

1 **Validation of new bioinformatic tools to identify expanded repeats: a non-reference**
2 **intronic pentamer expansion in *RFC1* causes CANVAS**

3
4 Haloom Rafehi^{1,2,33}, David J Szmulewicz^{3,4,33}, Mark F Bennett^{1,2,5}, Nara LM Sobreira⁶, Kate
5 Pope⁷, Katherine R Smith¹, Greta Gillies⁷, Peter Diakumis⁸, Egor Dolzhenko⁹, Michael A
6 Eberle⁹, María García Barcina¹⁰, David P Breen^{11,12,13}, Andrew M Chancellor¹⁴, Phillip D
7 Cremer^{15,16}, Martin B. Delatycki^{7,17}, Brent L Fogel¹⁸, Anna Hackett^{19,20}, G. Michael
8 Halmagyi^{21,22}, Solange Kapetanovic²³, Anthony Lang^{24,25}, Stuart Mossman²⁶, Weiyi Mu⁶,
9 Peter Patrikios²⁷, Susan L Perlman²⁸, Ian Rosemargy²⁹, Elsdon Storey³⁰, Shaun RD Watson³¹,
10 Michael A Wilson⁷, David Zee³², David Valle⁶, David J Amor^{7,17}, Melanie Bahlo^{1,2,34} and
11 Paul J Lockhart^{7,17,34*}

12
13 ¹ Population Health and Immunity Division, The Walter and Eliza Hall Institute of Medical
14 Research, 1G Royal Parade, Parkville, Victoria 3052, Australia

15 ² Department of Medical Biology, University of Melbourne, 1G Royal Parade, Parkville,
16 Victoria 3052, Australia

17 ³ Cerebellar Ataxia Clinic, Neuroscience Department, Alfred Health, Melbourne, Victoria
18 3004, Australia

19 ⁴ Balance Disorders and Ataxia Service, Royal Victorian Eye & Ear Hospital, East
20 Melbourne, Victoria 3002, Australia.

21 ⁵ Epilepsy Research Centre, Department of Medicine, University of Melbourne, Austin
22 Health, 245 Burgundy Street, Heidelberg, Victoria 3084, Australia

23 ⁶ McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of
24 Medicine, Baltimore, MD, 21205, USA

25 ⁷ Bruce Lefroy Centre, Murdoch Children's Research Institute, Flemington Rd, Parkville,
26 Victoria 3052, Australia

27 ⁸ University of Melbourne Centre for Cancer Research, Victorian Comprehensive Cancer
28 Centre, 305 Grattan Street, Melbourne, Victoria 3000, Australia

29 ⁹ Illumina Inc, 5200 Illumina Way, San Diego, CA 92122, USA

30 ¹⁰ Genetic Unit, Basurto University Hospital, OSI Bilbao-Basurto, avenida Montevideo 18
31 (48013 Bilbao), Spain

32 ¹¹ Centre for Clinical Brain Sciences, University of Edinburgh, Edinburgh, EH16 4SB,
33 Scotland

- 34 ¹² Anne Rowling Regenerative Neurology Clinic, University of Edinburgh, Edinburgh, EH16
35 4SB, Scotland
- 36 ¹³ Centre for Medical Informatics, Usher Institute of Population Health Sciences and
37 Informatics, University of Edinburgh, Edinburgh, EH16 4UX, Scotland
- 38 ¹⁴ Department of Neurology, Tauranga Hospital, Private Bag, Cameron Road, Tauranga
39 3171, New Zealand
- 40 ¹⁵ University of Sydney, New South Wales 2006, Australia
- 41 ¹⁶ Royal North Shore Hospital, Pacific Hwy, St Leonards, New South Wales 2065, Australia
- 42 ¹⁷ Department of Paediatrics, University of Melbourne, Royal Children's Hospital,
43 Flemington Rd, Parkville, Victoria 3052, Australia
- 44 ¹⁸ Departments of Neurology and Human Genetics, David Geffen School of Medicine
45 University of California, Los Angeles, CA 90095, USA
- 46 ¹⁹ Hunter Genetics, Hunter New England Health Service, Waratah, Newcastle, New South
47 Wales 2300, Australia
- 48 ²⁰ University of Newcastle, Newcastle, New South Wales 2300, Australia
- 49 ²¹ Neurology Department, Royal Prince Alfred Hospital, Camperdown New South Wales,
50 2050, Australia
- 51 ²² Central Clinical School, University of Sydney, Camperdown, New South Wales 2050,
52 Australia
- 53 ²³ Servicio de Neurología, Hospital de Basurto, Avenida de Montevideo 18, 48013, Bilbao,
54 Bizkaia , Spain
- 55 ²⁴ Edmond J. Safra Program in Parkinson's disease and the Morton and Gloria Shulman
56 Movement Disorders Clinic, Toronto Western Hospital, Toronto, ON M5T 2S8, Canada
- 57 ²⁵ Department of Medicine, Division of Neurology, University Health Network and the
58 University of Toronto, Toronto, ON M5S, Canada
- 59 ²⁶ Department of Neurology, Wellington Hospital, Wellington 6021, New Zealand
- 60 ²⁷ Sunshine Neurology, Maroochydore, Queensland 4558, Australia
- 61 ²⁸ Department of Neurology, David Geffen School of Medicine, University of California, Los
62 Angeles, CA 90095, USA
- 63 ²⁹ Riddiford Medical, Newtown, Wellington 6023, New Zealand
- 64 ³⁰ Department of Neuroscience, Central Clinical School, Monash University, Alfred Hospital
65 Campus, Commercial Road, Melbourne, Victoria 3004, Australia
- 66 ³¹ Institute of Neurological Sciences, Prince of Wales Hospital, Randwick, New South Wales
67 2031, Australia

68 ³² Department of Neurology, Johns Hopkins Hospital, Baltimore, MD 21287 USA

69

70 ³³ These authors contributed equally to this work

71 ³⁴ These authors contributed equally to this work

72

73 * Correspondence to Dr Paul Lockhart, Murdoch Children's Research Institute, Flemington
74 Rd, Parkville, Victoria 3052, Australia; paul.lockhart@mcri.edu.au

75

76 Running title: An intronic repeat expansion in RFC1 causes CANVAS

77

78 Keywords: CANVAS, ataxia, repeat expansions, short tandem repeats, whole genome
79 sequencing.

80 **ABSTRACT**

81 Genomic technologies such as Next Generation Sequencing (NGS) are revolutionizing
82 molecular diagnostics and clinical medicine. However, these approaches have proven
83 inefficient at identifying pathogenic repeat expansions. Here, we apply a collection of
84 bioinformatics tools that can be utilized to identify either known or novel expanded repeat
85 sequences in NGS data. We performed genetic studies of a cohort of 35 individuals from 22
86 families with a clinical diagnosis of cerebellar ataxia with neuropathy and bilateral vestibular
87 areflexia syndrome (CANVAS). Analysis of whole genome sequence (WGS) data with five
88 independent algorithms identified a recessively inherited intronic repeat expansion
89 [(AAGGG)_{exp}] in the gene encoding Replication Factor C1 (*RFC1*). This motif, not reported
90 in the reference sequence, localized to an Alu element and replaced the reference (AAAAG)₁₁
91 short tandem repeat. Genetic analyses confirmed the pathogenic expansion in 18 of 22
92 CANVAS families and identified a core ancestral haplotype, estimated to have arisen in
93 Europe over twenty-five thousand years ago. WGS of the four *RFC1* negative CANVAS
94 families identified plausible variants in three, with genomic re-diagnosis of SCA3, spastic
95 ataxia of the Charlevoix-Saguenay type and SCA45. This study identified the genetic basis of
96 CANVAS and demonstrated that these improved bioinformatics tools increase the diagnostic
97 utility of WGS to determine the genetic basis of a heterogeneous group of clinically
98 overlapping neurogenetic disorders.

99

100 INTRODUCTION

101 Repetitive DNA sequences constitute approximately one third of the genome and are
102 thought to contribute to diversity within and between species.¹ Microsatellites or short
103 tandem repeats (STRs) are mini-repeats of DNA, typically two to five base-pairs in length,
104 which are usually present in a concatamer of between five and fifty repeated elements. There
105 are thousands of STRs scattered through the human genome and recent studies have
106 suggested important roles for STRs in the regulation of gene expression.^{2;3} STRs display
107 considerable variability in length between individuals, which is presumed to have no
108 detrimental consequences for humans^{4;5} unless the repeat number is expanded beyond a
109 gene-specific threshold.^{6;7} Pathogenic repeat expansions (REs) have been shown to underlie
110 at least 30 inherited human diseases, the majority being disorders of the nervous system.⁸
111 These disorders, which variably have autosomal dominant, autosomal recessive and X-linked
112 inheritance, have an overall prevalence of ~1:20,000.⁹ They display a broad onset age and are
113 characterized by progressive cerebellar ataxia with dysarthria, oculomotor abnormalities,
114 cognitive dysfunction and other symptoms.¹⁰ Additional novel pathogenic REs likely remain
115 to be identified. For example, putative spinocerebellar ataxia (SCA) loci, including SCA25
116 (MIM: 608703) and SCA30 (MIM: 613371) remain to be identified, and unsolved hereditary
117 ataxias such as cerebellar ataxia with neuropathy and bilateral vestibular areflexia syndrome
118 (CANVAS, MIM: 614575) display extensive clinical similarities with known RE disorders.

119

120 CANVAS is a cerebellar ataxia with combined cerebellar, vestibular and
121 somatosensory dysfunction.^{11;12} Historically, individuals with CANVAS have been assigned
122 the diagnosis of idiopathic late onset cerebellar ataxia.¹³ More recently, CANVAS is
123 clinically recognized and has been incorporated into the contemporary research and teaching
124 of both cerebellar and vestibular diseases.^{14;15} Unifying the oto- and neuropathology,
125 CANVAS is a neuronopathy (ganglionopathy) affecting the vestibular¹⁶ and dorsal root
126 ganglia.¹⁷ The progression of these clinical features can be measured longitudinally using a
127 specific neurophysiological protocol.¹⁸ A characteristic radiological pattern of cerebellar
128 atrophy has also been described and verified on post-mortem pathology.¹¹ The characteristic
129 oculomotor abnormality seen in combined cerebellar and vestibular impairment is the
130 visually-enhanced vestibulo-ocular reflex (VVOR), and this can now be evaluated using a
131 commercially available instrumented assessment tool.¹⁹⁻²¹ Altogether, these advances have
132 allowed the formulation of diagnostic criteria to aid identification of CANVAS, contributing
133 both research and clinical benefits including improved prognostication and targeted

134 management.^{12, 14} While detailed clinical findings have driven gene discovery in RE disorders
135 such as Friedreich ataxia²² the underlying genetic cause(s) of CANVAS has, until very
136 recently, remained elusive (see below).

137 The majority of individuals and families with CANVAS have been identified in
138 individuals of European ancestry, although CANVAS has recently been reported in two
139 individuals of Japanese ethnicity, a 68-year old male²³ and a 76 year old female.²⁴ A genetic
140 cause of CANVAS is highly plausible given the observation of 13 affected siblings and
141 families with multiple affected individuals over several generations.¹² The pattern of
142 inheritance suggests an autosomal recessive trait, although autosomal dominant inheritance
143 with incomplete penetrance cannot be excluded. CANVAS symptoms overlap considerably
144 with SCA3 (also known as Machado-Joseph disease, MIM: 109150) and Friedreich ataxia
145 (MIM: 229300), both genetic forms of ataxia caused by the inheritance of a pathogenic RE.
146 These observations are consistent with the hypothesis that a novel pathogenic STR expansion
147 may underlie CANVAS.

148

149 Historically, the detection of REs has been time-consuming and expensive. Indeed, it
150 is only in recent years that computational methods have been developed to screen for RE in
151 short-read whole exome sequence (WES) and WGS data²⁵, leading to the discovery of novel,
152 disease causing REs. For example, a pentanucleotide RE was identified to underlie autosomal
153 dominant spinocerebellar ataxia 37 (SCA37; OMIM: 615945).²⁶ Moreover, pathogenic REs
154 of intronic pentamers (TTTCA)_n and (TTTTA)_n were identified as the cause of Benign Adult
155 Familial Myoclonus Epilepsy locus 1, 6 and 7 (BAFME1, OMIM: 618073; BAFME6,
156 OMIM: 618074; BAFME1, OMIM: 618075).²⁷

157

158 A number of bioinformatics tools now exist that allow screening of short-read
159 sequencing data for expanded STRs.²⁵ Initially, STR detection tools, such as lobSTR and
160 hipSTR, were limited to short STRs that were encompassed by a single sequencing read.
161 However, in the last two years, multiple methods have been released that can screen WES
162 and WGS datasets for REs without being limited by read length. These include
163 ExpansionHunter (EH)²⁸, exSTRA⁸, TREDPARSE²⁹, STRetch³⁰ and GangSTR.³¹ These are
164 all reference based methods - i.e. they rely on a catalogue of STR loci and motifs and are
165 therefore limited to detecting expansion of previously defined STRs, such as those catalogued
166 in the UCSC track. Moreover, the normal variability in STR length and repeat composition
167 remains poorly described, particularly for rare STRs or those larger than ~100bp. Therefore,

168 there is a need for bioinformatics tools that are unbiased to the limited catalogues of STR loci
169 available. Ideally, these tools will be able to search genome-wide for expanded repeat
170 sequences in NGS data, independent of prior knowledge of either the location or composition
171 of the RE. Here, we utilized a STR reference-free method called Expansion Hunter De Novo
172 (EHdn), in combination with multiple reference-based tools, to show that CANVAS is caused
173 by the homozygous inheritance of a novel and expanded intronic pentamer [(AAGGG)_{exp}] in
174 the gene encoding Replication Factor C Subunit 1 (*RFC1*). An independent study, published
175 while this work was under review, similarly identified the causal pentamer in *RFC1*. Cortese
176 and colleagues defined a small linkage region from ten families with CANVAS and the
177 causative RE was identified by WGS and visual inspection of the aligned read pairs inside the
178 linkage region.³²

179 **MATERIALS AND METHODS**

180 **Recruitment, linkage and next generation sequence data**

181 The Royal Children's Hospital Human Research Ethics Committee approved the
182 study (HREC 28097). Informed consent was obtained from all participants and clinical
183 details were collected from clinical assessments and review of medical records. Genomic
184 DNA was isolated from peripheral blood. Single nucleotide polymorphism (SNP) genotype
185 data were generated for two affected siblings from three families (CANVAS1, 2, 3) and all
186 six siblings from family CANVAS4 using the Illumina Infinium HumanOmniExpress
187 BeadChip genotyping array. SNP genotypes for individuals from CANVAS9 were extracted
188 from WES data.³³ Parametric multipoint linkage analysis was subsequently performed using
189 LINKDATAGEN and MERLIN^{34;35} specifying a rare recessive disease model with complete
190 penetrance, and overlapping linkage signals were detected using BEDtools.³⁶ WES was
191 performed on individuals from CANVAS9 using Agilent SureSelect XT Human All exon V5
192 + UTR on the Illumina HiSeq2000 platform at 50x mean coverage. WES was performed on
193 an additional 23 individuals from 15 families in collaboration with the Johns Hopkins Center
194 for Inherited Disease Research (CIDR) as part of the Baylor-Hopkins Center for Mendelian
195 Genomics (BHCMG). WGS was performed in two stages. Libraries for the first round of
196 samples, including two affected individuals from CANVAS1 and CANVAS9 and 31
197 individuals lacking a clinical diagnosis of CANVAS (subsequently referred to as controls
198 although some have a diagnosis other than CANVAS), were prepared using the TruSeq nano
199 PCR-based Library Preparation Kit and sequenced on the Illumina HiSeq X platform.
200 Libraries for the second round of WGS, including affected individuals with evidence of an
201 alternate RE motif (CANVAS2 and CANVAS8) or lacking the pathogenic RE in *RFC1*
202 (CANVAS11,13, 17 and 19), were prepared using the TruSeq PCR-free DNA HT Library
203 Preparation Kit and sequenced on the Illumina NovaSeq 6000 platform. PCR-free WGS data
204 from 69 unrelated Coriell controls²⁸ was obtained from Illumina. GTEx samples (SRA files,
205 133 WGS with matching cerebellar RNA-seq) were downloaded from the dbGAP
206 (phs000424.v7.p2).

207

208 **Alignment and variant calling**

209 Alignment and haplotype calling were performed based on the GATK best practice
210 pipeline. All WES and WGS datasets were aligned to the hg19 reference genome using
211 BWA-mem, then duplicate marking, local realignment and recalibration were performed with
212 GATK. Merged VCF files were annotated using vcfanno³⁷ and ANNOVAR.³⁸ Candidate

213 variant filtering was performed using CAVALIER, an R package for variant interpretation in
214 NGS data (<https://github.com/bahlolab/cavalier>). Standard variant calling was performed on
215 WGS data for CANVAS samples negative for the pathogenic RE in *RFC1*. Candidate
216 variants were defined as i) occurring in known ataxia genes, as defined by OMIM, ii) exonic,
217 with a minor allele frequency of less than 0.0001 in gnomAD (both genome and exome data)
218 and iii) predicted pathogenic by both SIFT and PolyPhen2. RNA-seq data was aligned to the
219 hg19 reference genome (ENSEMBL Homo_sapiens.GRCh37.75) using STAR.³⁹ Reads were
220 summarized by gene ID into a counts matrix using featureCounts⁴⁰ (quality score ≥ 10) and
221 converted to log10 of the counts per million using limma.⁴¹

222

223 **STR analysis**

224 Genome-wide screening for putative REs was performed using Expansion Hunter
225 Denovo (EHdn) version 0.6.2, an open-source method that is being developed by Illumina,
226 Inc, the Walter Eliza Hall Institute and others. EHdn operates by performing a genome-wide
227 search for read pairs where one mate has confident alignment (anchor) and the second mate
228 consists of repetition of a repeat motif (in-repeat read). The program reports the counts of in-
229 repeat reads with anchor mates stratified by the repeat motif and genomic position of their
230 anchor mate. For this analysis we defined a confidently-aligned read as one aligned with
231 MAPQ of 50 or above. The counts of in-repeat reads with anchor mates were subsequently
232 compared for each region in cases (CANVAS) and controls using a permutation test (10^6
233 permutations). The resulting p-values were used to rank candidate sites with higher counts in
234 individuals with CANVAS than in the controls for further computational validation. These
235 candidates were subsequently annotated with ANNOVAR.

236

237 Computational validation was performed using five independent STR detection tools
238 for short-read NGS after updating the STR catalogue reference files to incorporate the
239 identified motifs. The RE candidates were screened in the two individuals with CANVAS
240 and the 31 non-CANVAS controls using exSTRa and EH, then the top candidate
241 [(AAGGG)_{exp} STR in *RFC1*] was further validated with TREDPARSE, GangSTR and
242 STRetch. All tools were used with default parameters, with the following additional
243 parameters for EH: read-depth of 30 and min-anchor-mapq of 20. All five tools were also
244 used to screen for the (AAGGG)_{exp} *RFC1* STR in the 69 Coriell control WGS datasets. A
245 short-list of (AAGGG)_{exp} carriers was generated based on consensus calling from at least four
246 of the five tools.

247

248 Individuals diagnosed with CANVAS lacking the (AAGGG)_{exp} *RFC1* RE were
249 further screened with EHdn for novel STRs and for known pathogenic STRs using exSTRa
250 and EH. The WES datasets could not be analyzed for the (AAGGG)_{exp} *RFC1* RE as the
251 intronic locus (chr4:39350045-39350095, hg19) was not captured during library preparation.
252 However, the region was visualized using the Integrative Genomics Viewer (IGV) to identify
253 potential off-target reads which could provide supportive evidence for the presence of the
254 (AAGGG)_{exp} motif. Only samples with at least one read mapping at the STR in the *RFC1*
255 locus were considered.

256

257 **Haplotyping and mutation dating**

258 Haplotyping was performed on the WES data. Variants were filtered based on read depth
259 (≥ 30), including both exonic and non-exonic variants. A core haplotype was defined based on
260 sharing amongst a majority of affected individuals. A method based on haplotype sharing⁴²
261 was used to determine the most recent common ancestor (MRCA) from whom the core
262 haplotype was inherited, as well as dating additional sub-haplotypes shared by clusters of
263 individuals, which are likely to be individuals with a MRCA who is more recent than that for
264 the whole group (<https://shiny.wehi.edu.au/rafehi.h/mutation-dating/>).

265

266 **Molecular genetic studies**

267 We designed a PCR assay to test for presence of additional inserted sequence, not
268 present in the reference database, at the *RFC1* STR. The primers (Table S1) flank the STR
269 and are predicted to amplify a 253bp fragment using standard PCR conditions with a 30
270 second extension cycle. Presence of the pathogenic *RFC1* RE was tested by repeat-primed
271 PCR utilizing three primers; TPP_CANVAS_FAM_2F, 5R_TPP_M13R_CANVAS_RE_R
272 and TPP_M13R (Table S1). The FAM labelled forward primer is locus specific, while the
273 repeat-specific primer (5R_TPP) includes a tag M13R sequence. PCR was performed in a
274 20 μ l reaction with 20 ng genomic DNA, 0.8 μ M of both the FAM labelled forward primer
275 and TPP_M13R and 0.2 μ M 5R_TPP using GoTaq® Long PCR Polymerase (Promega). A
276 standard 60TD55 protocol was utilized (94°C denaturation for 30 s, 60TD55°C anneal for 30
277 s, and 72°C extension for 2 min), products were detected on an ABI3730xl DNA Analyzer
278 and visualized using PeakScanner 2 (Applied Biosystems).

279 RESULTS

280 Case recruitment

281 The workflow for this study is summarized in Figure 1. Individuals with a clinical
282 diagnosis of CANVAS were recruited following neurological assessment and investigation in
283 accordance with published guidelines.⁴³ While variable between cases, data leading to the
284 clinical diagnosis included evidence of combined cerebellar and bilateral vestibular
285 impairment, cerebellar atrophy on MRI, neurophysiological evidence of impaired sensory
286 nerve function and negative genetic testing for pathogenic RE at common SCA loci (typically
287 *SCA1*, 2, 3, 6 and 7) and *FRDA* (Friedreich ataxia, FRDA). In total, the cohort consisted of 35
288 individuals with a clinical diagnosis of CANVAS (Table 1). The individuals came from
289 eleven families with a single affected individual, seven families with affected sib pairs and
290 four larger/multigenerational families (Figure S1). A full clinical description of the cohort
291 will be reported in a forthcoming manuscript.

292

293 Linkage analysis

294 CANVAS typically presents in families with one or multiple affected individuals in a
295 single generation, consistent with a recessive inheritance. For example, in the second-degree
296 consanguineous family CANVAS9, four siblings were diagnosed with CANVAS and two
297 were classified as unaffected at the time phenotyping was performed (Figure 2A). Parametric
298 multipoint linkage analysis was performed on five CANVAS families (CANVAS1, -2, -3, -4
299 and -9, Figure S1) specifying a rare recessive disease model with complete penetrance. This
300 identified linkage regions with logarithm of odds (LOD) scores ranging from 0.6 for smaller
301 pedigrees (two affected siblings), to a statistically significant linkage region on chromosome
302 4 in CANVAS9 (LOD=3.25, Figure 2B). Intersection of the linkage regions from the five
303 families identified a single region on chromosome four (chr4:38887351-40463592, hg19,
304 combined LOD=7.04) common to all families (Figure 2C). CNV analysis utilizing PennCNV
305 did not identify any potential copy number variants in the minimal linkage region.⁴⁴ The
306 1.5MB shared region contains 42 genes, of which 14 are protein coding, none with any
307 association with ataxia in OMIM or the published literature (Table S2).

308

309 Large-scale WES analysis did not identify candidate pathogenic variants

310 WES was used to screen 27 affected individuals with CANVAS from 15 families for
311 potentially pathogenic rare variants (MAF < 0.001) shared across multiple pedigrees in a

312 homozygous or compound heterozygous inheritance pattern. No candidate mutations were
313 detected, either within the chromosome 4 linkage region or elsewhere in the genome.

314

315 **Identification of a novel (AAGGG)_{exp} RE in the linkage region**

316 The lack of candidate variants identified from the WES data suggested the possibility
317 of (i) intronic or intergenic mutations, or (ii) that CANVAS might be caused by a non-
318 standard mutation, such as a pathogenic RE of an STR. Therefore, WGS was performed on
319 two individuals from different pedigrees (CANVAS1 and CANVAS9) who share the chr4
320 linkage region. EHdn was used to perform a genome-wide screen for STRs in the two
321 individuals with CANVAS compared to WGS data from 31 unrelated controls. This
322 identified 19 regions with a p value < 0.005 (Table S3), although genome-wide significance
323 could not be achieved after adjustment for multiple testing due to the skewed ratio of the
324 number of cases to controls (2 versus 31). These candidate STRs were visualized with the
325 Integrative Genomics Viewer (IGV) tool, which suggested that the (AAGGG)_{exp} STR within
326 intron 2 of the gene encoding Replication Factor C1 (*RFC1*) was likely real and present in
327 both alleles in the affected individuals, consistent with the recessive inheritance pattern
328 hypothesized for CANVAS (Figure S2). In addition, this was the only candidate that (I) was
329 localized to the chr4 linkage region and (II) was able to be validated using existing STR
330 detection tools (see below). In both individuals with CANVAS, the novel (AAGGG)_{exp}
331 pentamer replaced an (AAAAG)₁₁ motif located at the same position in the reference genome
332 (chr4:39350045-39350095, hg19) and appeared to be significantly expanded compared to
333 controls. Visualization of the region in the UCSC genome browser identified that the
334 reference motif (AAAAG)₁₁ is the 3' end of an Alu element, AluSx3. In individuals with
335 CANVAS, the (AAAAG)₁₁ motif is substituted by the (AAGGG)_{exp} motif, with potential
336 interruptions to the Alu element (Figure 2D).

337

338 **Confirmation of (AAGGG)_{exp} STR in off-target WES reads**

339 While WGS was only performed in two individuals with CANVAS, the majority of
340 the cohort (n=27) was analyzed by WES. The putative pathogenic CANVAS RE is located in
341 intron 2 of *RFC1*, 2863bp downstream of exon 2 and 2952bp upstream of exon 3. Therefore
342 WES data is *a priori* assumed to be uninformative for this RE as it is not targeted during
343 DNA capture. However, given that WES data includes off-target reads, we hypothesized that
344 some reads might map to the *RFC1* RE locus. Visual assessment of the WES data in IGV
345 identified 14 individuals with informative reads; the maximum off-target read coverage at

346 this locus was two, with a median of one read. While three individuals only had reads that
347 correspond to the reference genome STR sequence (AAAAG), eleven affected individuals
348 from nine families had reads containing (AAGGG) repeats (Table 1). Furthermore, single
349 affected individuals from families CANVAS2 and CANVAS8 had single, independent reads
350 identifying (AAGGG) and (AAAGG) motifs at the *RFC1* STR locus. This observation raised
351 the possibility that CANVAS might result from pathogenic expansions of different
352 pentanucleotide motifs.

353

354 **Computational validation with existing STR detection tools**

355 Multiple tools have been developed in recent years that test for the presence of REs at
356 pre-defined STRs. Therefore, we inserted the novel *RFC1* STR motifs into the STR reference
357 files and used exSTRa, EH, TREDPARSE, STRetch and GangSTR to estimate the size of the
358 STR, and/or detect REs in the WGS data from the two original CANVAS samples
359 (CANVAS1 and 9) and seven additional individuals with CANVAS. The seven additional
360 CANVAS samples selected for WGS were those with WES evidence for an alternate
361 (AAAGG) motif (families CANVAS2 and 8), and those who did not appear to have a RE at
362 the *RFC1* locus based on the PCR/RP-PCR studies described below (CANVAS11, 13, 17 and
363 19 families). The library preparation this second round of WGS was PCR-free as PCR
364 amplification has previously been shown to affect RE detection.⁸ Using exSTRa, we
365 confirmed the homozygous inheritance of the (AAGGG)_{exp} RE in three individuals
366 (CANVAS1, 2 and 9, Figure 3). The empirical cumulative distribution function (ECDF)
367 pattern for CANVAS2 is consistent with the presence of one shorter and one longer
368 (AAGGG)_{exp} RE, while CANVAS8 appears to only have a single (AAGGG)_{exp} allele.
369 Screening of all datasets for the (AAAGG)_{exp} motif at the chr4 *RFC1* locus using exSTRa
370 identified an expansion of this STR only in CANVAS8, suggesting this individual is a
371 compound heterozygote, with an (AAGGG)_{exp} RE on one allele and an (AAAGG)_{exp} RE on
372 the other. Visualization in IGV confirmed the presence of the (AAAGG)_n motif embedded
373 within the reference STR: (AAAAG)₆-(AAAGG)_n-(AAAAG)₆ (Figure S3). This observation
374 raises the possibility that an expanded (AAAGG)_{exp} motif might also be associated with
375 CANVAS.

376

377 EH, TREDPARSE and GangSTR were used to estimate the length of the AAGGG
378 motif on each allele (Figure 3). The results were highly variable depending on the tool used.
379 EH reported minimum and maximum allele sizes ranging from 30 to 68 in individuals with

380 the RE. The allele size ranges estimated by GangSTR and TREPARE were 2 to 27 and 7 to
381 14, respectively. Furthermore, all three tools inferred the presence of two alleles, even in
382 individuals who carry a single allele, and hence do not appear to be distinguishing read
383 contributions between the alleles, also contributing to unreliable size estimates. Reads
384 comprised of the (AAGGG)_{exp} motif in particular also showed evidence of high read
385 sequencing error. Based on these results, we can infer that while the CANVAS samples were
386 all correctly identified as having homozygous RE at *RFC1*, estimates of expansion size are
387 inconsistent and appear likely to significantly underestimate the actual repeat size.

388

389 The consensus of the different tools was that CANVAS11, 13, 17 and 19 families did
390 not encode a pathogenic RE [either (AAGGG)_{exp} or (AAAGG)_{exp}] at the *RFC1* locus, which
391 we confirmed by PCR analyses (see below). However, the *RFC1* (AAGGG)_{exp} RE was
392 present in three of the control WGS datasets [two heterozygous, one homozygous, allele
393 frequency ~0.06 (4/62), Figure 3]. No control individuals were identified to carry the
394 (AAAGG)_{exp} motif. As with the CANVAS samples, the STR sizing estimates using the
395 different tools was inconsistent, therefore no conclusions could be drawn from this *in silico*
396 analysis regarding the relative size of the (AAGGG)_{exp} RE in controls compared to
397 individuals with CANVAS. We then analyzed a larger in-house collection of unrelated
398 control Coriell WGS samples (N=69) and again failed to identify the (AAAGG)_{exp} motif.
399 However, we identified six individuals heterozygous for the (AAGGG)_{exp} RE, representing a
400 frequency estimate of ~0.04 [(6/138) (Figure S4)]. Using the NGS QC software tool peddy,
401 we found evidence that two of these heterozygous individuals are of European ancestry and
402 that two further individuals are of admixed Native American ancestry. Finally, we accessed
403 WGS from GTEx for 133 individuals who have matching brain (cerebellum) RNA-seq. Our
404 analysis identified 11 heterozygous carriers of the (AAGGG)_{exp} RE, representing an
405 estimated allele frequency of ~0.04 (11/266), consistent with our in-house collection.

406

407 **Validation of the (AAGGG)_n RE as the causal variant for CANVAS**

408 We developed a PCR assay that amplifies across the repeat tract to rapidly screen for
409 the presence of a non-expanded allele at the *RFC1* STR. Although the screen does not
410 distinguish between the (AAAAG)₁₁ reference STR, the (AAGGG) STR, or any other
411 potential motif, amplification of an ~250bp fragment indicates at least one allele is not
412 expanded. Moreover, the presence of two distinct non-expanded products is indicative of a
413 heterozygous non-mutant state. Conversely, the complete absence of the PCR product

414 provides indirect evidence of a RE affecting both alleles of the *RFC1* locus. Analysis of all
415 available DNA samples from individuals with CANVAS suggested that the reference STR at
416 the *RFC1* locus was not present in 30 clinically diagnosed individuals from 18 (of 22)
417 CANVAS families (Table 1, Figure 4). Notably, unaffected individuals from *RFC1* positive
418 families carried at least one non-expanded *RFC1* allele (Figure S5). To directly confirm
419 expansion of the novel (AAGGG) motif in *RFC1*, we developed a locus specific repeat-
420 primed PCR assay, using a primer located adjacent to the *RFC1* repeat and an AAGGG-
421 specific primer. Consistent with the PCR assay, affected individuals from the 18 families
422 demonstrated a saw-toothed ‘ladder’ when the repeat-primed PCR products were analyzed by
423 capillary array (Table 1, Figure 4). These results suggest a homozygous RE underlies
424 CANVAS in these 18 families, and at least one pathogenic allele encodes the (AAGGG)_{exp}
425 RE. Molecular analysis of the DNA for the three in-house control individuals with the *in*
426 *silico* predicted (AAGGG)_{exp} motif (Figure 3) demonstrated a ~250bp product in both
427 heterozygous samples but no product in the homozygous individual. The repeat-primed assay
428 demonstrated a saw-toothed ladder in all three samples (Figure S6). Collectively, these
429 analyses suggested all three control individuals have at least one copy of the pathogenic
430 (AAGGG)_{exp} RE at *RFC1*, although the size of the RE cannot be determined by these
431 analyses.

432

433 In four families (CANVAS11, 13, 17 and 19) the presence of the expected reference
434 PCR amplicon and lack of a repeat-primed PCR product suggested the pathogenic *RFC1* RE
435 was not present on either allele. This implied that these individuals have a different
436 CANVAS-causing mutation in *RFC1*, or there is locus heterogeneity. A third possibility is
437 that they do not have CANVAS but instead a related ataxia. Therefore, we performed WGS
438 on these individuals and initially screened for known REs associated with ataxias using
439 exSTRa and EH. A CAG trinucleotide expansion in *ATXN3*, associated with spino-cerebellar
440 ataxia type 3 (SCA3; OMIM 109150 also known as Machado-Josephs disease) was identified
441 in CANVAS13 (Figure S7) and confirmed by diagnostic testing. The WGS was then screened
442 for novel or rare SNPs and indels in genes known to cause ataxia. No *de novo* or rare variants
443 were identified in *RFC1* however a potential genomic re-diagnosis was achieved in two
444 additional families. In CANVAS17 two variants [NM_001278055:c.12398delT,
445 p.(Phe4133Serfs*28) and NM_001278055:c.5306T>A, p.(Val1769Asp)] were identified in
446 the gene encoding saccin (*SACS*) and segregation analysis confirmed they were in trans.
447 Biallelic mutations in *SACS* cause spastic ataxia of the Charlevoix-Saguenay type (MIM:

448 270550). In CANVAS19, a heterozygous variant in the gene encoding FAT tumor suppressor
449 homolog 2 [*FAT2*, NM_001447.2:c.4370T>C, p.(Val1457Ala)] was identified. Heterozygous
450 mutations in *FAT2* have recently been associated with SCA45 (MIM: 604269).⁴⁵ No
451 potentially pathogenic variants were identified in CANVAS11, however a variant of
452 unknown significance was identified in the gene encoding Ataxin 7 [*ATXN7*,
453 NM_001177387.1:c.2827C>G, p.(Arg943Gly)]. CANVAS11 was also screened genome-
454 wide with EHDn for potentially pathogenic novel RE, however no additional candidate REs
455 were identified.

456

457 **A single founder event for the (AAGGG)_n RE in *RFC1***

458 We performed haplotype analysis to determine if the (AAGGG)_{exp} RE arose more
459 than once in human history. Analysis of haplotypes inferred from the WES data identified a
460 core ancestral haplotype, comprised of 27 SNPs (Figure 5A), that was shared by most
461 individuals except CANVAS14 (Table 1, Table S4). The core haplotype spans four genes
462 (*TMEM156*, *KLHL5*, *WDR19* and *RFC1*) and is 0.36 MB in size (chr4:38995374-39353137
463 (hg19)). Inspection of this region in the UCSC browser suggested that the core haplotype
464 overlaps with a region of strong linkage disequilibrium in European and Asian populations
465 (Han Chinese and Japanese from Tokyo), but not the Yoruba population (an ethnic group
466 from West Africa, Figure 5B). Using a DNA recombination and haplotype-based mutation
467 dating technique⁴², we estimate that the most recent common ancestor (MRCA) of the
468 CANVAS cohort lived approximately 25,880 (CI: 14080-48020) years ago (Figure 5C). This
469 age estimate corresponds to the size of the haplotype and LD block and is roughly equivalent
470 to the origin of modern Europeans as represented by the HAPMAP CEU cohort. Further
471 investigation of the haplotypes allowed us to infer a simple phylogeny based on identified
472 clusters of shared haplotypes extending beyond the core haplotype, suggesting that some
473 individuals have common ancestors more recent than that of the MRCA for the whole group.
474 This approach identified four subgroups. Group A had a MRCA dating back 5,600 (CI: 2120-
475 15520) years and group B (further divided into groups B1 and B2) have a MRCA dating back
476 4,180 years (CI: 2240-7940). Furthermore, one individual shared part of their haplotype with
477 both groups A and B, suggesting that group B is a distant branch of the MRCA of group A.
478 Another subgroup, C, has a MRCA that lived 1860 (CI: 560-7020) years ago. The final group
479 labelled N, do not have any additional sharing beyond the core haplotype.

480 Next, we compared the haplotype of the nine control samples (three in-house controls
481 and six from the Coriell collection) that carry the (AAGGG)_{exp} RE to the core haplotype

482 defined in the individuals with CANVAS. All controls shared at least part of the core
483 haplotype, again suggesting that the (AAGGG)_{exp} RE arose once in history. Finally, we
484 determined that nine of the 11 individuals from GTEx heterozygous for the (AAGGG)_{exp} RE
485 also shared the same core haplotype identified in individuals with CANVAS. The haplotype-
486 specific SNP rs2066782 (exon 18, chr4:39303925, A>G) enabled us to analyse the
487 expression of the (AAGGG)_{exp} *RFC1* allele in the cerebellum RNA-seq data and confirm that
488 the STR did not inhibit the expression of *RFC1* compared to the reference (AAAAG)₁₁ allele.
489 The remaining two carriers do not appear to share the core haplotype. As they do not have
490 heterozygous SNPs in their exons, allele specific expression could not be determined.

491 **DISCUSSION**

492 Since the first description of the syndrome of cerebellar ataxia with bilateral
493 vestibulopathy in 2004⁴⁶ and proposal of CANVAS as a distinct clinical entity in 2011¹¹ there
494 has been little progress made in delineating the etiology of the disorder. While most affected
495 individuals are described as idiopathic, reports of multiple affected sib pairs¹² and a family
496 with three affected individuals⁴⁷ have suggested that an autosomal recessive mode of
497 inheritance is most likely. The genetic basis of CANVAS has now been identified and
498 validated in two independent studies, one recently published by Cortese *et al*³² and this study.
499 Both studies utilized a similar study design, with linkage analysis to reduce the genomic
500 search space to a modest interval (<2Mb), but no plausible causal variant(s) could be
501 identified in WES data. WGS was then performed on multiple individuals and Cortese *et al*
502 successfully identified the RE by visual inspection of the aligned read pairs inside the linkage
503 region using the Integrative Genomics Viewer. In contrast, we utilized a bioinformatics
504 approach and performed genome-wide analysis of WGS data to identify potential RE and
505 then prioritized the RE located within the linkage interval. While both approaches were
506 successful, the bioinformatics approach to RE detection, as described in this study, is likely
507 more sensitive and practical, and can be applied even in the absence of a small, or indeed any,
508 linkage region. Furthermore, using a bioinformatics approach allows simultaneous testing of
509 other potentially causal RE due to differential diagnoses. For example, we quickly re-
510 diagnosed an affected individual with a pathogenic SCA3 RE.

511

512 Previously, the only variant associated with CANVAS was a heterozygous missense
513 variant in the gene encoding E74 Like ETS Transcription Factor 2 (*ELF2*), which segregated
514 with the disorder in three individuals in a single family.⁴⁷ It is now apparent that the majority
515 of individuals with CANVAS result from the homozygous inheritance of an expanded
516 intronic pentamer in *RFC1*. We found the (AAGGG)_{exp} in 30 of 31 individuals with a RE at
517 this locus. In only a single individual did we observe a different, presumably pathogenic
518 motif; CANVAS8 had one allele with the (AAGGG)_{exp}, whereas the second allele appeared
519 to consist of an (AAAGG)_{exp}. Notably, this alternate motif does not share the AAGGG
520 haplotype (Figure S3). Analysis of the core haplotype in the majority of individuals with
521 CANVAS suggests that the (AAGGG)_{exp} RE arose once, approximately 25,000 years ago,
522 most likely in Europe. While the majority of individuals in our cohort who carry the
523 (AAGGG)_n RE are of European ancestry, the RE is also present in non-European individuals,
524 including a Lebanese family and two carriers of admixed Native American ancestry. Given

525 the age of the CANVAS RE and recent human admixture it is likely that the locus may
526 underlie CANVAS in apparently non-European individuals, despite the disorder being highly
527 overrepresented in European populations.

528

529 Importantly, Cortese *et al* extend the clinical significance of the CANVAS RE by
530 demonstrating it is potentially a common cause of unsolved ataxia not meeting the diagnostic
531 criteria of CANVAS. Screening for homozygous inheritance of the (AAGGG)_{exp} RE in a
532 cohort of 150 individuals with sporadic late-onset ataxia diagnosed 33 individuals (22%).
533 This is consistent with the relatively high allele frequency of the (AAGGG)_{exp} we report in
534 this paper. Collectively, the two studies screened for the RE in a total of 537 clinically normal
535 samples, identifying 23 heterozygous and a single homozygous individual (allele frequency
536 $25/1074=0.023$). Given that the allele size and RE composition could not be determined in all
537 controls, it is possible that the unaffected homozygous individual we identified carries two
538 alleles smaller than the pathogenic range of >400 repeats reported by Cortese *et al*. However,
539 the individual is less than half the mean age of CANVAS onset (~60 years) and the lack of
540 phenotype suggests the clinical features are yet to manifest.

541

542 **Mechanism of pathogenicity**

543 There are multiple mechanisms by which RE can lead to pathogenicity, including
544 RNA toxicity, protein toxicity and loss or gain of function.⁷ It is not yet known how the
545 (AAGGG)_{exp} RE in *RFC1* causes CANVAS, however the homozygous inheritance pattern
546 suggests a loss-of-function mechanism, rather than RNA or protein toxicity. In heterozygous
547 carriers of the (AAGGG)_{exp} RE, our analysis of the GTEx RNA-seq data using haplotype
548 tagging SNPs suggested that the pathogenic (AAGGG)_{exp} allele did not inhibit expression of
549 *RFC1* compared to the reference (AAAAG)₁₁ allele. Interestingly, Cortese *et al* also were
550 unable to determine a mechanism of action. The (AAGGG)_{exp} RE did not appear to alter
551 expression levels of *RFC1* or surrounding genes as determined by bulk RNAseq and qRT-
552 PCR. Similarly, *RFC1* expression and protein levels appeared unchanged in peripheral or
553 brain tissue derived from individuals with CANVAS, and no AAGGG RNA foci deposits
554 were observed.³² While *RFC1* has not been previously associated with any disorder, it
555 appears extremely intolerant to LoF (pLI = 0.97; observed/expected = 0.18, CI 0.12-0.31).⁴⁸
556 In addition, siblings in the families studied carried the pathogenic RE in a heterozygous state
557 but did not manifest any signs of the disorder. This observation is analogous to Friedreich's

558 ataxia, a recessive genetic ataxia caused by loss of function (LoF) of FRDA due to a
559 pathogenic intronic RE.

560

561 *RFC1* encodes a subunit of replication factor C, a five-subunit protein complex
562 required for DNA replication and repair. Analysis of the Genotype-Tissue Expression
563 (GTEx) database demonstrated significant expression of *RFC1* in brain tissue, particularly the
564 cerebellum. Replication factor C catalyzes opening the protein ring of proliferating cell
565 nuclear antigen (PCNA), allowing it to encircle the DNA and function as a scaffold to recruit
566 proteins involved in DNA replication, repair and remodeling.⁴⁹ Mutations in multiple DNA
567 replication and repair genes such as *TCD1*, *PNKP*, *XRCC1* and *APTX* result in ataxia⁵⁰,
568 highlighting the central role of this pathway in these overlapping disorders. One of the best
569 known examples is the severe and early onset autosomal recessive disorder, ataxia
570 telangiectasia, which is caused by mutations in the gene encoding ATM serine/threonine
571 kinase (*ATM*), which is important for the repair of DNA double-strand breaks.⁵¹

572

573 The minimum pathogenic length and fine structure of the *RFC1* RE is currently
574 unclear. While Cortese et al reported a pathogenic range of ~400-2000, the individual repeat
575 composition and a more precise repeat length was not determined. The short-read NGS
576 technologies utilized in this study were unable to extend more than ~100bp into the repeat
577 sequence and efforts to amplify across the region using long range PCR were unsuccessful.
578 While the repeat-primed PCR assay indicates the presence of the (AAGGG)_{exp} motif, it does
579 not extend beyond ~250bp (50 repeat units). The application of long read sequencing
580 technologies currently being developed for RE disorders will be required to accurately
581 elucidate both the length of the pathogenic allele and the repeat composition. Both of these
582 parameters provide important clinical information regarding onset, progression and
583 pathogenicity in other genetic ataxias such as SCA1 and Friedreich ataxia.^{52; 53} Additional
584 studies will also be required to elucidate the nature of the *RFC1* STR in control individuals.
585 Cortese *et al* demonstrated considerable variability is present in the size and composition of
586 the STR, but details regarding the size and composition of both normal and pathogenic alleles
587 are yet to be fully determined. We show that the (AAGGG)_{exp} RE occurs within the 3 prime
588 end of the Alu element, AluSx3. Alu elements typically have A-rich tails and in the reference
589 sequence the *RFC1* Alu has an A-rich tail containing an (AAAAG)₁₁ STR. There is some
590 evidence that motifs that follow the pattern A_nG_m, especially (AAAG)_n and (AAAGG)_n,
591 display strong base-stacking interactions and are more likely to expand through replication

592 slippage.³¹ This suggests an inherent mitotic instability of A and G rich motifs, consistent
593 with what we observe in CANVAS. Notably, a number of pathogenic RE located with Alu
594 have previously been described, including SCA10, SCA31, SCA37 and Friedreich ataxia.^{22;}
595 26; 54; 55

596

597 **Genomic re-diagnosis in CANVAS**

598 Four of twenty two families enrolled in this study with a clinical diagnosis of
599 CANVAS did not harbor the RE or any other potentially pathogenic variants in the *RFC1*
600 locus. CANVAS13 was re-diagnosed with SCA3 after the WGS data was analyzed using our
601 computational pipeline for detecting known pathogenic REs. In addition to cerebellar ataxia,
602 individuals with SCA3 not uncommonly manifests a somatosensory impairment^{56; 57} and
603 vestibular involvement may be variably present⁵⁷, resulting in a phenotype indistinguishable
604 from CANVAS.⁴³ This molecular re-diagnosis highlights the power of modern STR detection
605 techniques to diagnose RE ataxias. In addition, NGS data provides the opportunity to
606 simultaneously identify non-RE mediated causes of ataxia. In CANVAS17 we identified
607 biallelic variants in *SACS* as the likely cause of disease. While individuals with spastic ataxia
608 of the Charlevoix-Saguenay type may present with the combination of cerebellar ataxia and a
609 peripheral neuropathy^{58; 59} as seen in CANVAS, to our knowledge vestibular involvement has
610 not previously been described, and so this potentially constitutes a novel manifestation of the
611 disease. In addition, a very plausible heterozygous variant was identified in *FAT2* in
612 CANVAS19. While classified as a VUS using ACMG guidelines,⁶⁰ the variant is only
613 observed once in gnomAD and was predicted pathogenic by multiple *in silico* algorithms.
614 Very recently, heterozygous point mutations affecting the last cadherin domain
615 (p.Lys3586Asn) or the linker region (p.Arg3649Gln) of *FAT2* have been associated with
616 SCA45, adding weight to classifying the variant (p.Val1457Ala) in the thirteenth cadherin
617 domain, as likely pathogenic. While the published clinical phenotype and mutation spectrum
618 in SCA45 is limited, in common with CANVAS, it is a late onset and slowly progressive
619 cerebellar ataxia.⁴⁵

620

621 **Strengths and limitations of current STR detection tools**

622 In this study, we implemented multiple computational tools to identify and validate
623 the presence of a novel (AAGGG)_{exp} RE in the majority of individuals with a clinical
624 diagnosis of CANVAS. In particular, the use of EHdn, with its non-reference based RE
625 discovery framework, was crucial in identifying a putative candidate, with the reference-

626 based STR detection tools facilitating the follow up analysis. Although all tools gave highly
627 variable estimated repeat sizes, which are likely to be significantly less than the actual repeat
628 size, they provided consistent evidence that the (AAGGG) motif was expanded. This level of
629 evidence is helpful before embarking on the potentially complex process of molecular
630 validation. In our analysis, only a single tool (GangSTR) failed to detect the alternate
631 (AAAGG)_{exp} RE. It is not clear why this was the case, although it could be related to the
632 more complicated [(AAAAG)₆-(AAAGG)_{exp}-(AAAAG)₆] repeat structure. Notably, EHdn
633 was able to identify the (AAAGG)_{exp} RE in CANVAS8-8, however genome wide
634 significance was not achieved and the motif was less highly ranked based on p-value than the
635 initial discovery of AAGGG (Table S3). This is likely due to the fact that the RE was only
636 present in one copy and power was reduced as there was only a single case and a small
637 number of controls. EHdn appears most effective as a discovery tool when used with multiple
638 cases or larger numbers of controls. The results of this study highlight the importance of
639 utilizing multiple tools to provide redundancy in the data analysis pipeline. We have now
640 updated the exSTRa package (see weblinks below) to include CANVAS and other recently
641 described pathogenic RE, providing additional utility to the research community to rapidly
642 identify these RE in their cohorts. An additional issue we encountered, which potentially
643 limited all tools, was the poor sequencing quality in reads containing the (AAGGG) motif
644 compared to other STRs.

645

646 In conclusion, in this study we show that a recessively inherited, ancient RE located in
647 intron 2 of *RFC1* is the predominant cause of CANVAS. Recently developed RE discovery
648 tools facilitated the identification and verification of this novel RE, in addition to identifying
649 other genetic causes of disease in the cohort. Despite the RE being located in an intron, we
650 demonstrate that previously generated WES data with low-coverage genome-wide off target
651 reads were helpful in providing increased statistical confidence in RE identification.
652 Therefore, reanalysis of previously generated WES datasets potentially offers a cost effective
653 approach to facilitating identification of novel intronic RE in discovery projects. Finally, we
654 anticipate that implementation of these tools into routine diagnostic pipelines has the
655 potential to significantly increase the current diagnostic rates of 36% and 17%, recorded for
656 clinical exome and targeted panel analyses of individuals with ataxia, respectively.^{61; 62}

657 **SUPPLEMENTAL DATA**

658 The supplemental data contain 7 figures and 4 tables.

659

660 **CONFLICTS OF INTEREST**

661 The authors declare no conflicts of interest.

662

663 **ACKNOWLEDGMENTS**

664 We acknowledge access to the dataset EGA00001003562 from the European Genome-
665 Phenome Archive. This work was supported by the Australian Government National Health
666 and Medical Research Council (Program Grant 1054618 to MB), the NIH (NINDS grant
667 R01NS082094 to BLF) and the Murdoch Children's Research Institute. MB was supported
668 by an NHMRC Senior Research Fellowship (1102971) and DPB was supported by a
669 Wellcome Clinical Research Career Development Fellowship. Additional funding was
670 provided by the Independent Research Institute Infrastructure Support Scheme and the
671 Victorian State Government Operational Infrastructure Program.

672

673 **WEB RESOURCES**

674 exSTRa: <https://github.com/bahlolab/exSTRa>

675 Genotype-Tissue Expression (GTEx) project: <https://gtexportal.org/home/>

676 Genome Aggregation Database (gnomAD): <http://gnomad.broadinstitute.org/>

677 Integrative Genomics Viewer (IGV): <http://software.broadinstitute.org/software/igv/>

678 Online Mendelian Inheritance in Man: <http://www.omim.org/>

679 UCSC Genome Bioinformatics database: <https://genome.ucsc.edu/>

680 Varsome: <https://varsome.com>

681

682 **ACCESSION NUMBERS**

683 The ClinVar details for the *RFC1* variants reported in this paper are accessible via submission
684 SUB5220746.

685

686 **LEGENDS**

687 **Figure 1: Overview of the CANVAS study and genetic investigations performed.**

688

689 **Figure 2: Linkage of the CANVAS locus to chromosome 4 and identification of**
690 **(AAGGG)_{exp} intronic insertion in *RFC1***

691 A. The pedigree of the family CANVAS9 highlights the apparent recessive inheritance
692 pattern. B. Linkage analysis of CANVAS9 identified significant linkage to chromosome 4
693 (LOD=3.25). C. Linkage regions for individual families CANVAS1, 2, 3, 4 and 9 are shown
694 in blue and the overlapping region shown in red (chr4:38887351-40463592, combined
695 LOD=7.04). D. STR analysis of WGS from two unrelated individuals with CANVAS
696 identified a novel expanded STR in the second intron of *RFC1*. The (AAAAG)₁₁ motif that is
697 present in the reference genome and part of an existing Alu element (AluSx3) is replaced by
698 the (AAGGG)_{exp} RE.

699

700 **Figure 3: Computational validation of the (AAGGG)_{exp} RE**

701 The (AAGGG)_{exp} RE at the coordinates chr4:39350045-39350095 was added to the reference
702 databases of the tools exSTRa, EH, GangSTR, TREDPARSE and STRetch and WGS data
703 from four unrelated individuals with CANVAS was analysed [CANVAS1 (orange),
704 CANVAS2 (blue), CANVAS8 (red) and CANVAS9 (green)]. The non-CANVAS controls
705 are presented in grey. Plots have been divided into PCR-based and PCR-free WGS (left and
706 right columns, respectively). The Y and X axes for ExpansionHunter, GangSTR and
707 TREDPARSE refer to the number of repeat units on the longer and shorter allele per
708 individual, respectively. The Y axis for the STRetch plot refers to the number of individuals.
709

710 **Figure 4: Genetic validation of the (AAGGG)_{exp} RE**

711 A. PCR analysis of the *RFC1* STR failed to produce the control ~253bp reference product in
712 18 of 22 CANVAS families. Representative images of the repeat-primed PCR for the
713 (AAGGG)_{exp} RE demonstrating a saw-toothed product with 5 base pair repeat unit size,
714 amplified from gDNA of individuals from CANVAS1 (B) and CANVAS9 (C). No product
715 was observed for the unaffected control (D) and no gDNA template negative control (E).
716

717 **Figure 5: The majority of individuals with CANVAS encode an ancestral haplotype**

718 A. Analysis of WES data identified an ancestral haplotype surrounding *RFC1* in all affected
719 individuals confirmed to carry the (AAGGG)_{exp} RE. B. The core haplotype (blue highlight)
720 was intersected with the linkage disequilibrium (LD) track in the UCSC browser (converted
721 to hg18 coordinates). The three LD tracks represent the Yoruba population (top track),
722 Europeans (middle) and Han Chinese and Japanese from Tokyo (bottom). Red areas indicate
723 strong linkage disequilibrium. The core CANVAS haplotype spans a large LD block in
724 Europeans, which is broken up into two LD blocks in Japanese and Chinese, suggesting an

725 ancient origin for the CANVAS repeat expansion allele. C. Haplotype sharing between
726 individuals with CANVAS was used to determine the age of the most recent common
727 ancestor (MRCA) of the cohort.

728 **Table 1: Clinical features and genetic analysis of *RFC1* locus in study participants.**

| Family | Participants (sex) | SNP array | WES | WGS | <i>RFC1</i> STR in WES | PCR wildtype allele | Repeat-primed PCR | Genetic Diagnosis | Haplotype | Ethnicity |
|----------|--------------------|-----------|-----|-----|------------------------|---------------------|-------------------|-------------------|--------------|--------------------|
| CANVAS1 | 2 (F) | ✓ | ✓ | ✓ | ND | ✗ | ✓ | CANVAS | A/other | European |
| CANVAS2 | 2 (M) | ✓ | ✓ | ✓ | AAGGG and AAAGG | ✗ | ✓ | CANVAS | A | European |
| CANVAS3 | 2 (F) | ✓ | ✓ | ✗ | AAGGG | ✗ | ✓ | CANVAS | A | European |
| CANVAS4 | 4 (3M,1F) | ✓ | ✓ | ✗ | AAGGG | ✗ | ✓ | CANVAS | A | Greek-Cypriot |
| CANVAS5 | 2 (M,F) | ✗ | ✗ | ✗ | ND | ✗ | ✓ | CANVAS | Not assessed | Not reported |
| CANVAS6 | 2 (M) | ✗ | ✓ | ✗ | AAGGG | ✗ | ✓ | CANVAS | A | Lithuanian/Latvian |
| CANVAS7 | 1 (M) | ✗ | ✓ | ✗ | ND | ✗ | ✓ | CANVAS | A | European-Maori |
| CANVAS8 | 1 (F) | ✗ | ✓ | ✓ | AAGGG and AAAGG | ✗ | ✓ | CANVAS | A/other | European |
| CANVAS9 | 4 (1M,3F) | ✗ | ✓ | ✓ | AAGGG | ✗ | ✓ | CANVAS | A | Lebanese |
| CANVAS10 | 1 (M) | ✗ | ✓ | ✗ | AAGGG | ✗ | ✓ | CANVAS | A | European |
| CANVAS11 | 1 (M) | ✗ | ✓ | ✓ | ND | ✓ | ✗ | ? | NA | Anglo-saxon |
| CANVAS12 | 1 (M) | ✗ | ✓ | ✗ | ND | ✗ | ✓ | CANVAS | A | Turkish |
| CANVAS13 | 1 (M) | ✗ | ✓ | ✓ | Reference | ✓ | ✗ | SCA3 | NA | Martinique |
| CANVAS14 | 1 (M) | ✗ | ✓ | ✗ | AAGGG | ✗ | ✓ | CANVAS | Other* | European |
| CANVAS16 | 1 (F) | ✗ | ✗ | ✗ | NA | ✗ | ✓ | CANVAS | Not assessed | European |
| CANVAS17 | 2 (M) | ✗ | ✓ | ✓ | Reference | ✓ | ✗ | SACS | NA | European |

| | | | | | | | | | | |
|----------|-----------|---|---|---|----|---|---|--------|--------------|----------------|
| CANVAS18 | 1 (F) | x | ✓ | x | ND | x | ✓ | CANVAS | A | European-Maori |
| CANVAS19 | 1 (F) | x | x | ✓ | NA | ✓ | x | SCA45 | Not assessed | European |
| CANVAS20 | 2 (1M,1F) | x | x | x | NA | x | ✓ | CANVAS | Not assessed | Spanish |
| CANVAS21 | 1 (M) | x | x | x | NA | x | ✓ | CANVAS | Not assessed | Indian |
| CANVAS22 | 1 (M) | x | x | x | NA | x | ✓ | CANVAS | Not assessed | Hungarian |
| CANVAS23 | 1 (U) | x | x | x | NA | x | ✓ | CANVAS | Not assessed | Not reported |

729

730 M=male, F=female, U=deidentified, NA=not applicable, ND=not detected, Other*= a different haplotype OR shortened A haplotpye

731 The gene reference sequences utilized were NC_000004 and NM_002913 (RFC1).

732

733

734

735

736 **REFERENCES**

737

- 738 1. McMurray, C.T. (2010). Mechanisms of trinucleotide repeat instability during human development. *Nat Rev*
739 *Genet* 11, 786-799.
- 740 2. Gymrek, M., Willems, T., Guilmatre, A., Zeng, H., Markus, B., Georgiev, S., Daly, M.J., Price, A.L.,
741 Pritchard, J.K., Sharp, A.J., et al. (2016). Abundant contribution of short tandem repeats to gene
742 expression variation in humans. *Nature genetics* 48, 22-29.
- 743 3. Quilez, J., Guilmatre, A., Garg, P., Highnam, G., Gymrek, M., Erlich, Y., Joshi, R.S., Mittelman, D., and
744 Sharp, A.J. (2016). Polymorphic tandem repeats within gene promoters act as modifiers of gene
745 expression and DNA methylation in humans. *Nucleic Acids Res* 44, 3750-3762.
- 746 4. Benson, G. (1999). Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res* 27, 573-
747 580.
- 748 5. Subramanian, S., Madgula, V.M., George, R., Mishra, R.K., Pandit, M.W., Kumar, C.S., and Singh, L.
749 (2003). Triplet repeats in human genome: distribution and their association with genes and other
750 genomic regions. *Bioinformatics* 19, 549-552.
- 751 6. La Spada, A.R., and Taylor, J.P. (2010). Repeat expansion disease: progress and puzzles in disease
752 pathogenesis. *Nat Rev Genet* 11, 247-258.
- 753 7. Hannan, A.J. (2018). Tandem repeats mediating genetic plasticity in health and disease. *Nat Rev Genet* 19,
754 286-298.
- 755 8. Tankard, R.M., Bennett, M.F., Degorski, P., Delatycki, M.B., Lockhart, P.J., and Bahlo, M. (2018). Detecting
756 Expansions of Tandem Repeats in Cohorts Sequenced with Short-Read Sequencing Data. *American*
757 *journal of human genetics* 103, 858-873.
- 758 9. Ruano, L., Melo, C., Silva, M.C., and Coutinho, P. (2014). The global epidemiology of hereditary ataxia and
759 spastic paraplegia: a systematic review of prevalence studies. *Neuroepidemiology* 42, 174-183.
- 760 10. Bird, T.D. (2018 Update). Hereditary Ataxia Overview. In *GeneReviews*(R), M.P. Adam, H.H. Ardinger,
761 R.A. Pagon, S.E. Wallace, L.J.H. Bean, K. Stephens, and A. Amemiya, eds. (Seattle (WA)).
- 762 11. Szmulewicz, D.J., Waterston, J.A., Halmagyi, G.M., Mossman, S., Chancellor, A.M., McLean, C.A., and
763 Storey, E. (2011). Sensory neuropathy as part of the cerebellar ataxia neuropathy vestibular areflexia
764 syndrome. *Neurology* 76, 1903-1910.
- 765 12. Szmulewicz, D.J., McLean, C.A., MacDougall, H.G., Roberts, L., Storey, E., and Halmagyi, G.M. (2014).
766 CANVAS an update: clinical presentation, investigation and management. *J Vestib Res* 24, 465-474.
- 767 13. Harding, A.E. (1981). "Idiopathic" late onset cerebellar ataxia. A clinical and genetic study of 36 cases. *J*
768 *Neurol Sci* 51, 259-271.
- 769 14. Szmulewicz, D.J. (2017). Combined Central and Peripheral Degenerative Vestibular Disorders: CANVAS,
770 Idiopathic Cerebellar Ataxia with Bilateral Vestibulopathy (CABV) and Other Differential Diagnoses
771 of the CABV Phenotype. *Curr Otorhinolaryngol Rep* 5, 167-174.
- 772 15. Cha, Y.H. (2012). Less common neuro-otologic disorders. *Continuum (Minneapolis)* 18, 1142-1157.
- 773 16. Szmulewicz, D.J., Merchant, S.N., and Halmagyi, G.M. (2011). Cerebellar ataxia with neuropathy and
774 bilateral vestibular areflexia syndrome: a histopathologic case report. *Otol Neurotol* 32, e63-65.

- 775 17. Szmulewicz, D.J., McLean, C.A., Rodriguez, M.L., Chancellor, A.M., Mossman, S., Lamont, D., Roberts,
776 L., Storey, E., and Halmagyi, G.M. (2014). Dorsal root ganglionopathy is responsible for the sensory
777 impairment in CANVAS. *Neurology* 82, 1410-1415.
- 778 18. Szmulewicz, D.J., Seiderer, L., Halmagyi, G.M., Storey, E., and Roberts, L. (2015). Neurophysiological
779 evidence for generalized sensory neuronopathy in cerebellar ataxia with neuropathy and bilateral
780 vestibular areflexia syndrome. *Muscle Nerve* 51, 600-603.
- 781 19. Szmulewicz, D.J., Waterston, J.A., MacDougall, H.G., Mossman, S., Chancellor, A.M., McLean, C.A.,
782 Merchant, S., Patrikios, P., Halmagyi, G.M., and Storey, E. (2011). Cerebellar ataxia, neuropathy,
783 vestibular areflexia syndrome (CANVAS): a review of the clinical features and video-oculographic
784 diagnosis. *Ann N Y Acad Sci* 1233, 139-147.
- 785 20. Petersen, J.A., Wichmann, W.W., and Weber, K.P. (2013). The pivotal sign of CANVAS. *Neurology* 81,
786 1642-1643.
- 787 21. Szmulewicz, D., MacDougall, H., Storey, E., Curthoys, I., and Halmagyi, M. (2014). A Novel Quantitative
788 Bedside Test of Balance Function: The Video Visually Enhanced Vestibulo-ocular Reflex (VVOR)
789 *Neurology* 82, S19.002.
- 790 22. Campuzano, V., Montermini, L., Molto, M.D., Pianese, L., Cossee, M., Cavalcanti, F., Monros, E., Rodius,
791 F., Duclos, F., Monticelli, A., et al. (1996). Friedreich's ataxia: autosomal recessive disease caused by
792 an intronic GAA triplet repeat expansion. *Science* 271, 1423-1427.
- 793 23. Taki, M., Nakamura, T., Matsuura, H., Hasegawa, T., Sakaguchi, H., Morita, K., Ishii, R., Mizuta, I., Kasai,
794 T., Mizuno, T., et al. (2018). Cerebellar ataxia with neuropathy and vestibular areflexia syndrome
795 (CANVAS). *Auris Nasus Larynx* 45, 866-870.
- 796 24. Maruta, K., Aoki, M., and Sonoda, Y. (2019). [Cerebellar ataxia with neuropathy and vestibular areflexia
797 syndrome (CANVAS): a case report]. *Rinsho Shinkeigaku* 59, 27-32.
- 798 25. Bahlo, M., Bennett, M.F., Degorski, P., Tankard, R.M., Delatycki, M.B., and Lockhart, P.J. (2018). Recent
799 advances in the detection of repeat expansions with short-read next-generation sequencing. *F1000Res*
800 7.
- 801 26. Seixas, A.I., Loureiro, J.R., Costa, C., Ordonez-Ugalde, A., Marcelino, H., Oliveira, C.L., Loureiro, J.L.,
802 Dhingra, A., Brandao, E., Cruz, V.T., et al. (2017). A Pentanucleotide ATTTTC Repeat Insertion in the
803 Non-coding Region of DAB1, Mapping to SCA37, Causes Spinocerebellar Ataxia. *American journal*
804 *of human genetics* 101, 87-103.
- 805 27. Ishiura, H., Doi, K., Mitsui, J., Yoshimura, J., Matsukawa, M.K., Fujiyama, A., Toyoshima, Y., Kakita, A.,
806 Takahashi, H., Suzuki, Y., et al. (2018). Expansions of intronic TTTCA and TTTTA repeats in benign
807 adult familial myoclonic epilepsy. *Nature genetics* 50, 581-590.
- 808 28. Dolzhenko, E., van Vugt, J., Shaw, R.J., Bekritsky, M.A., van Blitterswijk, M., Narzisi, G., Ajay, S.S.,
809 Rajan, V., Lajoie, B.R., Johnson, N.H., et al. (2017). Detection of long repeat expansions from PCR-
810 free whole-genome sequence data. *Genome research* 27, 1895-1903.
- 811 29. Tang, H., Kirkness, E.F., Lippert, C., Biggs, W.H., Fabani, M., Guzman, E., Ramakrishnan, S., Lavrenko,
812 V., Kakaradov, B., Hou, C., et al. (2017). Profiling of Short-Tandem-Repeat Disease Alleles in 12,632
813 Human Whole Genomes. *American journal of human genetics* 101, 700-715.

- 814 30. Dashnow, H., Lek, M., Phipson, B., Halman, A., Sadedin, S., Lonsdale, A., Davis, M., Lamont, P., Clayton,
815 J.S., Laing, N.G., et al. (2018). STRetch: detecting and discovering pathogenic short tandem repeat
816 expansions. *Genome biology* 19, 121.
- 817 31. Mousavi, N., Shleizer-Burko, S., and Gymrek, M. (2018). Profiling the genome-wide landscape of tandem
818 repeat expansions. *BioRxiv*, <https://doi.org/10.1101/361162>
- 819 32. Cortese, A., Simone, R., Sullivan, R., Vandrovцова, J., Tariq, H., Yan, Y.W., Humphrey, J., Jaunmuktane,
820 Z., Sivakumar, P., Polke, J., et al. (2019). Biallelic expansion of an intronic repeat in RFC1 is a
821 common cause of late-onset ataxia. *Nature genetics* 51, 649-658.
- 822 33. Smith, K.R., Bromhead, C.J., Hildebrand, M.S., Shearer, A.E., Lockhart, P.J., Najmabadi, H., Leventer, R.J.,
823 McGillivray, G., Amor, D.J., Smith, R.J., et al. (2011). Reducing the exome search space for mendelian
824 diseases using genetic linkage analysis of exome genotypes. *Genome biology* 12, R85.
- 825 34. Bahlo, M., and Bromhead, C.J. (2009). Generating linkage mapping files from Affymetrix SNP chip data.
826 *Bioinformatics* 25, 1961-1962.
- 827 35. Abecasis, G.R., Cherny, S.S., Cookson, W.O., and Cardon, L.R. (2002). Merlin--rapid analysis of dense
828 genetic maps using sparse gene flow trees. *Nature genetics* 30, 97-101.
- 829 36. Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features.
830 *Bioinformatics* 26, 841-842.
- 831 37. Pedersen, B.S., Layer, R.M., and Quinlan, A.R. (2016). Vcfanno: fast, flexible annotation of genetic
832 variants. *Genome biology* 17, 118.
- 833 38. Wang, K., Li, M., and Hakonarson, H. (2010). ANNOVAR: functional annotation of genetic variants from
834 high-throughput sequencing data. *Nucleic Acids Res* 38, e164.
- 835 39. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and
836 Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29, 15-21.
- 837 40. Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general purpose program for
838 assigning sequence reads to genomic features. *Bioinformatics* 30, 923-930.
- 839 41. Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., and Smyth, G.K. (2015). limma powers
840 differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* 43,
841 e47.
- 842 42. Gandolfo, L.C., Bahlo, M., and Speed, T.P. (2014). Dating rare mutations from small samples with dense
843 marker data. *Genetics* 197, 1315-1327.
- 844 43. Szmulewicz, D.J., Roberts, L., McLean, C.A., MacDougall, H.G., Halmagyi, G.M., and Storey, E. (2016).
845 Proposed diagnostic criteria for cerebellar ataxia with neuropathy and vestibular areflexia syndrome
846 (CANVAS). *Neurol Clin Pract* 6, 61-68.
- 847 44. Wang, K., Li, M., Hadley, D., Liu, R., Glessner, J., Grant, S.F., Hakonarson, H., and Bucan, M. (2007).
848 PennCNV: an integrated hidden Markov model designed for high-resolution copy number variation
849 detection in whole-genome SNP genotyping data. *Genome research* 17, 1665-1674.
- 850 45. Nibbeling, E.A.R., Duarri, A., Verschuuren-Bemelmans, C.C., Fokkens, M.R., Karjalainen, J.M., Smeets,
851 C., de Boer-Bergsma, J.J., van der Vries, G., Dooijes, D., Bampi, G.B., et al. (2017). Exome
852 sequencing and network analysis identifies shared mechanisms underlying spinocerebellar ataxia. *Brain*
853 140, 2860-2878.

- 854 46. Rinne, T., Bronstein, A.M., Rudge, P., Gresty, M.A., and Luxon, L.M. (1998). Bilateral loss of vestibular
855 function: clinical findings in 53 patients. *J Neurol* 245, 314-321.
- 856 47. Ahmad, H., Requena, T., Frejo, L., Cobo, M., Gallego-Martinez, A., Martin, F., Lopez-Escamez, J.A., and
857 Bronstein, A.M. (2018). Clinical and Functional Characterization of a Missense ELF2 Variant in a
858 CANVAS Family. *Front Genet* 9, 85.
- 859 48. Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T., O'Donnell-Luria, A.H.,
860 Ware, J.S., Hill, A.J., Cummings, B.B., et al. (2016). Analysis of protein-coding genetic variation in
861 60,706 humans. *Nature* 536, 285-291.
- 862 49. Zhang, G., Gibbs, E., Kelman, Z., O'Donnell, M., and Hurwitz, J. (1999). Studies on the interactions
863 between human replication factor C and human proliferating cell nuclear antigen. *Proc Natl Acad Sci U*
864 *S A* 96, 1869-1874.
- 865 50. Yoon, G., and Caldecott, K.W. (2018). Nonsyndromic cerebellar ataxias associated with disorders of DNA
866 single-strand break repair. *Handbook of clinical neurology* 155, 105-115.
- 867 51. Savitsky, K., Bar-Shira, A., Gilad, S., Rotman, G., Ziv, Y., Vanagaite, L., Tagle, D.A., Smith, S., Uziel, T.,
868 Sfez, S., et al. (1995). A single ataxia telangiectasia gene with a product similar to PI-3 kinase. *Science*
869 268, 1749-1753.
- 870 52. Kraus-Perrotta, C., and Lagalwar, S. (2016). Expansion, mosaicism and interruption: mechanisms of the
871 CAG repeat mutation in spinocerebellar ataxia type 1. *Cerebellum Ataxias* 3, 20.
- 872 53. Mateo, I., Llorca, J., Volpini, V., Corral, J., Berciano, J., and Combarros, O. (2003). GAA expansion size
873 and age at onset of Friedreich's ataxia. *Neurology* 61, 274-275.
- 874 54. Bushara, K., Bower, M., Liu, J., McFarland, K.N., Landrian, I., Hutter, D., Teive, H.A., Rasmussen, A.,
875 Mulligan, C.J., and Ashizawa, T. (2013). Expansion of the Spinocerebellar ataxia type 10 (SCA10)
876 repeat in a patient with Sioux Native American ancestry. *PloS one* 8, e81342.
- 877 55. Sato, N., Amino, T., Kobayashi, K., Asakawa, S., Ishiguro, T., Tsunemi, T., Takahashi, M., Matsuura, T.,
878 Flanigan, K.M., Iwasaki, S., et al. (2009). Spinocerebellar ataxia type 31 is associated with "inserted"
879 penta-nucleotide repeats containing (TGGA)n. *American journal of human genetics* 85, 544-557.
- 880 56. Jardim, L.B., Pereira, M.L., Silveira, I., Ferro, A., Sequeiros, J., and Giugliani, R. (2001). Neurologic
881 findings in Machado-Joseph disease: relation with disease duration, subtypes, and (CAG)n. *Arch*
882 *Neurol* 58, 899-904.
- 883 57. Gordon, C.R., Zivotofsky, A.Z., and Caspi, A. (2014). Impaired vestibulo-ocular reflex (VOR) in
884 spinocerebellar ataxia type 3 (SCA3): bedside and search coil evaluation. *J Vestib Res* 24, 351-355.
- 885 58. Gagnon, C., Desrosiers, J., and Mathieu, J. (2004). Autosomal recessive spastic ataxia of Charlevoix-
886 Saguenay: upper extremity aptitudes, functional independence and social participation. *Int J Rehabil*
887 *Res* 27, 253-256.
- 888 59. Vill, K., Muller-Felber, W., Glaser, D., Kuhn, M., Teusch, V., Schreiber, H., Weis, J., Klepper, J.,
889 Schirmacher, A., Blaschek, A., et al. (2018). SACS variants are a relevant cause of autosomal recessive
890 hereditary motor and sensory neuropathy. *Human genetics* 137, 911-919.
- 891 60. Richards, S., Aziz, N., Bale, S., Bick, D., Das, S., Gastier-Foster, J., Grody, W.W., Hegde, M., Lyon, E.,
892 Spector, E., et al. (2015). Standards and guidelines for the interpretation of sequence variants: a joint

- 893 consensus recommendation of the American College of Medical Genetics and Genomics and the
894 Association for Molecular Pathology. *Genet Med* 17, 405-424.
- 895 61. Sullivan, R., Yau, W.Y., O'Connor, E., and Houlden, H. (2019). Spinocerebellar ataxia: an update. *J Neurol*
896 266, 533-544.
- 897 62. Galatolo, D., Tessa, A., Filla, A., and Santorelli, F.M. (2018). Clinical application of next generation
898 sequencing in hereditary spinocerebellar ataxia: increasing the diagnostic yield and broadening the
899 ataxia-spasticity spectrum. A retrospective analysis. *Neurogenetics* 19, 1-8.
- 900

CANVAS study overview

Cohort collection

22 families with CANVAS:
11 sporadic cases
7 affected sibling pairs
4 multigeneration affected families

Linkage analysis

SNP chip: 4 families (CANVAS1,2,3,4)
WES: 1 family (CANVAS9)

Identify **homozygous** single overlapping linkage region:
chr4:38941465-40390306

Whole exome sequencing (WES) - large collaboration with CIDR

23 affected individuals from 15 families
No shared rare or de novo variants detected

Whole genome sequencing (WGS)

Two unrelated individuals with CANVAS
No shared rare or de novo variants detected

Identify novel RE expansion: homozygous inheritance of rare AAGGG intronic RE (chr4:39350045) in the gene *RFC1* - within the chr4 linkage region.

Validation by repeat primed PCR

Confirm homozygous AAGGG inheritance in 18 of 22 CANVAS families
4 families negative for AAGGG RE - prioritised for further WGS

Re-analysis of WES (CANVAS9 and CIDR)

Off target reads in WES protocol
Single read coverage at chr4:39350045 in 14 individuals

AAGGG expansion detected in 11 patients, from 9 different families

3 patients (2 families) contain evidence for the reference genome
2 patients (2 families) contain evidence for both AAGGG and AAAGG

WGS round 2

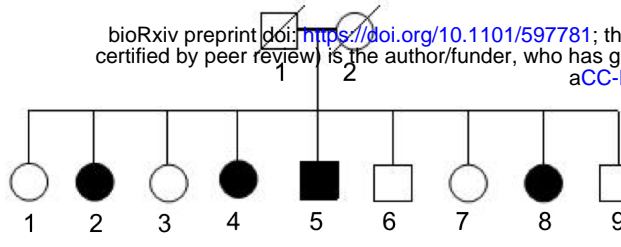
WGS: 5 patients negative for AAGGG in *RFC1*
- genomic re-diagnosis:
- SCA3, SACS (compound heterozygous),
- SCA45 (point mutation in *FAT2*)
- VOUS: point mutation in *ATXN7*

WGS: 2 patients with potential AAAGG/AAGGG RE
- CANVAS2-2: confirm AAGGG RE on both alleles
- **CANVAS8-8: confirm AAGGG on one allele, and AAAGG on second allele**

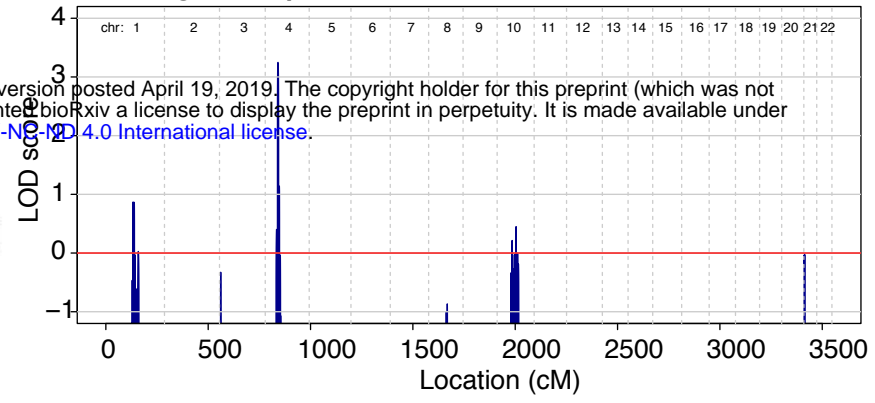
A CANVAS9

bioRxiv preprint doi: <https://doi.org/10.1101/597781>; this version posted April 19, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

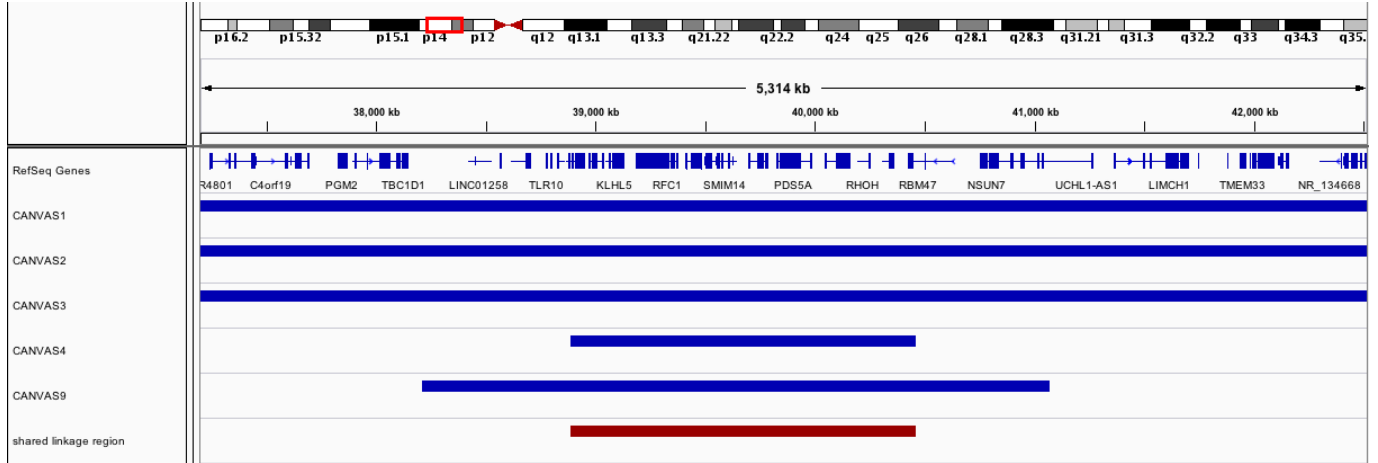
II



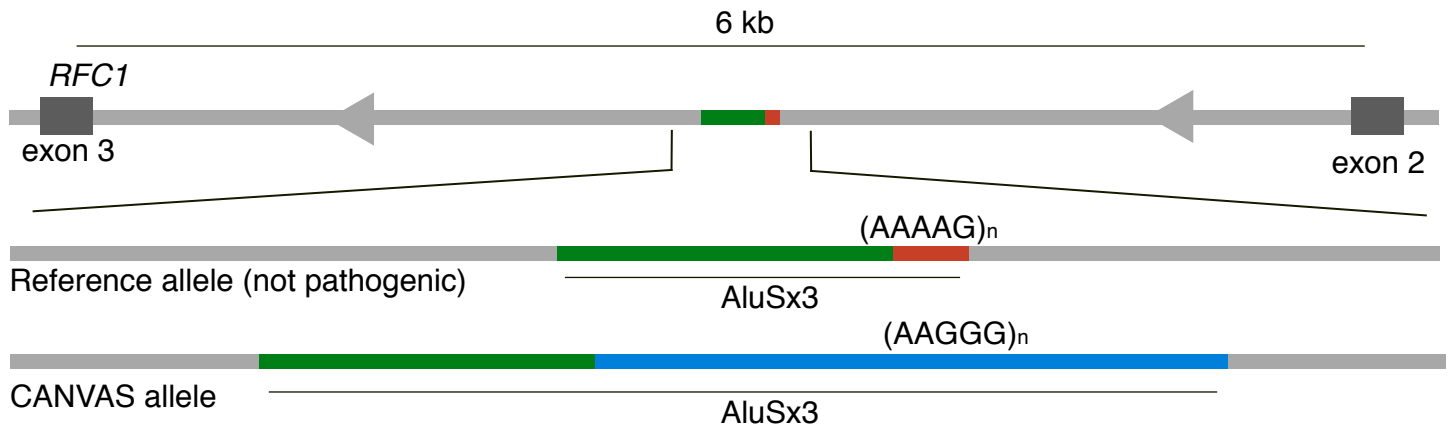
B genome parametric LOD score for CANVAS9



C

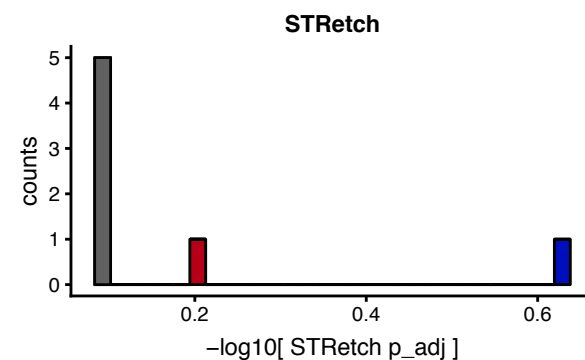
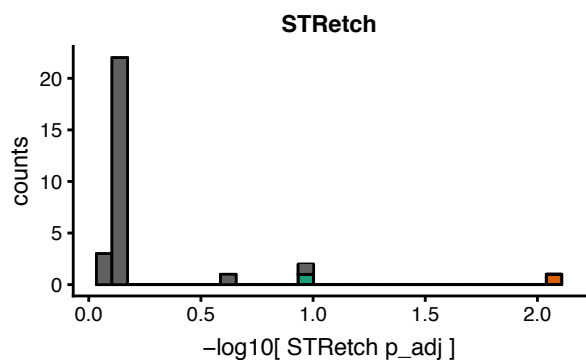
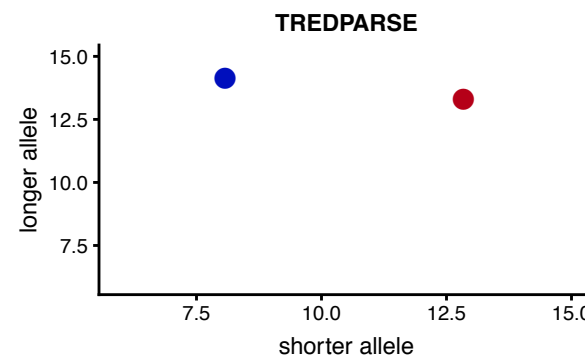
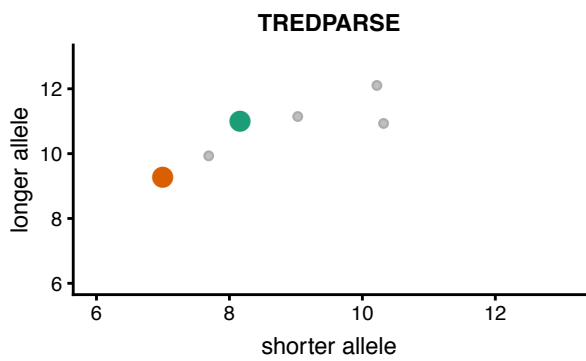
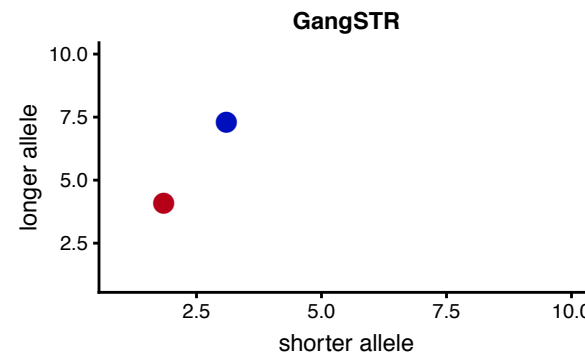
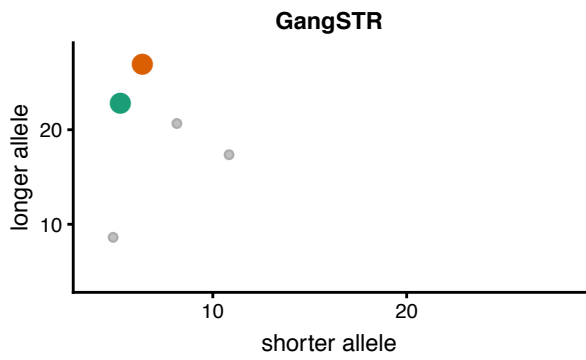
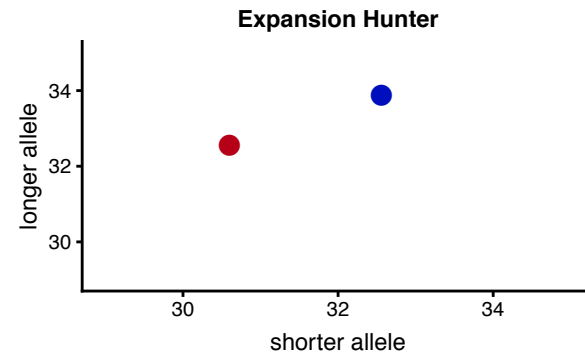
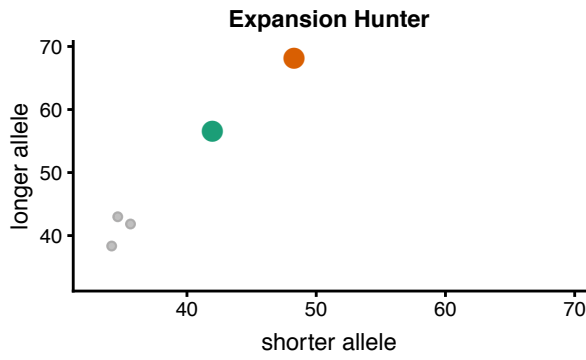
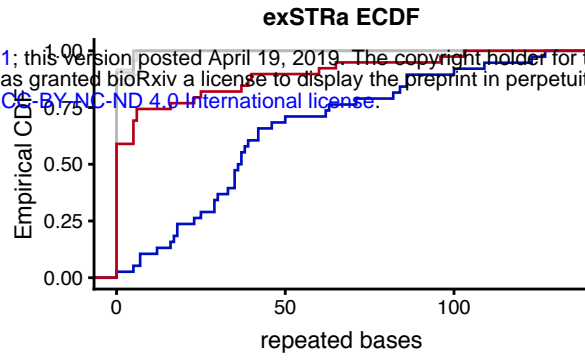
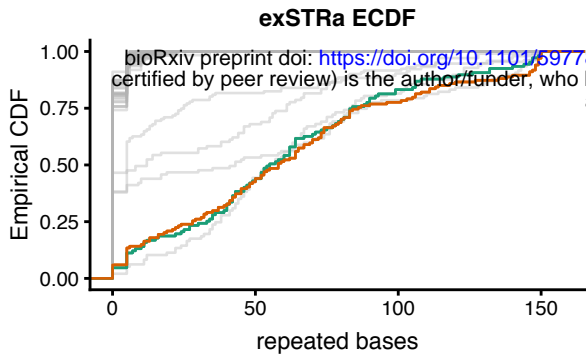


D



PCR-based WGS

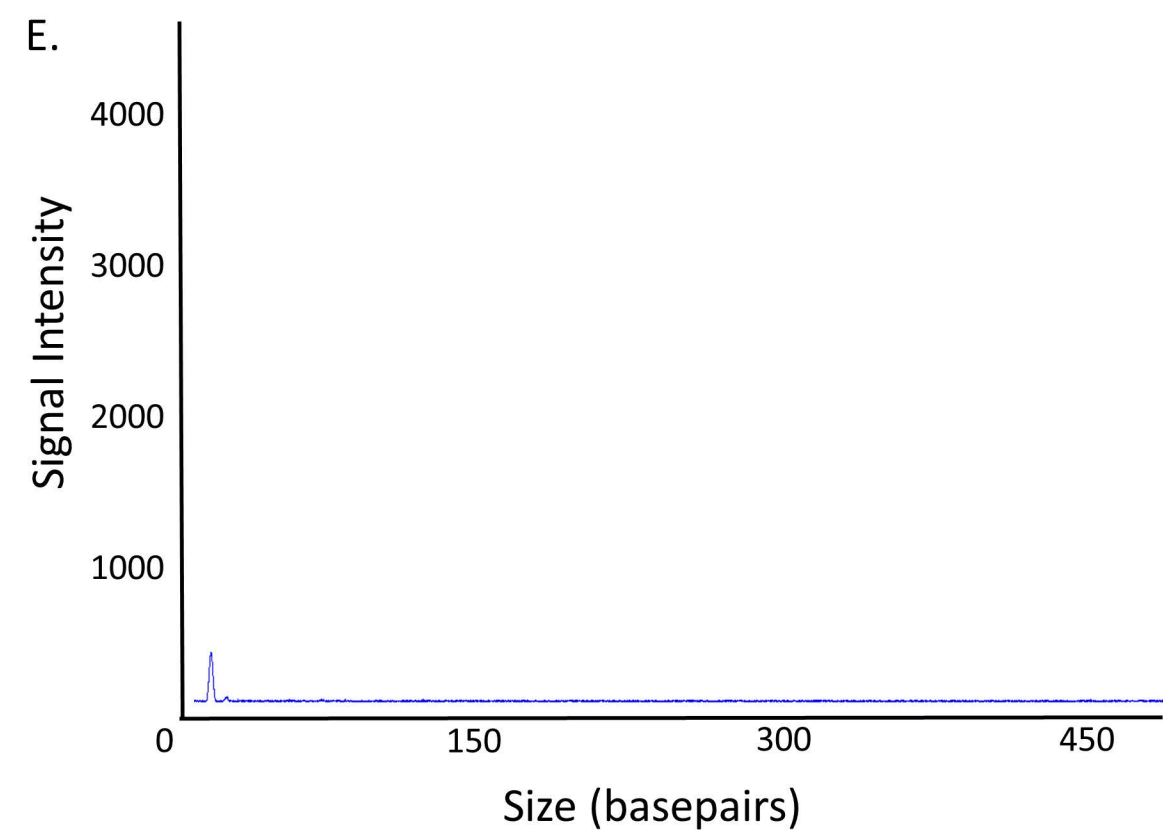
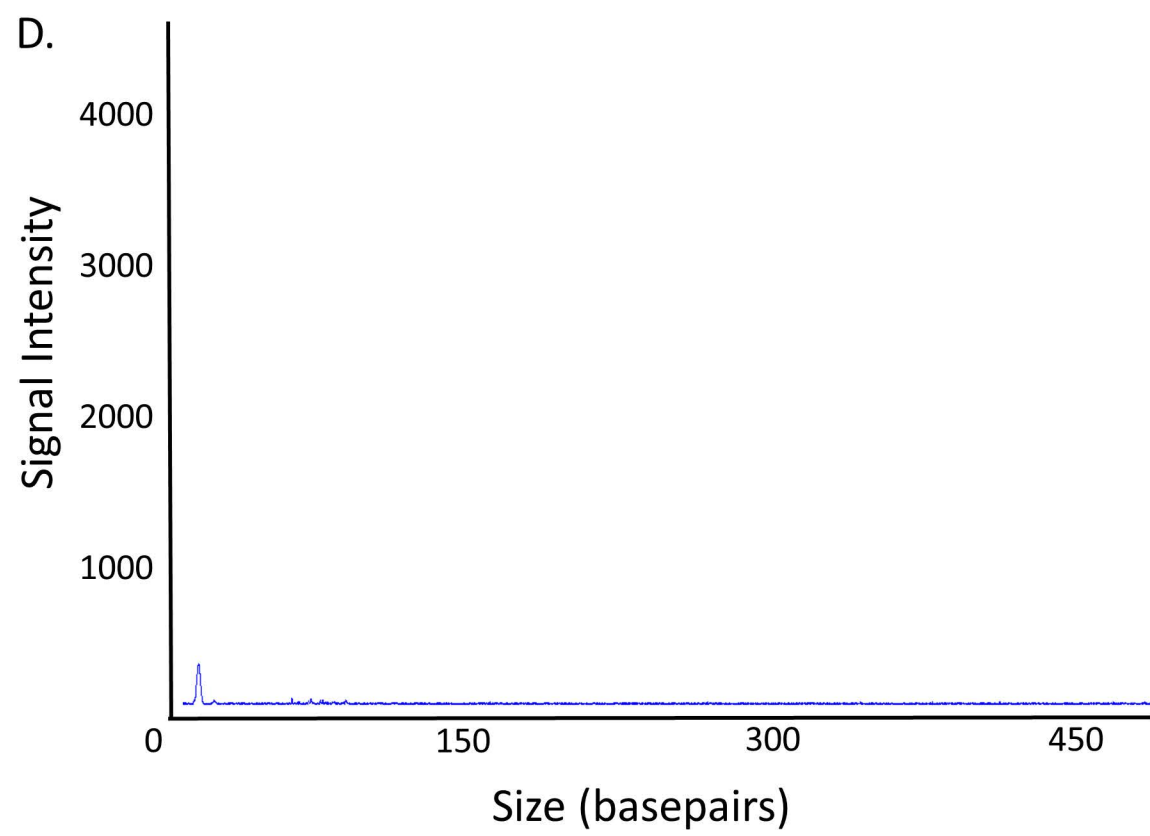
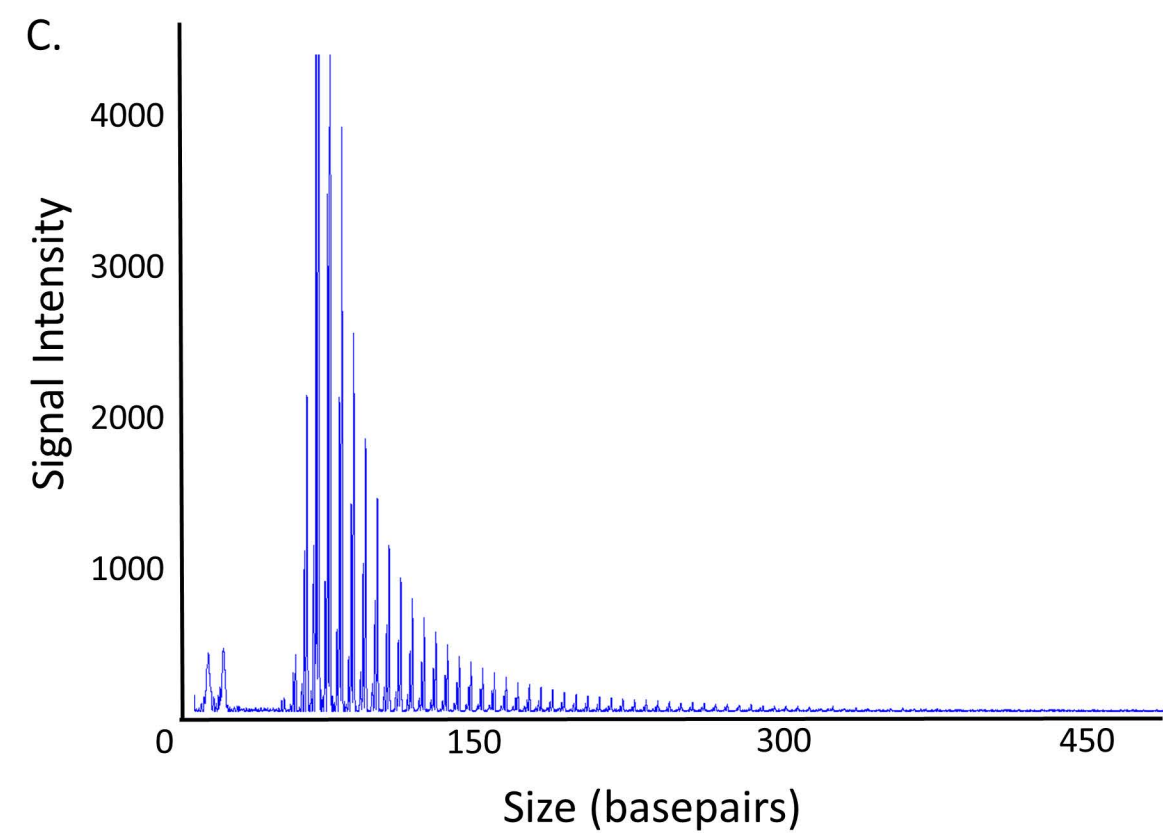
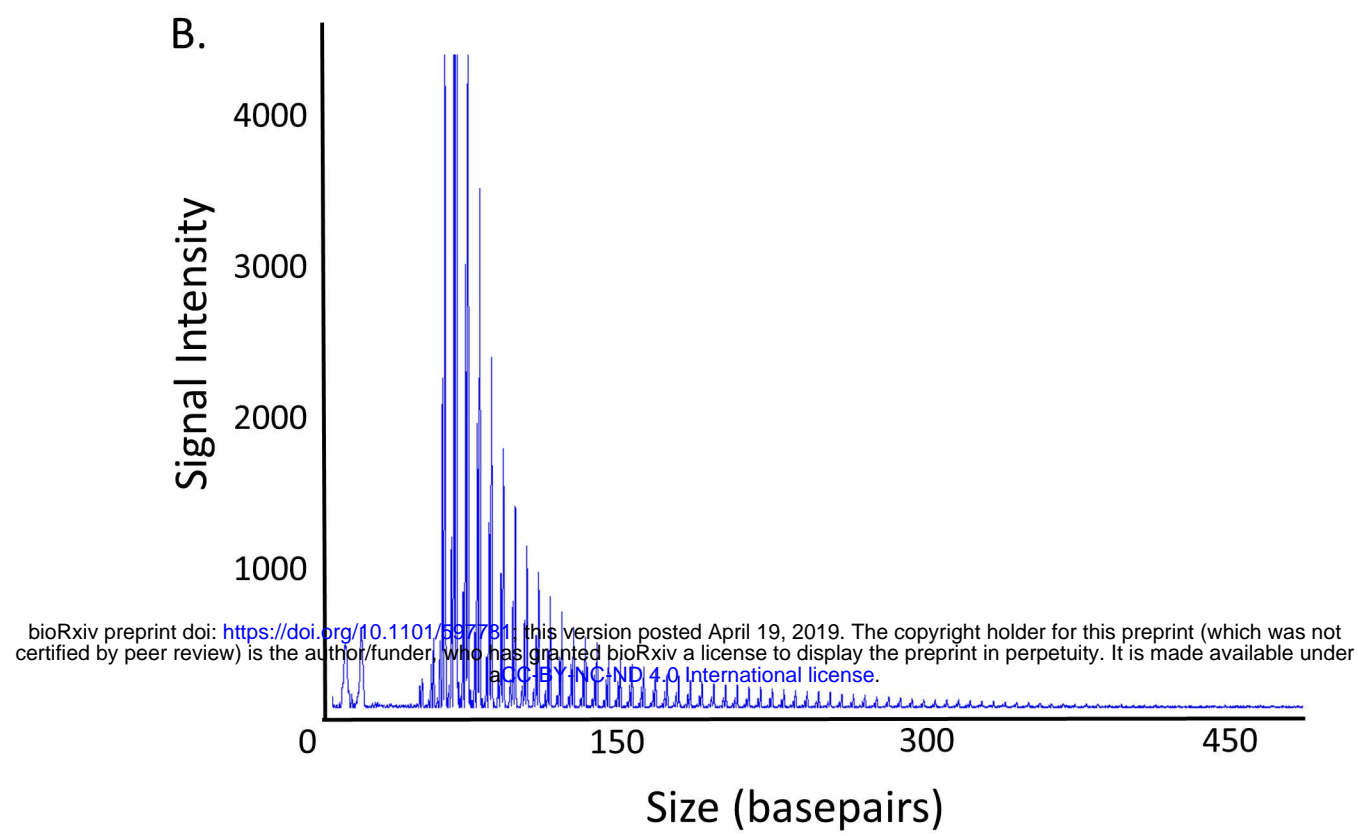
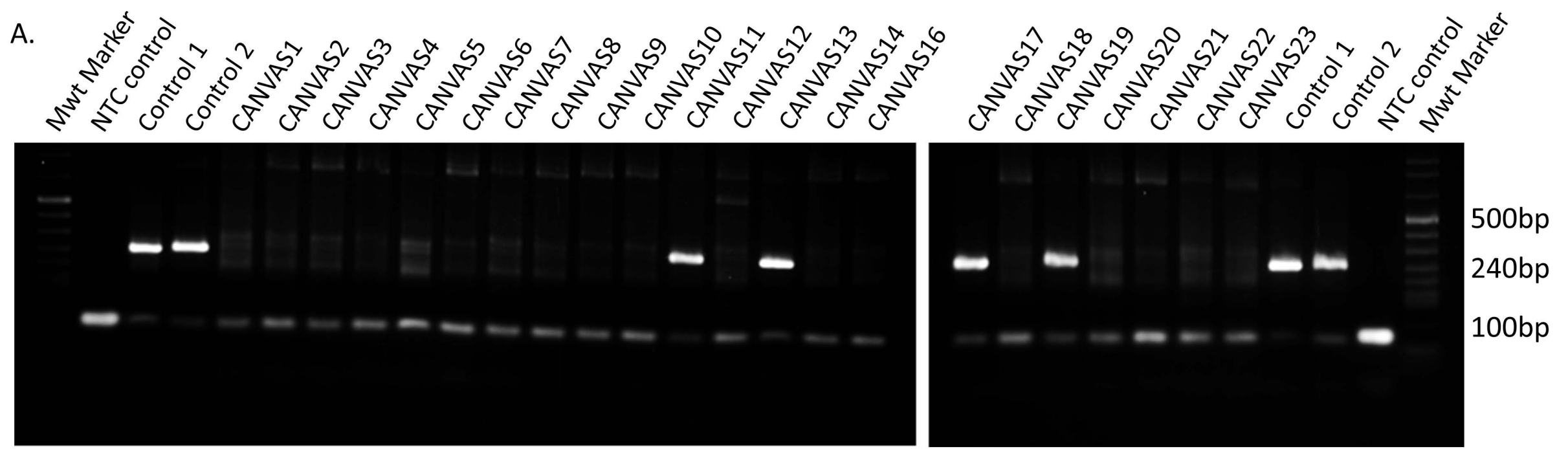
PCR-free WGS



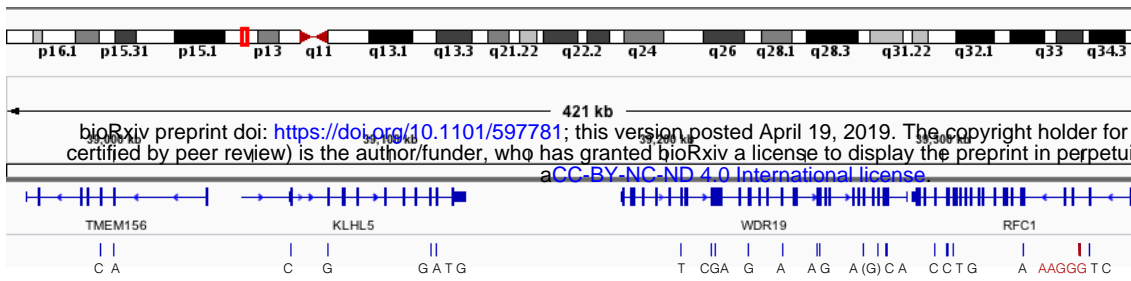
- CANVAS9
- CANVAS1
- non-CANVAS

- CANVAS2
- CANVAS8
- non-CANVAS

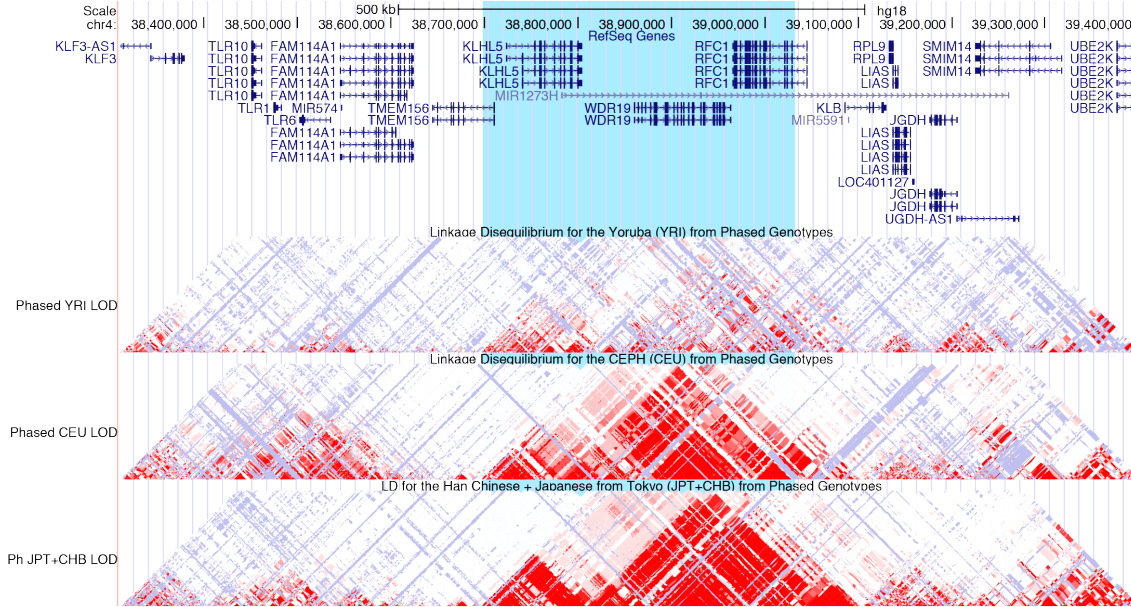
bioRxiv preprint doi: <https://doi.org/10.1101/597781>; this version posted April 19, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



A



B



C

