnusson M, Lundin S, Källner M. MultiQC: summarize analysis results for
633 multiple tools and samples in a single report. *Bioinforma Oxf Engl*. 2016 01;32(19):3047–8.
- 634 64. Reimand J, Arak T, Adler P, Kolberg L, Reisberg S, Peterson H, et al. g:Profiler—a web server
635 for functional interpretation of gene lists (2016 update). *Nucleic Acids Res*. 2016
636 08;44(W1):W83–89.
- 637 65. Paz I, Kosti I, Ares M, Cline M, Mandel-Gutfreund Y. RBPmap: a web server for mapping
638 binding sites of RNA-binding proteins. *Nucleic Acids Res*. 2014 Jul;42(Web Server
639 issue):W361–367.
- 640 66. Griffiths-Jones S, Bateman A, Marshall M, Khanna A, Eddy SR. Rfam: an RNA family
641 database. *Nucleic Acids Res*. 2003 Jan 1;31(1):439–41.
- 642 67. Podhorecka M, Skladanowski A, Bozko P. H2AX Phosphorylation: Its Role in DNA Damage
643 Response and Cancer Therapy. *J Nucleic Acids*. 2010 Aug 3;2010.
- 644 68. Sharma A, Singh K, Almasan A. Histone H2AX phosphorylation: a marker for DNA damage.
645 *Methods Mol Biol Clifton NJ*. 2012;920:613–26.

646

647

648 **FIGURES LEGENDS**

649

650 **Figure 1. a, Clinical spectrum of idiopathic late-onset ataxia from isolated cerebellar,**
 651 **vestibular and sensory variants to full-blown CANVAS.** ILOCA: idiopathic late-onset
 652 cerebellar ataxia. CANVAS: cerebellar ataxia, neuropathy, vestibular areflexia syndrome.
 653 **b, Pedigrees of CANVAS families.** Squares indicate males and circles females. Diagonal
 654 lines are used for deceased individuals. CANVAS patients are indicated with filled
 655 symbols. Black dots indicate genotyped individuals. Red dots indicate patients enrolled for
 656 whole-genome sequencing study.

657

658 **Figure 2. Identification of CANVAS locus. a,** Non-parametric multipoint linkage analysis
 659 identifies a unique locus associated with the disease in chromosomal region 4p14 with
 660 maximal HLOD score of 5.8. **b,** Schematic representation of shared haplotypes within
 661 single families. Light blue bars indicate a genomic region shared by affected siblings in a
 662 family and for which unaffected siblings are discordant. Two red dashed lines define a 1.7
 663 Mb region common to the different families. Single-nucleotide polymorphisms defining
 664 the haplotypes are represented on the top line. **c,** Fine-mapping inside the 1.7 Mb region
 665 identifies a recessive haplotype shared by all distinct families (green highlighted), except
 666 for individual fam 5b-2, who likely shares only one allele (light green highlighted). **d,**
 667 Schematic representation of the candidate 1.7 Mb region encompassing all 24 exons and
 668 flanking regions of Replication Factor Subunit C (*RFC1*) and the last exon and flanking
 669 intron of WD Repeat-Containing Protein 19 (*WDR19*).

670

671 **Fig. 3 | A recessive expansion of a mutated AAGGG repeated unit in intron 2 of *RFC1***
 672 **causes CANVAS and late-onset ataxia in familial and sporadic cases. a.** A reduced read
 673 depth of whole genome sequencing is observed in CANVAS patients (n=6) in a region
 674 corresponding to a short tandem AAAAG repeat in intron 2 of *RFC1*. STR: short tandem
 675 repeat. **b,** Visualization on IGV of reads aligned to the short repeat and flanking region
 676 show in patients (n=6) the presence of a mutated AAGGG repeat unit (representative
 677 image). Reads from both sides are interrupted and are unable to cover the entire length of
 678 the microsatellite region. Note that, as per IGV default setting, AAGGG repeated units not
 679 mapping to the (AAAAG)₁₁ reference sequence are soft-clipped and do not contribute to
 680 the coverage of the STR in Figure 3A, which is virtually absent. However, ≥20 reads
 681 containing the AAGGG repeated unit could be observed in each patient if soft-clipped
 682 reads are shown. **c,** Repeat-primed PCR (RPPCR) targeting the mutated AAGGG repeated
 683 unit. FAM-labelled PCR products are separated on an ABI3730 DNA Analyzer.
 684 Electropherograms are visualized on GENEMAPPER at 2,000 relative fluorescence units.
 685 Representative plots from a patient carrying the AAGGG repeat expansion and one non-
 686 carrier are shown. RPPCR experiments were repeated independently at least twice with
 687 similar results. **d,** Sanger sequencing of long-range PCR reactions confirms in patients the
 688 AAAAG to AAGGG nucleotide change of the repeated unit.

689

690 **Figure 4. Polymorphic configurations of the repeat expansion locus and allelic**
 691 **distribution in healthy controls**

692 **A.** Schematic representation of the repeat expansion locus in intron 2 of Replication factor
 693 C subunit 1 and its main allelic variants. **B.** Estimated allelic frequencies in 608
 694 chromosomes from 304 healthy controls. **C.** Average size and standard deviation of
 695 $(AAAGG)_{exp}$ and $(AAAGG)_{exp}$ expansions in healthy controls and $(AAAGG)_{exp}$ in CANVAS
 696 patients

697

698 **Figure 5. Pathology of cerebellar degeneration in a patient with CANVAS carrying the**
 699 **recessive AAGGG repeat expansion**

700 **(A-E)** Haematoxylin and Eosin (H&E) stained sections and **(A1-E1)** Sections
 701 immunostained for p62.

702 In a control brain **(A)**, age-matched for the patient with CANVAS syndrome, there is well
 703 preserved density of Purkinje cells (yellow arrow) and also granule cell layer is densely
 704 populated with small neurocytes (green asterisk). In CANVAS syndrome **(B)** there is
 705 severe, widespread depletion of Purkinje cells with associated prominent Bergmann gliosis
 706 (blue arrow), whilst cell density in the granule cell layer is well preserved. In a patient with
 707 genetically confirmed Friedreich's ataxia **(C)**, there is patchy depletion of Purkinje cells
 708 associated with Bergmann gliosis and unremarkable appearance of the granule cell layer.
 709 In a patient with genetically confirmed spinocerebellar ataxia 17 (SCA17) **(D)**, there is
 710 widespread Purkinje cell loss with only occasional Purkinje cells remaining; also in this
 711 patient granule cell layer is densely populated with small neurocytes. In a patient with
 712 frontotemporal dementia due to *C9orf72* expansion **(E)**, the Purkinje cell loss is patchy and
 713 granule cell layer is unremarkable. Immunostaining for p62 shows no pathological
 714 cytoplasmic or intranuclear inclusions in the cerebellar cortex in the control patient **(A1)**,
 715 the patient with CANVAS syndrome **(B1)** and also in the patient with Friedreich's ataxia
 716 **(C1)**. In SCA17 patient, there are scattered discrete intranuclear p62 immunoreactive
 717 inclusions in the small neurones within granule cell layer **(D1)**; high-power view of a
 718 representative intranuclear inclusion is demonstrated in the inset within **D1**). In the patient
 719 with *C9orf72* expansion, there are frequent characteristic perinuclear p62 positive
 720 inclusions in the granule cell layer **(E1)** and high-power view of a representative inclusion
 721 is shown in the inset within **E1**). Scale bar: 100 μ m in A-E, 30 μ m in A1-E1 and 5 μ m in insets
 722 in **D1** and **E1**.

723

724

725 **Figure 6. RFC1 expression is not affected by the AAGGG repeat expansion. A.** Plots
 726 showing expression levels of *RFC1* and *FXN* as Fragments Per Kilobase Million (FPKM) in
 727 controls (Ctrl), patients with Friedreich's ataxia (FRDA) and one CANVAS patient. **B.**
 728 Mapping on *RFC1* transcript 1 of the primers used for assessment by qRT-PCR of *RFC1*
 729 *mRNA* (cF1-cR1 and cF2-cR2) and pre-*mRNA* (cF1/iR1) expression. Blue arrows indicate
 730 primers mapping to exonic and intronic regions of canonical *RFC1* transcript. Primers
 731 spanning across exonic junctions are connected by dotted lines. A red triangle indicates the

732 site of the AAGGG repeat expansion. C. Expression levels of the canonical coding *RFC1*
733 mRNA as measured by qRT-PCR using two separate set of primers cF1-cR1 and cF2-cR2.
734 D. *RFC1* protein levels as measured by Western blotting using the polyclonal antibody
735 (ab193559) and normalized to beta-actin in fibroblasts (FBs), lymphoblasts (LBLs), and
736 post-mortem cerebellum (CBM) and frontal cortex (FCX) from patients with CANVAS
737 compared to healthy controls (Ctrl) and Friedreich's Ataxia (FRDA) cases. Bar graphs show
738 mean \pm SD

739

740

741 **Supplementary Table 4. Lists of differently expressed genes and exons from RNAseq**
742 **experiments.** Differential gene expression was assessed with DESeq2 (1.8.2)
743 (**LBLs_DESeq_hits**) and differential splicing was assessed with DEXSeq
744 (**LBLs_DEXSeq_hits**) running on R (3.3.2) (R project for statistical computing) between
745 CANVAS (n=2) and control (n=3) lymphoblasts. The thresholds for significance for
746 differential expression and splicing were set at a Benjamini-Hochberg false discovery rate
747 of 10%. Motif analysis was conducted on 49 alternatively spliced exons in lymphoblasts
748 identified by unambiguous sequences with known strand using RBPmap
749 (**LBLs_RBP_motif_hit_counts**). Normalised count values in brain samples for previously
750 identified differently expressed genes and exons in lymphoblasts are provided in
751 **Normcount_Brain_LBLs_DESeq_hits** and **Normcount_Brain_LBLs_DEXSeqhit** tables,
752 respectively. CANVAS cerebellar ataxia, neuropathy, vestibular areflexia syndrome, CBM
753 cerebellum, Ctrl control, FCX frontal cortex, FDR false discovery rate, FRDA Friedreich's
754 ataxia, LBLs lymphoblasts

755

756

757

758 TABLES

759

760 **Table 1. Clinical features of patients with familial or sporadic late-onset ataxia carrying**
761 **the recessive AAGGG repeat expansion in *RFC1*.**

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

762 cMAP compound motor action potential, CT computed tomography, LL lower limbs, MRI magnetic
763 resonance imaging, NS not significant, SAP sensory action potential, UL upper limbs.

764

765

766 **METHODS**767 **Patients**

768 For the initial linkage study, we enrolled 29 individuals, 23 affected and 6 unaffected,
769 from 11 families with a clinical diagnosis of CANVAS across four Centres: National
770 Hospital for Neurology and Neurosurgery (London, UK), C. Mondino National
771 Neurological Institute (Pavia, Italy), C. Besta Neurological Institute and Department of
772 Neurology, School of Medicine (Ribeirão Preto, Brazil).

773 An additional 150 patients with sporadic CANVAS or late onset ataxia (onset after
774 35 years of age) were identified from the neurogenetic database of the National Hospital
775 for Neurology and Neurosurgery (London, UK). For the experimental procedures,
776 patients' samples are generally referred to as CANVAS and no distinction between
777 samples from patients with full-blown CANVAS or other more limited variants of late-
778 onset ataxia is made. A skin biopsy was performed in five (fam 1-3, fam 2-2, fam 5a-2, fam
779 5b-2, fam 6b-1) genetically confirmed subjects and six age and gender matched controls.
780 Fibroblast cultures were maintained according to standard procedures (47). Epstein-Barr
781 virus-transformed lymphoblast cultures from two patients (fam 8-2, fam 11-2) were
782 generated and maintained. Epstein-Barr virus-transformed lymphoblast cultures from
783 three age and gender matched healthy controls were provided by the European Collection
784 of Authenticated Cell Cultures (ECACC) (Salisbury, UK)

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

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

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

798

799 **SNP genotyping and linkage analysis**

800 Genotype calls were generated by the UCL genomics genotyping facility using
801 InfiniumCoreExome arrays (Illumina, Carlsbad, CA, USA). Raw data were processed and
802 QC'ed using GenomeStudio (Illumina, Carlsbad, CA, USA). All individual passed the 99%
803 call rate threshold and were included in the subsequent analysis using PLINK 1.9 software
804 (48). Uninformative markers or markers with missing genotypes > 10% were removed and
805 the resulting dataset was further pruned to remove markers in high linkage equilibrium.
806 Finally, the dataset was thinned to include 1cM spaced markers covering all autosomes. In

807 total 3476 markers were included. For fine-mapping analyses all available informative
808 markers were included.

809 Parametric linkage analysis was performed using MERLIN (49) assuming a highly
810 penetrant recessive model of inheritance and disease allele frequency less than 1:10,000.
811 MERLIN software was also used to obtain the most likely haplotypes in the candidate
812 region. All genotyped individuals were included for haplotype analysis.

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

817

818 **Whole Genome Sequencing**

819 Whole Genome Sequencing was performed by deCODE genetics, Inc. (deCODE
820 genetics, Reykjavik, Iceland). Paired-end sequencing reads (100bp) were generated using
821 HiSeq4000 (Illumina, San Diego, CA, USA) and aligned to GRCH37 using Burrows-
822 Wheeler Aligner (50). The mean coverage per sample was 35x. Variants were called
823 according to the GATK UnifiedGenotyper (51) workflow and annotated using ANNOVAR
824 (52). Variants were prioritised based on segregation, minor allele frequency (<0.0001 in
825 the 1000 Genomes Project (53), NHLBI GO Exome Sequencing project (Exome Variant
826 Server, NHLBI GO Exome Sequencing Project (ESP), Seattle, WA (URL:
827 <http://evs.gs.washington.edu/EVS/>) [September 2017], or gnomAD (54), evolutionary
828 conservation and in-silico prediction of pathogenicity for coding variants. Copy number
829 analysis was performed using LUMPY (55) with default parameters. The candidate region
830 on chromosome 4 was also visually inspected for any copy number or structural variants
831 using IGV (56).

832

833 **Repeat-primed PCR**

834 Repeat-primed PCR was performed in order to provide qualitative assessment of
835 the presence of an expanded AAGGG repeat as well expansions of the reference AAAAG
836 allele or the AAAGG variant. The repeat-primed PCR was designed such that the reverse
837 primers binds at different points within the repeat expansion to produce multiple
838 amplicons of incremental size. 25 to 27 nucleotides flanking the repeat were added in order
839 to increase binding affinity of the reverse primer to the polymorphic (A/AA/-) 3' end of the
840 microsatellite and flanking region and give preferential amplification of the larger PCR
841 product, thus allowing sizing of the expansion in some cases. Primers sequences,
842 concentrations and PCR thermocycling conditions are provided in **supplementary table 3**

843 Reverse primers were used in equimolar concentrations. Fragment length analysis
844 was performed on an ABI 3730xl genetic analyser (Applied Biosystems, Foster City, CA,
845 USA), and data were analysed using GeneMapper software. Expansions with a
846 characteristic "saw-tooth" pattern were identified and put forward for Southern blotting
847 where sufficient DNA allowed.

848

849 Southern Blotting

850 Five µg of gDNA was digested for three hours with EcoRI (10U) prior to
851 electrophoresis. DNA was transferred to positively charged nylon membrane (Roche
852 Applied Science) by capillary blotting and was crosslinked by exposure by ultraviolet light.
853 Digoxigenin (DIG)-labelled probes were prepared by PCR amplification of a genomic
854 fragment cloned into a pGEM®-T Easy Vector using PCR DIG Probe Synthesis Kit (Roche
855 Applied Science). Primer pairs used for cloning of gDNA fragment and PCR amplification
856 of digoxigenin-labelled probe and PCR conditions are shown in **Supplementary Table 3**.
857 Filter hybridization was undertaken as recommended in the DIG Application Manual
858 (Roche Applied Science) except for the supplementation of DIG Easy Hyb buffer with 100
859 mg/ml denatured fragmented salmon sperm DNA. After prehybridization at 46°C for three
860 hours, hybridization was allowed to proceed at 46°C overnight. A total of 600 µl of PCR
861 products containing the labelled oligonucleotide probe was used in 30 ml of hybridization
862 solution. Membranes were washed initially in 23 standard sodium citrate (SSC) and 0.1%
863 sodium dodecyl sulfate (SDS), while the oven was being ramped from 48°C to 65°C and
864 then washed three times in fresh solution at 65°C for 15 min. Detection of the hybridized
865 probe DNA was carried out as recommended in the DIG Application Manual with CSPD
866 ready-to-use (Roche Applied Science) as a chemiluminescent substrate. Signals were
867 visualized on Fluorescent Detection Film (Roche Applied Science) after 1 hr. All samples
868 were electrophoresed against DIG-labelled DNA molecular-weight markers II and III
869 (Roche Applied Science). Pentanucleotide repeat number was estimated after subtraction
870 of the wild-type allele fragment size (5,037 bp). Sizes of the detected bands were recorded
871 for each individual and number of expanded repeated unit was estimated using the
872 formula repeated pentanucleotides unit = (size of the expanded band in bp – 5000 bp)/5.

873

874 Neuropathological examination

875 The formalin fixed cerebellar tissue was embedded in paraffin wax, from which 5µm
876 thick sections were cut for routine haematoxylin and eosin staining and
877 immunohistochemistry. The sections were immunostained for p62 (Abcam, ab56416,
878 1:500), TDP43 (Novus Biologicals, 2E2-D3, 1:500), α-synuclein (Abcam, 4D6, 1:1000),
879 phospho-Tau (AT-8, Innogenetics, 1:100) and anti βA4 (DAKO 6F3D, 1:50).
880 Immunostaining, together with appropriate controls, was performed on a Roche Ventana
881 Discovery automated staining platform following the manufacturer's guidelines, using
882 biotinylated secondary antibodies and streptavidin-conjugated horseradish peroxidase
883 and diaminobenzidine as the chromogen. Assessment of neuronal density in the cerebellar
884 cortex was performed semi-quantitatively. Nerve and muscle biopsy specimens were
885 performed and analysed according to standard procedures (57,58). In brief, all nerve
886 biopsies were examined after processing for paraffin histology (immunostaining for
887 neurofilaments was performed with SMI31 antibody (Sternberger, 1:5000) and in resin
888 blocks (semithin resin sections were stained with methylene blue azure – basic fuchsin).
889 The muscle biopsies were examined with routine histochemical stains after freezing in
890 isopentane cooled in liquid nitrogen.

891

892 qRT-PCR

893

894 Total RNA was extracted from fibroblasts, lymphoblasts and brain regions using 1
895 ml of Qiazol (Qiagen) and 200 μ l chloroform. Aqueous phase was loaded and purified on
896 columns using the RNeasy Lipid Tissue Mini kit (Qiagen) and treated with RNase-free
897 DNase I (Qiagen). cDNA was synthesized using 500 ng of total RNA for all samples, with
898 a Superscript III first strand cDNA synthesis kit (Invitrogen) and an equimolar mixture of
899 oligo dT and random hexamer primers. Real-time qRT-PCR was carried out using Power
900 SYBR Green Master Mix (Applied Biosystems) and measured using a QuantStudio 7 Flex
901 Real-Time PCR platform (Applied Biosystems). Glyceraldehyde 3-phosphate
902 dehydrogenase (GAPDH) was used as housekeeping gene to normalize across different
903 samples. Amplified transcripts were quantified using the comparative Ct method and
904 presented as normalized fold expression change ($2^{-\Delta\Delta C_t}$). Oligonucleotide sequences and
905 thermocycling conditions are provided in **Supplementary Table 3**.

905

906 Western blotting

907

908 Cells and tissues were lysed in radioimmunoprecipitation assay (RIPA) buffer
909 supplemented with complete EDTA-free protease inhibitor cocktail (Roche). Brain lysates
910 were homogenised on ice using a tissue ruptor with disposable probes (Qiagen). Protein
911 lysate concentrations were measured by the BCA protein assay (Bio-Rad). After adding 5 μ l
912 of sample buffer (Bio-Rad) and 2 μ l of NuPAGE reducing agent (Invitrogen) and boiling at
913 95 °C for 5 min, 15-30 μ g proteins for each sample were separated on 4-12% SDS-
914 polyacrylamide gel (Bio-Rad) in MES buffer and transferred onto nitrocellulose
915 membranes (GE-Healthcare) using a Turbo Transfer Pack (Bio-Rad). After blocking in 5%
916 milk, immunoblotting was performed incubating over night at 4°C with the following
917 primary antibodies: anti-RFC1 (GTX129291, GeneTex 1:1000), anti- β -actin (A2228, Sigma,
918 1:2000). Secondary antibodies were as follows: IRDye-800CW or IRDye-680CW conjugated
919 goat anti-rabbit, donkey anti-mouse, IgG (Li-COR Bioscience). Signals of *RFC1* bands were
920 normalized to those of the corresponding β -actin bands as internal controls. Signals were
921 digitally acquired by using an Odyssey Fc infrared scanner (Li-COR Bioscience) and
922 quantified using Image Studio software (Li-COR Bioscience).

922

923 RNA-sequencing

924

925 Reads were aligned to the hg38 human genome build using STAR (2.4.2a) (59). BAM
926 files were sorted, and duplicate reads flagged using NovoSort (1.03.09) (Novocraft). The
927 aligned reads overlapping human exons (Ensembl 82) were counted using HTSeq (0.1)
928 (60). For each gene and each sample, the fragments per kilobase of exon per million
929 mapped reads (FPKM) was calculated. Any gene with a mean FPKM across all samples in
930 a dataset < 1 was discarded from further analysis. Differential gene expression was
931 assessed with DESeq2 (1.8.2) (61) and differential splicing was assessed with DEXSeq (62),
932 running on R (3.3.2) (R project for statistical computing). The thresholds for significance
for differential expression and splicing were set at a Benjamini-Hochberg false discovery

933 rate of 10%. Quality control reports were collated using MultiQC (63). Gene ontology
934 enrichment testing was done using g:Profiler (64) with both GO and KEGG ontologies,
935 with minimum term size of 5 genes and all p-values Bonferroni corrected for multiple
936 testing. Motif analysis was conducted on 49 alternatively spliced exons in lymphoblasts
937 identified by unambiguous sequences with known strand using RBPmap (65). Prediction
938 of non-coding RNAs sequences in intron 2 of *RFC1* was tested by Rfam (66)

939

940 **Statistical analyses**

941 Clinical variables were compared between familial and sporadic cases with two-
942 tailed Student's t test (continuous variables) and Chi² (categorical variables). Correlation
943 between repeat expansion size and age of onset of neuropathy was calculated using
944 Pearson's correlation coefficient. FPKM of *FXN* and *RFC1* was compared using the two-
945 tailed Student's t test. The relative expression of *RFC1* transcript 1 versus *GAPDH* as
946 measured by qRT-PCR was compared with two-tailed Student's t test. Statistical analysis
947 of the results of the western blot analysis was performed with two-tailed Student's t test
948 after confirmation of equality of variances. P values of < 0.05 were considered to be
949 significant.

950

951

952 **Cloning of *RFC1* repeat expansion locus**

953 The *RFC1* locus containing the AAGGG repeat expansion was amplified by long-range
954 PCR from genomic DNA from a CANVAS patient carrying the homozygous AAGGG
955 repeat expansion and a healthy control carrying two (AAAAG)₁₁ alleles. PCR products
956 were cloned into the pcDNA3.1/TOPO vector (Invitrogen) according to manufacturer's
957 instructions. Primers and thermocycling conditions are provided in **Supplementary**
958 **Table 3**. The size of the insert was determined by digestion with BstXI. Integrity of
959 repeats and their orientation was confirmed by DNA sequencing (Eurofins Genomics,
960 Louisville, KY, USA), which revealed uninterrupted 94x (CCCTT) and 54x (AAGGG)
961 repeats in mutant clones, as well as 11x (CTTTT) and 11x (AAAAG) repeat sequences in
962 wild-type clone. Once confirmed, the four clones used for experimental procedures were
963 amplified using a maxi-prep plasmid purification system.

964

965 **RNA *in situ* hybridization**

966 Paraffin-embedded formalin-fixed post-mortem Vermis sections from a CANVAS case,
967 2 healthy and 2 cerebellar degeneration age-matched controls were deparaffinized in
968 xylene twice for 10min, then rehydrated in 100%, 90% and 70% ethanol, then in
969 phosphate-buffered saline (PBS). About 10⁵ SH-SY5Y cells were seeded on coverslips in
970 24-well plates and transfected using lipofectamine 3000 (Invitrogen) with plasmids
971 expressing wild-type sense (TTTTTC)₁₁, wild-type anti-sense (AAAAG)₁₁, mutant sense
972 (TTCCC)₉₄ or mutant anti-sense (AAGGG)₅₄ repeat sequences and were analyzed after
973 24 hours. Cells were fixed in 4% methanol-free paraformaldehyde (Pierce) for 10 min at
974 room temperature, dehydrated in a graded series of alcohols, air dried and rehydrated

975 in PBS, permeabilised for 10 min in 0.1% Triton X100 in PBS, briefly washed in 2×SSC
976 and incubated for 30 min in pre-hybridisation solution (40 % formamide, 2×SSC,
977 1 mg/ml tRNA, 1 mg/ml salmon sperm DNA, 0.2 % BSA, 10 % dextran sulphate, 2 mM
978 ribonucleoside vanadyl complex) at 57 °C. Hybridisation solution (40 % formamide,
979 2×SSC, 1mg/ml tRNA, 1 mg/ml salmon sperm DNA, 0.2 % BSA, 10 % dextran sulphate,
980 2 mM ribonucleoside vanadyl complex, 0.2 ng/μl (AAGGG)₅ or (CCCTT)₅ LNA probe,
981 5' TYE563-labeled (Exiqon), was heated at 95 °C for 10 min prior to incubation with
982 sections for 1 h at 57 °C. Cells were washed for 30 min at 57 °C with high-stringency
983 buffer (2x SSC, 0.2% Triton X100, 40% formamide) and then for 20 min each, in 0.2x SSC
984 buffer. Nuclei were stained by DAPI. Coverslips were then dehydrated in 70% then
985 100% EtOH and mounted onto slides in Vectashield mounting medium. Images were
986 acquired using an LSM710 confocal microscope (Zeiss) using a plan-apochromat 63x oil
987 immersion objective.

988

989 **Response to DNA damage**

990 Fibroblasts were grown in 10 cm dishes in Dulbecco's modified Eagle's medium
991 supplemented with 10% fetal bovine serum. Asynchronous cell cultures were grown to
992 approximately 80% confluency and treated with UV, methyl methanesulfonate or
993 untreated. For UV irradiation, cells were washed with PBS, and exposed to 30 or 120
994 J/m² UV light (254 nm) using a Stratalinker UV crosslinker®. For genotoxin treatment,
995 methyl methanesulfonate (Sigma-Aldrich, St. Louis, MO, USA) was added to the culture
996 media to give a final concentration of 1mM and cells were exposed for 8 hours. After UV
997 irradiation or genotoxin treatment cells were allowed to recover for 24 hours prior to
998 analysis.

999 Cells were homogenized in RIPA Buffer containing 50 mM Tris pH 7.4, 150 mM NaCl,
1000 1% Triton X-100, 0.5% Na deoxycholate, 0.1% SDS, 1 mM EDTA, and protease inhibitor.
1001 Samples were sonicated and centrifuged before protein levels were quantified using a
1002 BCA assay (Thermo Fisher Scientific Pierce, Rockford, IL, USA). For Western blot
1003 analysis, protein (5 μg) was size separated by SDS-PAGE, transferred to nitrocellulose
1004 membranes, and subjected to standard immunoblotting procedures using the following
1005 antibodies: γH2AX (Abcam, USA; 1:1000), β-Actin (Sigma-Aldrich, St. Louis, MO, USA;
1006 1:1000). γH2AX has been extensively used as a marker for DNA double strand breaks
1007 (DSBs). It is one of the initial markers of DSB being common to all DNA repair pathways.
1008 Secondary HRP-conjugated antibodies were purchased from PorteinTech and used at a
1009 1:2000 concentration. Antibody staining was detected by ECL (Thermo Fisher Scientific
1010 Pierce, Rockford, IL, USA) and visualized by X-ray film.

1011 Cell viability was assessed using CellTiter-Glo® Luminescent Cell Viability Assay
1012 following manufacturers protocol. For cell-viability assessment, 20,000cells/well were
1013 seeded in 96-well plates prior to treatment and treated as previously described.

1014

1015 **Life Sciences Reporting Summary.** Further information on experimental design is
1016 available in the Life Sciences Reporting Summary.

1017

1018 **Data availability**

1019 The genotyping microarray data and sequence data obtained by whole-genome
1020 sequencing and RNA sequencing are available on request from the corresponding
1021 authors (A.C.; H.H.). They are not publicly available because some of the study
1022 participants did not give full consent for releasing data publicly. Since whole-genome
1023 sequence data are protected by the Personal Information Protection Law, availability of
1024 these data is under the regulation by the institutional review board. The data obtained
1025 RNA sequencing have been deposited on SRA under accession number SUB5043763.

1026