

# **Validation of new tools to identify expanded repeats: an intronic pentamer expansion in *RFC1* causes CANVAS**

Haloom Rafehi<sup>1,3^</sup>, David J Szmulewicz<sup>2^</sup>, Mark F Bennett<sup>1,3,4</sup>, Nara LM Sobreira<sup>5</sup>, Kate Pope<sup>6</sup>, Katherine R Smith<sup>7</sup>, Greta Gillies<sup>6</sup>, Peter Diakumis<sup>8</sup>, Egor Dolzhenko<sup>9</sup>, Mike Eberle<sup>9</sup>, María García Barcina<sup>10</sup>, David P Breen<sup>11,12,13</sup>, Andrew M Chancellor<sup>14</sup>, Phillip D Cremer<sup>15,16</sup>, Martin B. Delatycki<sup>6,17</sup>, Brent L Fogel<sup>18</sup>, Anna Hackett<sup>19,20</sup>, G. Michael Halmagyi<sup>21,22</sup>, Solange Kapetanovic<sup>23</sup>, Anthony Lang<sup>24,25</sup>, Stuart Mossman<sup>26</sup>, Weiyi Mu<sup>5</sup>, Peter Patrikios<sup>27</sup>, Susan L Perlman<sup>28</sup>, Ian Rosemargy<sup>29</sup>, Elsdon Storey<sup>30</sup>, Shaun RD Watson<sup>31</sup>, Michael A Wilson<sup>6</sup>, David Zee<sup>32</sup>, David Valle<sup>5</sup>, David J Amor<sup>6,17</sup>, Melanie Bahlo<sup>1,3\*</sup> and Paul J Lockhart<sup>6,17\*,#</sup>

<sup>1</sup> Population Health and Immunity Division, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, Victoria 3052, Australia

<sup>2</sup> Cerebellar Ataxia Clinic, Neuroscience Department, Alfred Health, Melbourne, Australia

<sup>3</sup> Department of Medical Biology, University of Melbourne, 1G Royal Parade, Parkville, Victoria 3052, Australia

<sup>4</sup> Epilepsy Research Centre, Department of Medicine, University of Melbourne, Austin Health, 245 Burgundy Street, Heidelberg, Victoria 3084, Australia

<sup>5</sup> McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, 21205, USA

<sup>6</sup> Bruce Lefroy Centre, Murdoch Children's Research Institute, Flemington Rd, Parkville, Victoria 3052, Australia

<sup>7</sup> Placeholder, Parkville, VIC, Australia

<sup>8</sup> University of Melbourne Centre for Cancer Research, Victorian Comprehensive Cancer Centre, 305 Grattan Street, Melbourne VIC 3000, Australia

<sup>9</sup> Illumina Inc, 5200 Illumina Way, San Diego, CA 92122, USA

<sup>10</sup> Genetic Unit, Basurto University Hospital, OSI Bilbao-Basurto, avenida Montevideo 18 (48013 Bilbao), Spain

<sup>11</sup> Centre for Clinical Brain Sciences, University of Edinburgh, Edinburgh, Scotland

<sup>12</sup> Anne Rowling Regenerative Neurology Clinic, University of Edinburgh, Edinburgh, Scotland

<sup>13</sup> Centre for Medical Informatics, Usher Institute of Population Health Sciences and Informatics, University of Edinburgh, Edinburgh, UK

<sup>14</sup> Department of Neurology, Tauranga Hospital, Private Bag, Cameron Road, Tauranga 3171, New Zealand

<sup>15</sup> University of Sydney, NSW 2006

<sup>16</sup> Royal North Shore Hospital, Pacific Hwy, St Leonards NSW 2065, Australia

<sup>17</sup> Department of Paediatrics, University of Melbourne, Royal Children's Hospital, Flemington Rd, Parkville, Victoria 3052, Australia

<sup>18</sup> Departments of Neurology and Human Genetics, David Geffen School of Medicine University of California, Los Angeles (UCLA)

<sup>19</sup> Hunter Genetics, Hunter New England Health Service, Waratah, Newcastle, NSW, 2300, Australia

<sup>20</sup> University of Newcastle, Newcastle, NSW, 2300, Australia

<sup>21</sup> Neurology Department, Royal Prince Alfred Hospital Sydney

<sup>22</sup> Central Clinical School, University of Sydney

<sup>23</sup> Servicio de Neurología, Hospital de Basurto, Bilbao, Vizcaya, España

<sup>24</sup> Edmond J. Safra Program in Parkinson's disease and the Morton and Gloria Shulman  
Movement Disorders Clinic, Toronto Western Hospital

<sup>25</sup> Department of Medicine, Division of Neurology, University Health Network and the  
University of Toronto

<sup>26</sup> Department of Neurology, Wellington Hospital, Wellington 6021, New Zealand

<sup>27</sup> Sunshine Neurology, Maroochydore QLD 4558, Australia

<sup>28</sup> Department of Neurology, David Geffen School of Medicine, University of California, Los  
Angeles, CA 90095, USA

<sup>29</sup> Riddiford Medical, Newtown, Wellington 6023 New Zealand

<sup>30</sup> Department of Neuroscience, Central Clinical School, Monash University, Alfred Hospital  
Campus, Commercial Road, Melbourne, Vic 3004, Australia

<sup>31</sup> Institute of Neurological Sciences, Prince of Wales Hospital, Barker St, Randwick NSW  
2031, Australia

<sup>32</sup> Department of Neurology, Johns Hopkins Hospital, Baltimore, MD 21287 USA

<sup>33</sup> McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of  
Medicine, Baltimore, MD, 21205

<sup>^</sup> These authors contributed equally to this work

<sup>\*</sup>These authors contributed equally to this work

<sup>#</sup>Correspondence to A/Prof Paul Lockhart, Murdoch Children's Research Institute,  
Flemington Rd, Parkville, Victoria 3052, Australia; [paul.lockhart@mcri.edu.au](mailto:paul.lockhart@mcri.edu.au)

Running title: An intronic repeat expansion in RFC1 causes CANVAS

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

## ABSTRACT

Genomic technologies such as Next Generation Sequencing (NGS) are revolutionizing molecular diagnostics and clinical medicine. However, these approaches have proven inefficient at identifying pathogenic repeat expansions. Several new tools can interrogate the catalogue of known short tandem repeat (STR) loci to identify disease-causing expansions but these are limited to detecting expansion of previously defined STRs. Here, we describe a reference-free method called Expansion Hunter De Novo (EHdn), which can be utilized to identify either known or novel expanded repeat sequences in NGS data. We performed genetic studies of a cohort of 35 individuals from 22 families with a clinical diagnosis of cerebellar ataxia with neuropathy and bilateral vestibular areflexia syndrome (CANVAS). Analysis of whole genome sequence (WGS) data with EHdn identified a recessively inherited intronic (AAGGG)<sub>n</sub> repeat expansion in the gene encoding Replication Factor C1 (*RFC1*). This motif, not reported in the reference sequence, localized to an Alu element and replaced the reference (AAAAG)<sub>11</sub> short tandem repeat. Genetic analyses confirmed the pathogenic expansion in 18 of 22 CANVAS families and identified a core ancestral haplotype, estimated to have arisen in Europe over twenty five thousand years ago. WGS of the four *RFC1* negative families identified plausible variants in three, with genomic re-diagnosis of SCA3, spastic ataxia of the Charlevoix-Saguenay type and SCA45. This study identified the genetic basis of CANVAS and demonstrated that these improved bioinformatics tools increase the diagnostic utility of WGS to determine the genetic basis of a heterogeneous group of clinically overlapping neurogenetic disorders.

# INTRODUCTION

Repetitive DNA sequences constitute approximately one third of the genome and are thought to contribute to diversity within and between species.<sup>1</sup> Microsatellites or short tandem repeats (STRs) are mini-repeats of DNA, typically two to five base-pairs in length, which are usually present in a concatamer of between five and fifty repeated elements. There are thousands of STRs scattered through the human genome and recent studies have suggested important roles for STRs in the regulation of gene expression.<sup>2; 3</sup> STRs display considerable variability in length between individuals, which is presumed to have no detrimental consequences for humans<sup>4; 5</sup> unless the repeat number is expanded beyond a gene-specific threshold.<sup>6; 7</sup> Pathogenic repeat expansions (REs) have been shown to underlie at least 30 inherited human diseases, the majority being disorders of the nervous system.<sup>8</sup> These disorders, which variably have autosomal dominant, autosomal recessive and X-linked inheritance, have an overall prevalence of ~1:20,000.<sup>9</sup> They display a broad onset age and are characterized by progressive cerebellar ataxia with dysarthria, oculomotor abnormalities, cognitive dysfunction and other symptoms.<sup>10</sup> Additional novel pathogenic REs likely remain to be identified. For example, putative spinocerebellar ataxia (SCA) loci, including SCA25 (MIM: 608703) and SCA30 (MIM: 613371) remain to be identified, and unsolved hereditary ataxias such as cerebellar ataxia with neuropathy and bilateral vestibular areflexia syndrome (CANVAS, MIM: 614575) display extensive clinical similarities with known RE disorders.

CANVAS is a cerebellar ataxia with combined cerebellar, vestibular and somatosensory dysfunction.<sup>11; 12</sup> Historically, individuals with CANVAS have been assigned the diagnosis of idiopathic late onset cerebellar ataxia.<sup>13</sup> More recently, CANVAS is clinically recognized and has been incorporated into the contemporary research and teaching of both cerebellar and vestibular diseases.<sup>14; 15</sup> Unifying the oto- and neuropathology, CANVAS is a neuronopathy

(ganglionopathy) affecting the vestibular<sup>16</sup> and dorsal root ganglia.<sup>17</sup> The progression of these clinical features can be measured longitudinally using a specific neurophysiological protocol.<sup>18</sup> A characteristic radiological pattern of cerebellar atrophy has also been described and verified on post-mortem pathology.<sup>11</sup> The characteristic oculomotor abnormality seen in combined cerebellar and vestibular impairment is the visually-enhanced vestibulo-ocular reflex (VVOR), and this can now be evaluated using a commercially available instrumented assessment tool.<sup>19-21</sup> Altogether, these advances have allowed the formulation of diagnostic criteria to aid identification of CANVAS, contributing both research and clinical benefits including improved prognostication and targeted management.<sup>12; 14</sup> While detailed clinical findings have driven gene discovery in RE disorders such as Friedreich ataxia<sup>22</sup> the underlying genetic cause(s) of CANVAS has, until very recently, remained elusive (see below).

The majority of individuals and families with CANVAS have been reported in individuals from European populations, although CANVAS has recently been reported in two individuals of Japanese ethnicity, a 68-year old male<sup>23</sup> and a 76 year old female.<sup>24</sup> A genetic cause of CANVAS is highly plausible given the observation of 13 affected siblings and families with multiple affected individuals over several generations.<sup>12</sup> The pattern of inheritance suggests an autosomal recessive trait, although autosomal dominant inheritance with incomplete penetrance cannot be excluded. CANVAS symptoms overlap considerably with SCA3 (also known as Machado-Joseph disease, MIM: 109150) and Friedreich ataxia (MIM: 229300), both genetic forms of ataxia caused by the inheritance of a pathogenic RE. These observations are consistent with the hypothesis that a novel pathogenic STR expansion may underlie CANVAS.

Historically, the detection of REs has been time-consuming and expensive. Indeed, it is only in recent years that new computational methods have been developed to screen for RE in short-read whole exome sequence (WES) and WGS data<sup>25</sup>, leading to the discovery of novel, disease causing REs. For example, a pentanucleotide RE was identified to underlie autosomal dominant spinocerebellar ataxia 37 (SCA37; OMIM: 615945).<sup>26</sup> Moreover, pathogenic REs of intronic pentamers (TTTCA)<sub>n</sub> and (TTTTA)<sub>n</sub> were identified as the cause of Benign Adult Familial Myoclonus Epilepsy locus 1, 6 and 7 (BAFME1, OMIM: 618073; BAFME6, OMIM: 618074; BAFME1, OMIM: 618075).<sup>27</sup>

Multiple tools now exist that allow screening of short-read sequencing data for expanded STRs.<sup>25</sup> Initially, STR detection tools such as lobSTR and hipSTR were limited to short STRs that were encompassed by a single read. However, in the last two years, multiple methods have been released that can screen WES and WGS datasets for REs without being limited by read length. These include ExpansionHunter (EH)<sup>28</sup>, exSTRa<sup>8</sup>, TREDPARSE<sup>29</sup>, STRetch<sup>30</sup> and GangSTR.<sup>31</sup> These are all reference based methods - i.e. they rely on a catalogue of STR loci and motifs and are therefore limited to detecting expansion of previously defined STRs, such as those catalogued in the UCSC track. Here, we describe a reference-free method called Expansion Hunter De Novo (EHdn) and show that CANVAS is caused by the homozygous inheritance of an expanded intronic (AAGGG)<sub>n</sub> pentamer in the gene encoding *RFC1*. Our workflow is summarized in Figure 1. A parallel study published very recently similarly identified the causal pentamer in *RFC1*.<sup>32</sup> Cortese and colleagues defined a small linkage region from ten families with CANVAS. WGS was performed on six unrelated individuals and the causative RE was identified by visual inspection of the aligned read pairs inside the linkage region.



# MATERIALS AND METHODS

## Recruitment, linkage and next generation sequence data

The Royal Children's Hospital Human Research Ethics Committee approved the study (HREC 28097). Informed consent was obtained from all participants and clinical details were collected from clinical assessments and review of medical records. Genomic DNA was isolated from peripheral blood. Single nucleotide polymorphism (SNP) genotype data were generated for two affected siblings from three families (CANVAS1, 2, 3) and all six siblings from family CANVAS4 using the Illumina Infinium HumanOmniExpress BeadChip genotyping array. SNP genotypes for individuals from CANVAS9 were extracted from WES data.<sup>33</sup> Parametric multipoint linkage analysis was subsequently performed using LINKDATAGEN and MERLIN<sup>34; 35</sup> specifying a rare recessive disease model with complete penetrance, and overlapping linkage signals were detected using BEDtools.<sup>36</sup> WES was performed on individuals from CANVAS9 using Agilent SureSelect XT Human All exon V5 + UTR on the Illumina HiSeq2000 platform at 50x mean coverage. WES was performed on an additional 23 individuals from 15 families in collaboration with the Johns Hopkins Center for Inherited Disease Research (CIDR) as part of the Baylor-Hopkins Center for Mendelian Genomics (BHCMG). WGS was performed in two stages. Libraries for the first round of samples, including two affected individuals from CANVAS1 and CANVAS9 and 31 individuals lacking a clinical diagnosis of CANVAS (subsequently referred to as controls although some have a diagnosis other than CANVAS), were prepared using the TruSeq nano PCR-based Library Preparation Kit and sequenced on the Illumina HiSeq X platform. Libraries for the second round of WGS, including affected individuals with evidence of an alternate RE motif (CANVAS2 and CANVAS8) or lacking the pathogenic RE in *RFC1* (CANVAS11,13, 17 and 19), were prepared using the TruSeq PCR-free DNA HT Library Preparation Kit and sequenced on the Illumina NovaSeq 6000 platform. PCR-free WGS data

from 69 unrelated Coriell controls<sup>28</sup> was obtained from Illumina. GTEx samples (SRA files, 133 WGS with matching cerebellar RNA-seq) were downloaded from the dbGAP (phs000424.v7.p2).

## Alignment and variant calling

Alignment and haplotype calling were performed based on the GATK best practice pipeline. All WES and WGS datasets were aligned to the hg19 reference genome using BWA-mem, then duplicate marking, local realignment and recalibration were performed with GATK. Merged VCF files were annotated using vcfanno<sup>37</sup> and ANNOVAR.<sup>38</sup> Candidate variant filtering was performed using CAVALIER, an R package for variant interpretation in NGS data (<https://github.com/bahlolab/cavalier>). Standard variant calling was performed on WGS data for CANVAS samples negative for the pathogenic RE in *RFC1*. Candidate variants were defined as i) occurring in known ataxia genes, as defined by OMIM, ii) exonic, with a minor allele frequency of less than 0.0001 in gnomAD (both genome and exome data) and iii) predicted pathogenic by both SIFT and PolyPhen2. RNA-seq data was aligned to the hg19 reference genome (ENSEMBL Homo\_sapiens.GRCh37.75) using STAR.<sup>39</sup> Reads were summarized by gene ID into a counts matrix using featureCounts<sup>40</sup> (quality score  $\geq 10$ ) and converted to log10 of the counts per million using limma.<sup>41</sup>

## STR analysis

Genome-wide screening for putative REs was performed using Expansion Hunter Denovo (EHdn) version 0.6.2. EHdn operates by performing a genome-wide search for read pairs where one mate has confident alignment (anchor) and the second mate consists of repetition of a repeat motif (in-repeat read). The program reports the counts of in-repeat reads with anchor mates stratified by the repeat motif and genomic position of their anchor mate.

For this analysis we defined a confidently-aligned read as one aligned with MAPQ of 50 or above. The counts of in-repeat reads with anchor mates were subsequently compared for each region in cases (CANVAS) and controls using a permutation test ( $10^6$  permutations). The resulting p-values were used to rank candidate sites with higher counts in individuals with CANVAS than in the controls for further computational validation. These candidates were subsequently annotated with ANNOVAR.

Validation was performed using all available reference-based STR detection tools for short-read NGS. The RE candidates were screened in the two individuals with CANVAS and the 31 non-CANVAS controls using exSTRa and EH, then the top candidate [(AAGGG)<sub>n</sub> STR in *RFC1*] was further validated with TREDPARSE, GangSTR and STRetch. All tools were used with default parameters, with the following additional parameters for EH: read-depth of 30 and min-anchor-mapq of 20. All five tools were also used to screen for the (AAGGG)<sub>n</sub> *RFC1* STR in the 69 Coriell control WGS datasets. A short-list of AAGGG carriers was generated based on consensus calling from at least four of the five tools.

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

## Haplotyping and mutation dating

Haplotyping was performed on the WES data. Variants were filtered based on read depth ( $\geq 30$ ), including both exonic and non-exonic variants. A core haplotype was defined based on sharing amongst a majority of affected individuals. A method based on haplotype sharing<sup>42</sup> was used to determine the most recent common ancestor (MRCA) from whom the core haplotype was inherited, as well as dating additional sub-haplotypes shared by clusters of individuals, which are likely to be individuals with a MRCA who is more recent than that for the whole group (<https://shiny.wehi.edu.au/rafehi.h/mutation-dating/>).

## Molecular genetic studies

We designed a PCR assay to test for presence of the non-expanded reference allele at the *RFC1* STR. The primers (Table S1) flank the STR and amplify a 253bp fragment using standard PCR conditions with a 30 second extension cycle. The pathogenic *RFC1* RE was amplified by repeat-primed PCR with three primers; TPP\_CANVAS\_FAM\_2F, 5R\_TPP\_M13R\_CANVAS\_RE\_R and TPP\_M13R (Table S1). The FAM labelled forward primer is locus specific, while the repeat-specific primer (5R\_TPP) includes a tag M13R sequence. PCR was performed in a 20 $\mu$ l reaction with 20 ng genomic DNA, 0.8  $\mu$ M of both the FAM labelled forward primer and TPP\_M13R and 0.2  $\mu$ M 5R\_TPP using GoTaq® Long PCR Polymerase (Promega). A standard 60TD55 protocol was utilized (94°C denaturation for 30 s, 60TD55°C anneal for 30 s, and 72°C extension for 2 min), products were detected on an ABI3730xl DNA Analyzer and visualized using PeakScanner 2 (Applied Biosystems).

## RESULTS

### Case recruitment

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

### Linkage analysis

CANVAS typically presents in families with one or multiple affected individuals in a single generation, consistent with a recessive inheritance. For example, in the second-degree consanguineous family CANVAS9, four siblings were diagnosed with CANVAS and two were classified as unaffected at the time phenotyping was performed (Figure 2A). Linkage analysis was initially performed on five CANVAS families (CANVAS1, -2, -3, -4 and -9, Figure S1) using a nonparametric model, which suggested recessive inheritance (results not shown). Additional linkage analysis was performed on the five pedigrees with a parametric recessive inheritance model. This identified linkage regions with logarithm of odds (LOD) scores ranging from 0.6 for smaller pedigrees (two affected siblings), to a statistically significant linkage region on chromosome 4 in CANVAS9 (LOD=3.25, Figure 2B). Intersection of the linkage regions from the five families identified a single region on

chromosome four (chr4:38887351-40463592, hg19) common to all families (Figure 2C). The 1.5MB shared region contains 42 genes, of which 14 are protein coding, none with any association with ataxia in OMIM or the published literature (Table S2).

### **Large-scale WES analysis did not identify candidate pathogenic variants**

WES was used to screen 27 affected individuals with CANVAS from 15 families for potentially pathogenic rare variants (MAF < 0.001) shared across multiple pedigrees in a homozygous or compound heterozygous inheritance pattern. No candidate mutations were detected, either within the chromosome 4 linkage region or elsewhere in the genome.

### **Identification of a novel (AAGGG)<sub>n</sub> RE in the linkage region**

The lack of candidate variants identified from the WES data suggested the possibility of (i) intronic or intergenic mutations, or (ii) that CANVAS might be caused by a non-standard mutation, such as a pathogenic RE of an STR. Therefore, WGS was performed on two individuals from different pedigrees (CANVAS1 and CANVAS9) who share the chr4 linkage region and 31 unrelated controls. EHdn was used to perform a genome-wide screen for STRs in the two individuals with CANVAS compared to the controls. This identified 19 regions with a p value < 0.005 (Table S3), although genome-wide significance could not be achieved after adjustment for multiple testing due to the skewed ratio of the number of cases to controls (2 versus 31). These candidate STRs were visualized with the Integrative Genomics Viewer (IGV) tool, which suggested that only the candidate (AAGGG)<sub>n</sub> STR within intron 2 of the gene encoding Replication Factor C1 (*RFC1*) was real and present in both alleles in the affected individuals, consistent with the recessive inheritance pattern hypothesized for CANVAS (Figure S2). In addition, this was the only candidate that (I) was localized to the chr4 linkage region and (II) was able to be validated using existing STR

detection tools (see below). In both individuals with CANVAS, the novel (AAGGG)<sub>n</sub> pentamer replaced an (AAAAG)<sub>11</sub> motif located at the same position in the reference genome (chr4:39350045-39350095, hg19) and appeared to be significantly expanded compared to controls. Visualization of the region in the UCSC genome browser identified that the reference motif (AAAAG)<sub>n</sub> is the 3' end of an Alu element, AluSx3. In individuals with CANVAS, the (AAAAG)<sub>n</sub> motif is substituted by the (AAGGG)<sub>n</sub> motif, with potential interruptions to the Alu element (Figure 2D).

### **Confirmation of (AAGGG)<sub>n</sub> STR in off-target WES reads**

While WGS was only performed in two individuals with CANVAS, the majority of the cohort (n=27) was analyzed by WES. The putative pathogenic CANVAS RE is located in intron 2 of *RFC1*, 2863bp downstream of exon 2 and 2952bp upstream of exon 3. Therefore WES data is *a priori* assumed to be uninformative for this RE as it is not targeted during DNA capture. However, given that WES can often have off-target reads, we hypothesized that off-target reads might be captured at the (AAGGG)<sub>n</sub> locus. Visual assessment of the WES data in IGV identified 14 of 27 individuals with at least one read mapping within the *RFC1* STR locus. The maximum off-target read coverage at this locus was two, with a median of one read. While three individuals only had reads that correspond to the reference genome sequence (AAAAG)<sub>n</sub>, eleven affected individuals from nine families had reads containing (AAGGG)<sub>n</sub> repeats (Table 1). Furthermore, single affected individuals from families CANVAS2 and CANVAS8 showed evidence of both (AAGGG)<sub>n</sub> and (AAAGG)<sub>n</sub> motifs at the *RFC1* STR locus. This observation raised the possibility that CANVAS might result from pathogenic expansions of different pentanucleotide motifs.

### **Validation with existing STR detection tools**

Multiple tools have been developed in recent years that test for the presence of REs at pre-defined STRs. Therefore, we inserted the novel *RFCI* STR into the STR reference files and used exSTRa, EH, TREDPARSE, STRetch and GangSTR to estimate the size of the STR, and/or detect REs in the WGS data from the two original CANVAS samples (CANVAS1 and 9) and seven additional individuals with CANVAS. The seven additional CANVAS samples selected for WGS were those with WES evidence for an alternate (AAAGG)<sub>n</sub> motif (families CANVAS2 and 8), and those who did not appear to have a RE at the (AAGGG) locus (families CANVAS11, 13, 17 and 19). The analysis cohorts were divided into PCR-based (CANVAS1, 9 and 31 controls) and PCR-free (CANVAS2, 8, 11, 13, 17 and 19) WGS samples since PCR protocols have previously been shown to affect RE detection.<sup>8</sup> Using exSTRa, we confirmed the homozygous inheritance of the (AAGGG)<sub>n</sub> motif in three individuals (CANVAS1, 2 and 9, Figure 3). The ECDF pattern for CANVAS2 is consistent with the presence of one shorter and one longer (AAGGG)<sub>n</sub> RE, while CANVAS8 appears to have only a single (AAGGG)<sub>n</sub> allele. Screening for the (AAAGG)<sub>n</sub> motif at the chr4 *RFCI* locus identified an expansion for this STR only in CANVAS8. Visualization in IGV confirmed the presence of the (AAAGG)<sub>n</sub> motif embedded within the reference STR: (AAAAG)<sub>n</sub>-(AAAGG)<sub>n</sub>-(AAAAG)<sub>n</sub> (Figure S3). This observation raised the possibility that an expanded (AAAGG)<sub>n</sub> motif might also be associated with CANVAS. This analysis also confirmed that families CANVAS11, 13, 17 and 19 do not share the (AAGGG)<sub>n</sub> RE at this locus.

EH, TREDPARSE and GangSTR were used to estimate the CANVAS RE allele sizes (Figure 3), which were highly variable depending on the tool used. EH reported larger allele sizes for individuals with CANVAS (range of [68,30]) compared to GangSTR (range of [27,2]) and TREDPARSE (range of [14,7]). Furthermore, all three tools inferred the presence



of two alleles, even in individuals who carry a single allele, and hence do not appear to be distinguishing read contributions between the alleles, also contributing to unreliable size estimates. Reads comprised of the (AAGGG)<sub>n</sub> motif in particular also showed evidence of high read sequencing error. Based on these results, we can infer that while the CANVAS samples were all correctly identified as having homozygous RE at *RFC1*, estimates of expansion size are inconsistent and appear likely to significantly underestimate the actual repeat size.

The consensus of the different tools was that the *RFC1* (AAGGG)<sub>n</sub> STR was present in three of the control WGS datasets [two heterozygous, one homozygous, allele frequency ~0.06 (4/62), Figure 3]. No individuals were identified to carry the (AAAGG)<sub>n</sub> motif. As with the CANVAS samples, the STR sizing estimates using the different tools was inconsistent, therefore no conclusions could be drawn from this *in silico* analysis regarding the relative size of the (AAGGG)<sub>n</sub> STR in controls compared to individuals with CANVAS. We then analyzed a larger in-house collection of unrelated control Coriell WGS samples (N=69) and again failed to identify the (AAAGG)<sub>n</sub> motif. However, we identified six individuals heterozygous for the (AAGGG)<sub>n</sub> STR, representing a frequency estimate of ~0.04 [(6/138) (Figure S4)]. Using the NGS QC software tool peddy, we found evidence that two of these heterozygous individuals are of European ancestry and that two further individuals are of admixed Native American ancestry. Finally, we accessed WGS from GTEx for 133 individuals who have matching brain (cerebellum) RNA-seq. Our analysis identified 11 heterozygous carriers of the (AAGGG)<sub>n</sub> STR, representing an estimated allele frequency of ~0.04 (11/266), consistent with our in-house collection.

## **Validation of the (AAGGG)<sub>n</sub> RE as the causal variant for CANVAS**

We developed a PCR assay to rapidly screen for the presence of a non-expanded allele at the *RFC1* STR. Although the screen does not distinguish between the (AAAAG)<sub>n</sub> reference or (AAGGG)<sub>n</sub> novel motif, amplification of an ~250bp fragment indicates one or more alleles is not expanded. Conversely, the complete absence of the PCR product provides indirect evidence of a homozygous RE. Analysis of all available DNA samples from individuals with CANVAS suggested that the reference STR at the *RFC1* locus was not present in any of the 30 clinically diagnosed individuals from 18 (of 22) CANVAS families (Table 1, Figure 4). Notably, all nine individuals from these families without a diagnosis of CANVAS carried at least one non-expanded *RFC1* allele (data not shown). To directly confirm expansion of the novel (AAGGG)<sub>n</sub> motif in *RFC1*, we developed a locus specific repeat-primed PCR assay. Consistent with the PCR assay, all affected individuals from the 18 families demonstrated a saw-toothed ‘ladder’ when the repeat-primed PCR products were analyzed by capillary array (Table 1, Figure 4). These results suggest homozygous loss of the reference sequence and insertion of a novel (AAGGG)<sub>n</sub> RE underlies CANVAS in these 18 families, although the size of the RE remains undetermined. Molecular analysis of the DNA for the three in-house control individuals with the *in silico* predicted (AAGGG)<sub>n</sub> motif (Figure 3) demonstrated a ~250bp product in both heterozygous samples but no product in the homozygous individual. The repeat-primed assay demonstrated a saw-toothed ladder in all three samples (data not shown). Collectively, these analyses suggested all three control individuals have at least one copy of the novel (AAGGG)<sub>n</sub> motif at *RFC1*, although the size of the RE cannot be determined.

In four families (CANVAS11, 13, 17 and 19) the lack of a repeat-primed PCR product and presence of the expected reference PCR amplicon suggested the pathogenic expanded (AAGGG)<sub>n</sub> RE was not present on either allele. This implied that these individuals

have a different CANVAS-causing mutation in *RFC1*, or there is locus heterogeneity. A third possibility is that they do not have CANVAS but instead a related ataxia. Therefore, we performed WGS on these individuals and initially screened for known REs associated with ataxias. A CAG trinucleotide expansion in *ATXN3*, associated with spino-cerebellar ataxia type 3 (SCA3; OMIM 109150 also known as Machado-Josephs disease) was identified in CANVAS13 (Figure S5) and confirmed by diagnostic testing. The remaining three cases were screened genome-wide with EHdn for potentially pathogenic RE of (AAGGG)<sub>n</sub> and other STR motifs, however no further candidate RE were identified. The WGS was then screened for novel or rare SNPs and indels in genes known to cause ataxia. No *de novo* or rare variants were identified in *RFC1* however a potential genomic re-diagnosis was achieved in two additional families. In CANVAS17 two variants [NM\_001278055:c.12398delT, p.(Phe4133Serfs\*28) and NM\_001278055:c.5306T>A, p.(Val1769Asp)] were identified in the gene encoding saccin (*SACS*) and segregation analysis confirmed they were in trans. Biallelic mutations in *SACS* cause spastic ataxia of the Charlevoix-Saguenay type (MIM: 270550). In CANVAS19, a heterozygous variant in the gene encoding FAT tumor suppressor homolog 2 [*FAT2*, NM\_001447.2:c.4370T>C, p.(Val1457Ala)] was identified. Heterozygous mutations in *FAT2* have recently been associated with SCA45 (MIM: 604269).<sup>44</sup> No potentially pathogenic variants were identified in CANVAS11, however a variant of unknown significance was identified in the gene encoding Ataxin 7 [*ATXN7*, NM\_001177387.1:c.2827C>G, p.(Arg943Gly)].

### **A single founder event for the (AAGGG)<sub>n</sub> RE in *RFC1***

We performed haplotype analysis to determine if the (AAGGG)<sub>n</sub> RE arose more than once in human history. Analysis of haplotypes inferred from the WES data identified a core ancestral haplotype, comprised of 27 SNPs (Figure 5A), that was shared by most individuals

except CANVAS14 (Table 1, Table S4). The core haplotype spans four genes (*TMEM156*, *KLHL5*, *WDR19* and *RFC1*) and is 0.36 MB in size (chr4:38995374-39353137 (hg19)). Inspection of this region in the UCSC browser suggested that the core haplotype overlaps with a region of strong linkage disequilibrium in European and Asian populations (Han Chinese and Japanese from Tokyo), but not the Yoruba population (an ethnic group from West Africa, Figure 5B). Using a DNA recombination and haplotype-based mutation dating technique<sup>42</sup>, we estimate that the most recent common ancestor (MRCA) of the CANVAS cohort lived approximately 25,880 (CI: 14,080-48,020) years ago (Figure 5C). This age estimate corresponds to the size of the haplotype and LD block and is roughly equivalent to the origin of modern Caucasians as represented by the HAPMAP Caucasian cohort. Further investigation of the haplotypes allowed us to infer a simple phylogeny based on identified clusters of shared haplotypes extending beyond the core haplotype, suggesting that some individuals have common ancestors more recent than that of the MRCA for the whole group. This approach identified four subgroups. Group A had a MRCA dating back 5,600 (CI: 2120-15520) years and group B (further divided into groups B1 and B2) have a MRCA dating back 4,180 years (CI: 2240-7940). Furthermore, one individual shared part of their haplotype with both groups A and B, suggesting that group B is a distant branch of the MRCA of group A. Another subgroup, C, has a MRCA that lived 1860 (CI: 560-7020) years ago. The final group labelled N, do not have any additional sharing beyond the core haplotype.

Next, we compared the haplotype of all nine in-house control samples that carry the (AAGGG)<sub>n</sub> STR (3 in-house controls and 6 from the Coriell collection) to the core haplotype defined in the individuals with CANVAS. All controls shared at least part of the core haplotype, again suggesting that the (AAGGG)<sub>n</sub> STR arose once in history. Finally, we determined that nine of the 11 individuals from GTEx heterozygous for the (AAGGG)<sub>n</sub> STR

also shared the same core haplotype identified in the individuals with CANVAS. Haplotype-specific SNPs enabled us to analyse the expression of the (AAGGG)<sub>n</sub> *RFC1* allele in the cerebellum RNA-seq data and confirm that the STR did not inhibit the expression of *RFC1* compared to the reference (AAAAG)<sub>n</sub> allele. The remaining two carriers do not appear to share the core haplotype. As they do not have heterozygous SNPs in their exons, exon specific expression could not be determined.

## DISCUSSION

Since the first description of the syndrome of cerebellar ataxia with bilateral vestibulopathy in 2004<sup>45</sup> and proposal of CANVAS as a distinct clinical entity in 2011<sup>11</sup> there has been little progress made in delineating the etiology of the disorder. While most affected individuals are described as idiopathic, reports of multiple affected sib pairs<sup>12</sup> and a family with three affected individuals<sup>46</sup> have suggested that an autosomal recessive mode of inheritance is most likely. The genetic basis of CANVAS has now been identified and validated in two independent studies, one recently published by Cortese et al<sup>32</sup> and this study. Both studies utilized a similar study design, with linkage analysis to reduce the genomic search space to a modest interval (<2Mb), but no plausible causal variant(s) could be identified in WES data. WGS was then performed on multiple individuals and Cortese et al successfully identified the RE by visually inspection of the aligned read pairs inside the linkage region using the Integrative Genomics Viewer. In contrast, we utilized a bioinformatics approach and performed genome-wide analysis of WGS data to identify potential RE and then prioritized further investigation of the only RE located within the linkage interval. While both approaches were successful, the bioinformatics approach to RE detection, as described in this study, is likely more sensitive and practical, and can be applied even in the absence of a small, or indeed any, linkage region. Furthermore, using a bioinformatics approach allows simultaneous testing of other potentially causal RE due to differential diagnosis. For example, we quickly re-diagnosed an affected individual with a pathogenic SCA3 RE.

Importantly, Cortese et al extend the clinical significance of the CANVAS RE by demonstrating it is potentially a common cause of unsolved ataxia not meeting the diagnostic criteria of CANVAS. Screening for the homozygous inheritance of the AAGGG motif in a

cohort of 150 individuals with sporadic late-onset ataxia diagnosed 33 (22%). This is consistent with the relatively high allele frequency of the AAGGG STR (~0.03) we report in this paper and our estimate that the expansion first arose, presumably in Caucasians, ~25,000 years ago.

Previously, the only variant associated with CANVAS was a heterozygous missense variant in the gene encoding E74 Like ETS Transcription Factor 2 (*ELF2*), which segregated with the disorder in three individuals in a single family.<sup>46</sup> Here we confirm that CANVAS is caused by the homozygous inheritance of a novel and expanded intronic (AAGGG)<sub>n</sub> pentamer in *RFC1*. We found this motif in 30 of 31 individuals with a RE at this locus. In only a single individual did we observe a different, presumably pathogenic motif; CANVAS8 had one allele with an expanded (AAGGG)<sub>n</sub> motif, whereas the second allele appeared to consist of an expanded novel (AAAGG)<sub>n</sub> motif. Notably, this alternate motif does not share the AAGGG haplotype (Figure S3). Analysis of the core haplotype in the majority of individuals with CANVAS suggests that the (AAGGG)<sub>n</sub> STR arose once, approximately 25,000 years ago, most likely in Europe. While the majority of individuals in our cohort who carry the (AAGGG)<sub>n</sub> RE are of European ancestry, the RE is also present in non-European individuals, including a Lebanese family with CANVAS and two carriers of admixed Native American ancestry. Given the age of the CANVAS RE and recent human admixture it is likely that the locus may underlie CANVAS in apparently non-European individuals, despite the disorder being highly overrepresented in Caucasian populations.

### **Mechanism of pathogenicity**

There are multiple mechanisms by which RE can lead to pathogenicity, including RNA toxicity, protein toxicity and loss or gain of function.<sup>7</sup> It is not yet known how the

(AAGGG)<sub>n</sub> RE in *RFC1* causes CANVAS, however the homozygous inheritance pattern suggests a loss-of-function mechanism, rather than RNA or protein toxicity. Interestingly, Cortese et al were unable to determine a mechanism of action. The AAGGG RE did not appear to alter *RFC1* gene expression nor protein levels, and no AAGGG RNA foci deposits were observed in cells from individuals with CANVAS. While the gene has not been previously associated with any disorder it appears extremely intolerant to LoF (pLI = 0.97; observed/expected = 0.18, CI 0.12-0.31).<sup>47</sup> In addition, siblings in the families we studied carried the pathogenic RE in a heterozygous state but did not manifest any signs of the disorder. This observation is analogous to Freidreich ataxia, a recessive genetic ataxia caused by loss of function (LoF) of *FRDA* due to a pathogenic intronic RE.

*RFC1* encodes a subunit of replication factor C, a five-subunit protein complex required for DNA replication and repair. Analysis of the Genotype-Tissue Expression (GTEx) database demonstrated significant expression of *RFC1* in brain tissue, particularly the cerebellum. Replication factor C catalyzes opening the protein ring of proliferating cell nuclear antigen (PCNA), allowing it to encircle the DNA and function as a scaffold to recruit proteins involved in DNA replication, repair and remodeling.<sup>48</sup> Mutations in multiple DNA replication and repair genes such as *TCD1*, *PNKP*, *XRCC1* and *APTX* result in ataxia<sup>49</sup>, highlighting the central role of this pathway in these overlapping disorders. One of the best known examples is the severe and early onset autosomal recessive disorder, ataxia telangiectasia, which is caused by mutations in the gene encoding ATM serine/threonine kinase (*ATM*), which is important for the repair of DNA double-strand breaks.<sup>50</sup>

The minimum pathogenic length and fine structure of the *RFC1* RE is currently unknown. The short-read NGS technologies utilized in this study were unable to extend more



than ~100bp into the repeat sequence and efforts to amplify across the region using long range PCR were unsuccessful. While the repeat-primed PCR assay indicates the presence of an expanded (AAGGG) motif, it does not extend beyond ~250bp (50 repeat units). The application of long read sequencing technologies will be required to elucidate both the length of the pathogenic allele and the repeat composition. Both of these parameters provide important clinical information regarding onset, progression and pathogenicity in other genetic ataxias such as SCA1 and Friedreich ataxia.<sup>51; 52</sup> In addition, they will be required to elucidate the nature of the (AAGGG)<sub>n</sub> *RFC1* STR in control individuals-for example do these represent a fully expanded pathogenic allele or potentially an allele where the (AAGGG)<sub>n</sub> motif has replaced the reference (AAAAG)<sub>11</sub> STR, but has not expanded to a pathogenic size. We show that the (AAGGG)<sub>n</sub> STR occurs within the 3' end of the Alu element, AluSx3. Alu elements typically have A-rich tails and in the reference sequence the *RFC1* Alu has an A-rich tail containing an (AAAAG)<sub>11</sub> STR. There is some evidence that motifs that follow the pattern A<sub>n</sub>G<sub>m</sub>, especially (AAAG)<sub>n</sub> and (AAAGG)<sub>n</sub>, display strong base-stacking interactions and are more likely to expand through replication slippage.<sup>31</sup> This suggests an inherent instability of A and G rich motifs, consistent with what we observe in CANVAS. Notably, a number of pathogenic RE located with Alu have previously been described, including SCA10, SCA31, SCA37 and Friedreich ataxia.<sup>22; 26; 53; 54</sup>

## Genomic re-diagnosis in CANVAS

Four of twenty-two families enrolled in this study with a clinical diagnosis of CANVAS did not harbor the RE or any other potentially pathogenic variants in the *RFC1* locus. CANVAS13 was re-diagnosed with SCA3 after the WGS data was analyzed using our computational pipeline for detecting known pathogenic STRs. In addition to cerebellar ataxia, individuals with SCA3 not uncommonly manifests a somatosensory impairment<sup>55; 56</sup> and

vestibular involvement may be variably present<sup>56</sup>, resulting in a phenotype indistinguishable from CANVAS.<sup>43</sup> This molecular re-diagnosis highlights the power of modern STR detection techniques to diagnose RE ataxias. In addition, NGS data provides the opportunity to simultaneously identify non-RE mediated causes of ataxia. In CANVAS17 we identified biallelic variants in *SACS* as the likely cause of disease. While individuals with spastic ataxia of the Charlevoix-Saguenay type may present with the combination of cerebellar ataxia and a peripheral neuropathy<sup>57; 58</sup> as seen in CANVAS, to our knowledge vestibular involvement has not previously been described, and so this potentially constitutes a novel manifestation of the disease. In addition, a very plausible heterozygous variant was identified in *FAT2* in CANVAS19. While classified as a VUS using ACMG guidelines,<sup>59</sup> the variant is only observed once in gnomAD and was predicted pathogenic by six *in silico* algorithms. Very recently, heterozygous point mutations affecting the last cadherin domain (p.Lys3586Asn) or the linker region (p.Arg3649Gln) of *FAT2* have been associated with SCA45, adding weight to classifying the variant, a p.Val1457Ala substitution in the thirteenth cadherin domain, as likely pathogenic. While the published clinical phenotype and mutation spectrum in SCA45 is limited, in common with CANVAS, it is a late onset and slowly progressive cerebellar ataxia.<sup>44</sup>

## Strengths and limitations of current STR detection tools

In this study, we implemented multiple computational tools to identify and validate the presence of a novel (AAGGG)<sub>n</sub> RE in individuals with CANVAS. In particular, the use of EHdn, with its non-reference based RE discovery framework, was crucial in identifying a putative candidate, with the reference-based STR detection tools facilitating the follow up analysis. Although all tools gave highly variable estimated repeat sizes, which are likely to be significantly less than the actual repeat size, they provided consistent evidence that the

CANVAS (AAGGG)<sub>n</sub> motif was expanded. This level of evidence is helpful before embarking on the potentially complex process of molecular validation. In our analysis, only a single tool (GangSTR) failed to detect the alternate (AAAGG)<sub>n</sub> RE expansion. It is not clear why this was the case, although it could be related to the more complicated [(AAAAG)<sub>n</sub>-(AAAGG)<sub>n</sub>-(AAAAG)<sub>n</sub>] repeat structure. This observation highlights the importance of utilizing multiple tools to provide redundancy in the data analysis pipeline. An additional issue we encountered, which potentially limited all tools, was the poor sequencing quality in reads containing the (AAGGG)<sub>n</sub> motif compared to other STRs.

In conclusion, in this study we show that CANVAS is caused by a recessively inherited, ancient RE located in intron 2 of *RFC1*. Recently developed RE discovery tools facilitated the identification and verification of this novel RE, in addition to identifying other genetic causes of disease in the cohort. Despite the RE being located in an intron, we demonstrate that previously generated WES data with low-coverage genome-wide off target reads were helpful in providing increased statistical confidence in RE identification. Therefore, reanalysis of previously generated WES datasets potentially offers a cost effective approach to facilitating identification of novel intronic RE in discovery projects. Finally, we anticipate that implementation of these tools into routine diagnostic pipelines has the potential to significantly increase the current diagnostic rates of 36% and 17%, recorded for clinical exome and targeted panel analyses of individuals with ataxia, respectively.<sup>60; 61</sup>

## **SUPPLEMENTAL DATA**

The supplemental data contain 5 figures and 4 tables.

## **CONFLICTS OF INTEREST**

The authors declare no conflicts of interest.

## **ACKNOWLEDGMENTS**

We would like to thank Egor Dolzhenko and Michael Eberle (Illumina) for access to the dataset EGA00001003562 from the European Genome-Phenome Archive. This work was supported by the Australian Government National Health and Medical Research Council (Program Grant 1054618 to MB), the NIH (NINDS grant R01NS082094 to BLF) and the Murdoch Children's Research Institute. MB was supported by an NHMRC Senior Research Fellowship (1102971) and DPB was supported by a Wellcome Clinical Research Career Development Fellowship. Additional funding was provided by the Independent Research Institute Infrastructure Support Scheme and the Victorian State Government Operational Infrastructure Program.

## **WEB RESOURCES**

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

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

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

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

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

Varsome: <https://varsome.com>

## ACCESSION NUMBERS

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

## LEGENDS

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

**Figure 2: Linkage of the CANVAS locus to chromosome 4 and identification of (AAGGG)<sub>n</sub> intronic insertion in *RFC1***

A. The pedigree of the family CANVAS9 highlights the apparent recessive inheritance pattern. B. Linkage analysis of CANVAS9 identified significant linkage to chromosome 4 (LOD=3.25). C. Linkage regions for individual families CANVAS1, 2, 3, 4 and 9 are shown in blue and the overlapping region shown in red (chr4:38887351-40463592). D. STR analysis of WGS from two unrelated individuals with CANVAS identified a novel expanded STR in the second intron of *RFC1*. The (AAAAG)<sub>n</sub> motif that is present in the reference genome and part of an existing Alu element (AluSx3) is replaced by an expanded (AAGGG)<sub>n</sub> STR.

**Figure 3: Computational validation of the (AAGGG)<sub>n</sub> STR**

The (AAGGG)<sub>n</sub> STR at the coordinates chr4:39350045-39350095 was added to the reference databases of the tools exSTRa, EH, GangSTR, TREDPARSE and STRetch and WGS data from four unrelated individuals with CANVAS was analysed [CANVAS1 (green), CANVAS2 (red), CANVAS8 (blue) and CANVAS9 (orange)]. The non-CANVAS controls

are presented in grey. Plots have been divided into PCR-based and PCR-free WGS (left and right columns, respectively).

#### **Figure 4: Genetic validation of the (AAGGG)<sub>n</sub> STR**

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

#### **Figure 5:**

A. Analysis of WES data identified an ancestral haplotype surrounding *RFC1* in all affected individuals confirmed to carry the (AAGGG)<sub>n</sub> RE. B. The core haplotype (blue highlight) was intersected with the linkage disequilibrium (LD) track in the UCSC browser (converted to hg18 coordinates). The three LD tracks represent the Yoruba population (top track), Europeans (middle) and Han Chinese and Japanese from Tokyo (bottom). Red areas indicate strong linkage disequilibrium. The core CANVAS haplotype spans a large LD block in Caucasians, which is broken up into two LD blocks in Japanese and Chinese, suggesting an ancient origin for the CANVAS repeat expansion allele. C. Haplotype sharing between individuals with CANVAS was used to determine the age of the most recent common ancestor (MRCA) of the cohort.

**Table 1: Clinical features and genetic analysis of *RFC1* STR 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	Caucasian
CANVAS2	2 (M)	✓	✓	✓	AAGGG and AAAGG	✗	✓	CANVAS	A	Caucasian
CANVAS3	2 (F)	✓	✓	✗	AAGGG	✗	✓	CANVAS	A	Caucasian
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	Caucasian-Maori
CANVAS8	1 (F)	✗	✓	✓	AAGGG and AAAGG	✗	✓	CANVAS	A/other	Caucasian
CANVAS9	4 (1M,3F)	✗	✓	✓	AAGGG	✗	✓	CANVAS	A	Lebanese
CANVAS10	1 (M)	✗	✓	✗	AAGGG	✗	✓	CANVAS	A	Caucasian
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*	Caucasian
CANVAS16	1 (F)	✗	✗	✗	NA	✗	✓	CANVAS	Not assessed	Caucasian
CANVAS17	2 (M)	✗	✓	✓	Reference	✓	✗	SACS	NA	Caucasian
CANVAS18	1 (F)	✗	✓	✗	ND	✗	✓	CANVAS	A	Caucasian-Maori

CANVAS19	1 (F)	✖	✖	✓	NA	✓	✖	SCA45	Not assessed	Caucasian
CANVAS20	2 (1M,1F)	✖	✖	✖	NA	✖	✓	CANVAS	Not assessed	Spanish
CANVAS21	1 (M)	✖	✖	✖	NA	✖	✓	CANVAS	Not assessed	Indian
CANVAS22	1 (M)	✖	✖	✖	NA	✖	✓	CANVAS	Not assessed	Hungarian
CANVAS23	1 (U)	✖	✖	✖	NA	✖	✓	CANVAS	Not assessed	Not reported

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

The gene reference sequences utilized were NC\_000004 and NM\_002913 (RFC1).



## REFERENCES

1. McMurray, C.T. (2010). Mechanisms of trinucleotide repeat instability during human development. *Nat Rev Genet* 11, 786-799.
2. Gymrek, M., Willems, T., Guilmatre, A., Zeng, H., Markus, B., Georgiev, S., Daly, M.J., Price, A.L., Pritchard, J.K., Sharp, A.J., et al. (2016). Abundant contribution of short tandem repeats to gene expression variation in humans. *Nature genetics* 48, 22-29.
3. Quilez, J., Guilmatre, A., Garg, P., Highnam, G., Gymrek, M., Erlich, Y., Joshi, R.S., Mittelman, D., and Sharp, A.J. (2016). Polymorphic tandem repeats within gene promoters act as modifiers of gene expression and DNA methylation in humans. *Nucleic Acids Res* 44, 3750-3762.
4. Benson, G. (1999). Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res* 27, 573-580.
5. Subramanian, S., Madgula, V.M., George, R., Mishra, R.K., Pandit, M.W., Kumar, C.S., and Singh, L. (2003). Triplet repeats in human genome: distribution and their association with genes and other genomic regions. *Bioinformatics* 19, 549-552.
6. La Spada, A.R., and Taylor, J.P. (2010). Repeat expansion disease: progress and puzzles in disease pathogenesis. *Nat Rev Genet* 11, 247-258.
7. Hannan, A.J. (2018). Tandem repeats mediating genetic plasticity in health and disease. *Nat Rev Genet* 19, 286-298.
8. Tankard, R.M., Bennett, M.F., Degorski, P., Delatycki, M.B., Lockhart, P.J., and Bahlo, M. (2018). Detecting Expansions of Tandem Repeats in Cohorts Sequenced with Short-Read Sequencing Data. *American journal of human genetics* 103, 858-873.
9. Ruano, L., Melo, C., Silva, M.C., and Coutinho, P. (2014). The global epidemiology of hereditary ataxia and spastic paraplegia: a systematic review of prevalence studies. *Neuroepidemiology* 42, 174-183.
10. Bird, T.D. (2018 Update). Hereditary Ataxia Overview. In *GeneReviews*((R)), M.P. Adam, H.H. Ardinger, R.A. Pagon, S.E. Wallace, L.J.H. Bean, K. Stephens, and A. Amemiya, eds. (Seattle (WA)).
11. Szmulewicz, D.J., Waterston, J.A., Halmagyi, G.M., Mossman, S., Chancellor, A.M., McLean, C.A., and Storey, E. (2011). Sensory neuropathy as part of the cerebellar ataxia neuropathy vestibular areflexia syndrome. *Neurology* 76, 1903-1910.

12. Szmulewicz, D.J., McLean, C.A., MacDougall, H.G., Roberts, L., Storey, E., and Halmagyi, G.M. (2014). CANVAS an update: clinical presentation, investigation and management. *J Vestib Res* 24, 465-474.
13. Harding, A.E. (1981). "Idiopathic" late onset cerebellar ataxia. A clinical and genetic study of 36 cases. *J Neurol Sci* 51, 259-271.
14. Szmulewicz, D.J. (2017). Combined Central and Peripheral Degenerative Vestibular Disorders: CANVAS, Idiopathic Cerebellar Ataxia with Bilateral Vestibulopathy (CABV) and Other Differential Diagnoses of the CABV Phenotype. *Curr Otorhinolaryngol Rep* 5, 167–174.
15. Cha, Y.H. (2012). Less common neuro-otologic disorders. *Continuum (Minneap Minn)* 18, 1142-1157.
16. Szmulewicz, D.J., Merchant, S.N., and Halmagyi, G.M. (2011). Cerebellar ataxia with neuropathy and bilateral vestibular areflexia syndrome: a histopathologic case report. *Otol Neurotol* 32, e63-65.
17. Szmulewicz, D.J., McLean, C.A., Rodriguez, M.L., Chancellor, A.M., Mossman, S., Lamont, D., Roberts, L., Storey, E., and Halmagyi, G.M. (2014). Dorsal root ganglionopathy is responsible for the sensory impairment in CANVAS. *Neurology* 82, 1410-1415.
18. Szmulewicz, D.J., Seiderer, L., Halmagyi, G.M., Storey, E., and Roberts, L. (2015). Neurophysiological evidence for generalized sensory neuronopathy in cerebellar ataxia with neuropathy and bilateral vestibular areflexia syndrome. *Muscle Nerve* 51, 600-603.
19. Szmulewicz, D.J., Waterston, J.A., MacDougall, H.G., Mossman, S., Chancellor, A.M., McLean, C.A., Merchant, S., Patrikios, P., Halmagyi, G.M., and Storey, E. (2011). Cerebellar ataxia, neuropathy, vestibular areflexia syndrome (CANVAS): a review of the clinical features and video-oculographic diagnosis. *Ann N Y Acad Sci* 1233, 139-147.
20. Petersen, J.A., Wichmann, W.W., and Weber, K.P. (2013). The pivotal sign of CANVAS. *Neurology* 81, 1642-1643.
21. Szmulewicz, D., MacDougall, H., Storey, E., Curthoys, I., and Halmagyi, M. (2014). A Novel Quantitative Bedside Test of Balance Function: The Video Visually Enhanced Vestibulo-ocular Reflex (VVOR) *Neurology* 82, S19.002.
22. Campuzano, V., Montermini, L., Molto, M.D., Pianese, L., Cossee, M., Cavalcanti, F., Monros, E., Rodius, F., Duclos, F., Monticelli, A., et al. (1996). Friedreich's ataxia:

autosomal recessive disease caused by an intronic GAA triplet repeat expansion.

Science 271, 1423-1427.

23. Taki, M., Nakamura, T., Matsuura, H., Hasegawa, T., Sakaguchi, H., Morita, K., Ishii, R., Mizuta, I., Kasai, T., Mizuno, T., et al. (2018). Cerebellar ataxia with neuropathy and vestibular areflexia syndrome (CANVAS). *Auris Nasus Larynx* 45, 866-870.
24. Maruta, K., Aoki, M., and Sonoda, Y. (2019). [Cerebellar ataxia with neuropathy and vestibular areflexia syndrome (CANVAS): a case report]. *Rinsho Shinkeigaku* 59, 27-32.
25. Bahlo, M., Bennett, M.F., Degorski, P., Tankard, R.M., Delatycki, M.B., and Lockhart, P.J. (2018). Recent advances in the detection of repeat expansions with short-read next-generation sequencing. *F1000Res* 7.
26. Seixas, A.I., Loureiro, J.R., Costa, C., Ordonez-Ugalde, A., Marcelino, H., Oliveira, C.L., Loureiro, J.L., Dhingra, A., Brandao, E., Cruz, V.T., et al. (2017). A Pentanucleotide ATTTC Repeat Insertion in the Non-coding Region of DAB1, Mapping to SCA37, Causes Spinocerebellar Ataxia. *American journal of human genetics* 101, 87-103.
27. Ishiura, H., Doi, K., Mitsui, J., Yoshimura, J., Matsukawa, M.K., Fujiyama, A., Toyoshima, Y., Kakita, A., Takahashi, H., Suzuki, Y., et al. (2018). Expansions of intronic TTTCa and TTTTA repeats in benign adult familial myoclonic epilepsy. *Nature genetics* 50, 581-590.
28. Dolzhenko, E., van Vugt, J., Shaw, R.J., Bekritsky, M.A., van Blitterswijk, M., Narzisi, G., Ajay, S.S., Rajan, V., Lajoie, B.R., Johnson, N.H., et al. (2017). Detection of long repeat expansions from PCR-free whole-genome sequence data. *Genome research* 27, 1895-1903.
29. Tang, H., Kirkness, E.F., Lippert, C., Biggs, W.H., Fabani, M., Guzman, E., Ramakrishnan, S., Lavrenko, V., Kakaradov, B., Hou, C., et al. (2017). Profiling of Short-Tandem-Repeat Disease Alleles in 12,632 Human Whole Genomes. *American journal of human genetics* 101, 700-715.
30. Dashnow, H., Lek, M., Phipson, B., Halman, A., Sadedin, S., Lonsdale, A., Davis, M., Lamont, P., Clayton, J.S., Laing, N.G., et al. (2018). STRetch: detecting and discovering pathogenic short tandem repeat expansions. *Genome biology* 19, 121.
31. Mousavi, N., Shleizer-Burko, S., and Gymrek, M. (2018). Profiling the genome-wide landscape of tandem repeat expansions. *BioRxiv*, <https://doi.org/10.1101/361162>
32. Cortese, A., Simone, R., Sullivan, R., Vandrovicova, J., Tariq, H., Yan, Y.W., Humphrey, J., Jaunmuktane, Z., Sivakumar, P., Polke, J., et al. (2019). Biallelic expansion of an

- intronic repeat in RFC1 is a common cause of late-onset ataxia. *Nature genetics* 51, 649-658.
33. Smith, K.R., Bromhead, C.J., Hildebrand, M.S., Shearer, A.E., Lockhart, P.J., Najmabadi, H., Leventer, R.J., McGillivray, G., Amor, D.J., Smith, R.J., et al. (2011). Reducing the exome search space for mendelian diseases using genetic linkage analysis of exome genotypes. *Genome biology* 12, R85.
  34. Bahlo, M., and Bromhead, C.J. (2009). Generating linkage mapping files from Affymetrix SNP chip data. *Bioinformatics* 25, 1961-1962.
  35. Abecasis, G.R., Cherny, S.S., Cookson, W.O., and Cardon, L.R. (2002). Merlin--rapid analysis of dense genetic maps using sparse gene flow trees. *Nature genetics* 30, 97-101.
  36. Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26, 841-842.
  37. Pedersen, B.S., Layer, R.M., and Quinlan, A.R. (2016). Vcfanno: fast, flexible annotation of genetic variants. *Genome biology* 17, 118.
  38. Wang, K., Li, M., and Hakonarson, H. (2010). ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res* 38, e164.
  39. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29, 15-21.
  40. Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* 30, 923-930.
  41. Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., and Smyth, G.K. (2015). limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* 43, e47.
  42. Gandolfo, L.C., Bahlo, M., and Speed, T.P. (2014). Dating rare mutations from small samples with dense marker data. *Genetics* 197, 1315-1327.
  43. Szmulewicz, D.J., Roberts, L., McLean, C.A., MacDougall, H.G., Halmagyi, G.M., and Storey, E. (2016). Proposed diagnostic criteria for cerebellar ataxia with neuropathy and vestibular areflexia syndrome (CANVAS). *Neurol Clin Pract* 6, 61-68.
  44. Nibbeling, E.A.R., Duarri, A., Verschuuren-Bemelmans, C.C., Fokkens, M.R., Karjalainen, J.M., Smeets, C., de Boer-Bergsma, J.J., van der Vries, G., Dooijes, D.,

- Bampi, G.B., et al. (2017). Exome sequencing and network analysis identifies shared mechanisms underlying spinocerebellar ataxia. *Brain* 140, 2860-2878.
45. Rinne, T., Bronstein, A.M., Rudge, P., Gresty, M.A., and Luxon, L.M. (1998). Bilateral loss of vestibular function: clinical findings in 53 patients. *J Neurol* 245, 314-321.
46. Ahmad, H., Requena, T., Frejo, L., Cobo, M., Gallego-Martinez, A., Martin, F., Lopez-Escamez, J.A., and Bronstein, A.M. (2018). Clinical and Functional Characterization of a Missense ELF2 Variant in a CANVAS Family. *Front Genet* 9, 85.
47. Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T., O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B., et al. (2016). Analysis of protein-coding genetic variation in 60,706 humans. *Nature* 536, 285-291.
48. Zhang, G., Gibbs, E., Kelman, Z., O'Donnell, M., and Hurwitz, J. (1999). Studies on the interactions between human replication factor C and human proliferating cell nuclear antigen. *Proc Natl Acad Sci U S A* 96, 1869-1874.
49. Yoon, G., and Caldecott, K.W. (2018). Nonsyndromic cerebellar ataxias associated with disorders of DNA single-strand break repair. *Handbook of clinical neurology* 155, 105-115.
50. Savitsky, K., Bar-Shira, A., Gilad, S., Rotman, G., Ziv, Y., Vanagaite, L., Tagle, D.A., Smith, S., Uziel, T., Sfez, S., et al. (1995). A single ataxia telangiectasia gene with a product similar to PI-3 kinase. *Science* 268, 1749-1753.
51. Kraus-Perrotta, C., and Lagalwar, S. (2016). Expansion, mosaicism and interruption: mechanisms of the CAG repeat mutation in spinocerebellar ataxia type 1. *Cerebellum Ataxias* 3, 20.
52. Mateo, I., Llorca, J., Volpini, V., Corral, J., Berciano, J., and Combarros, O. (2003). GAA expansion size and age at onset of Friedreich's ataxia. *Neurology* 61, 274-275.
53. Bushara, K., Bower, M., Liu, J., McFarland, K.N., Landrian, I., Hutter, D., Teive, H.A., Rasmussen, A., Mulligan, C.J., and Ashizawa, T. (2013). Expansion of the Spinocerebellar ataxia type 10 (SCA10) repeat in a patient with Sioux Native American ancestry. *PloS one* 8, e81342.
54. Sato, N., Amino, T., Kobayashi, K., Asakawa, S., Ishiguro, T., Tsunemi, T., Takahashi, M., Matsuura, T., Flanigan, K.M., Iwasaki, S., et al. (2009). Spinocerebellar ataxia type 31 is associated with "inserted" penta-nucleotide repeats containing (TGGAA)<sub>n</sub>. *American journal of human genetics* 85, 544-557.

55. Jardim, L.B., Pereira, M.L., Silveira, I., Ferro, A., Sequeiros, J., and Giugliani, R. (2001). Neurologic findings in Machado-Joseph disease: relation with disease duration, subtypes, and (CAG)n. *Arch Neurol* 58, 899-904.
56. Gordon, C.R., Zivotofsky, A.Z., and Caspi, A. (2014). Impaired vestibulo-ocular reflex (VOR) in spinocerebellar ataxia type 3 (SCA3): bedside and search coil evaluation. *J Vestib Res* 24, 351-355.
57. Gagnon, C., Desrosiers, J., and Mathieu, J. (2004). Autosomal recessive spastic ataxia of Charlevoix-Saguenay: upper extremity aptitudes, functional independence and social participation. *Int J Rehabil Res* 27, 253-256.
58. Vill, K., Muller-Felber, W., Glaser, D., Kuhn, M., Teusch, V., Schreiber, H., Weis, J., Klepper, J., Schirmacher, A., Blaschek, A., et al. (2018). SACS variants are a relevant cause of autosomal recessive hereditary motor and sensory neuropathy. *Human genetics* 137, 911-919.
59. Richards, S., Aziz, N., Bale, S., Bick, D., Das, S., Gastier-Foster, J., Grody, W.W., Hegde, M., Lyon, E., Spector, E., et al. (2015). Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet Med* 17, 405-424.
60. Sullivan, R., Yau, W.Y., O'Connor, E., and Houlden, H. (2019). Spinocerebellar ataxia: an update. *J Neurol* 266, 533-544.
61. Galatolo, D., Tessa, A., Filla, A., and Santorelli, F.M. (2018). Clinical application of next generation sequencing in hereditary spinocerebellar ataxia: increasing the diagnostic yield and broadening the ataxia-spasticity spectrum. A retrospective analysis. *Neurogenetics* 19, 1-8.

# 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

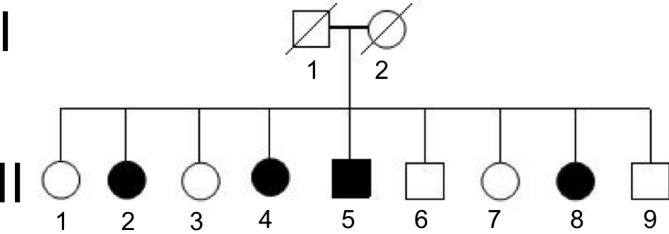
WGS: 5 patients negative for AAGGG in *RFC1*

- genomic re-diagnosis:
  - SCA3, SACS (compound heterozygous),
  - SCA45 (point mutation in *FAT2*)
  - VOUS: point mutation in *ATXN7*

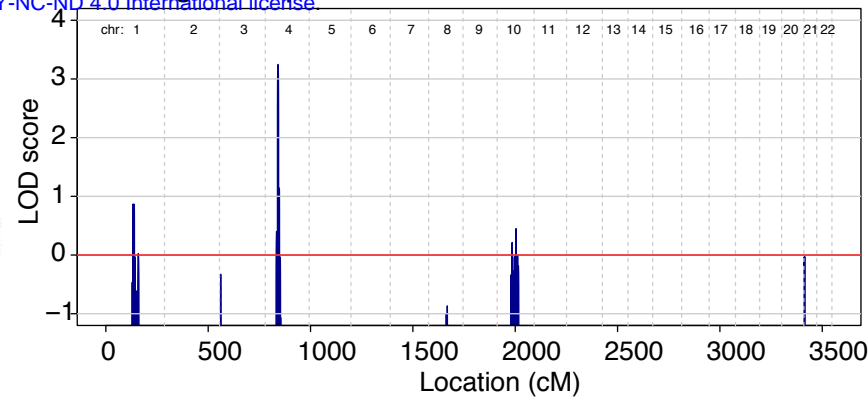
WGS: 2 patients with potential AAAGG/AAGGG RE

- CANVAS2-2: confirm AAGGG RE on both alleles
- **CANVAS8-8: confirm AAGGG on one allele, and AAAGG on second allele**

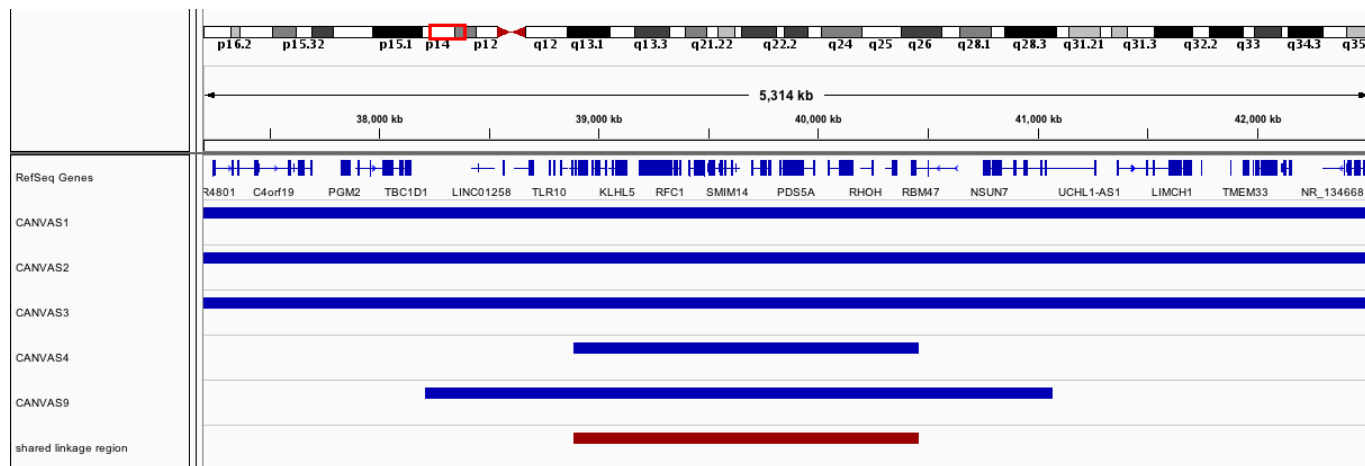
## A CANVAS9



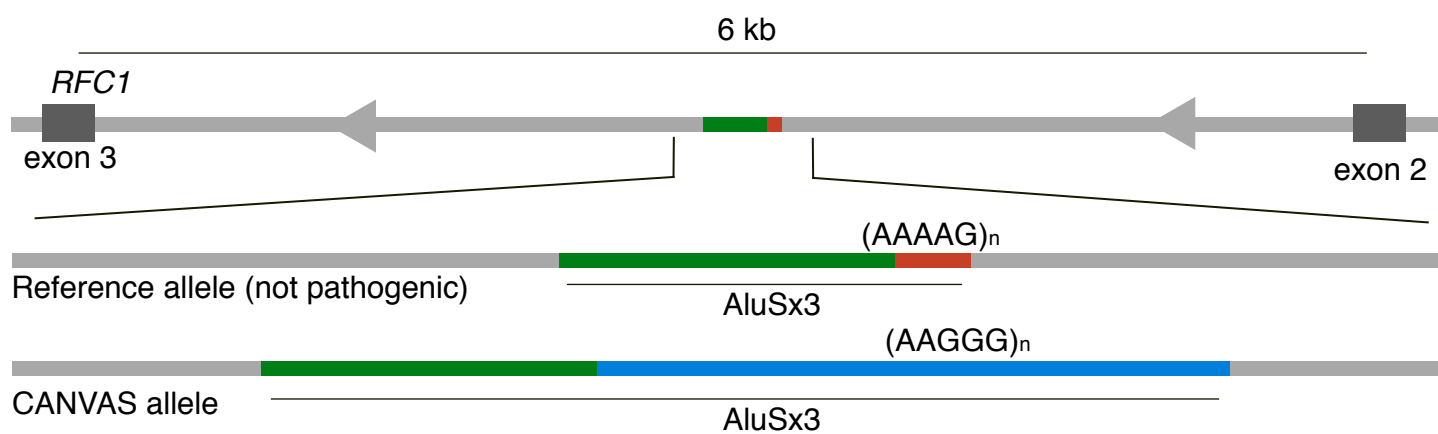
## B genome parametric LOD score for CANVAS9



## C



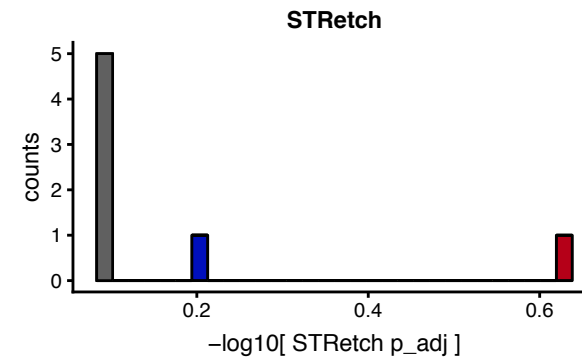
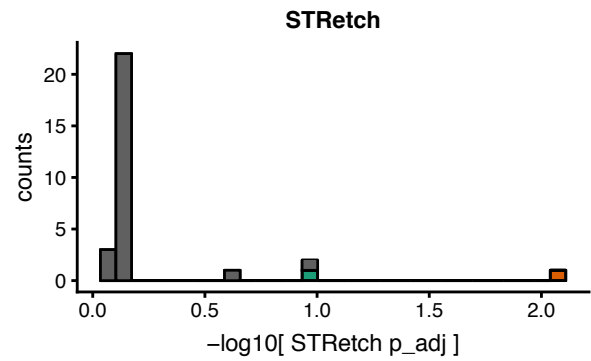
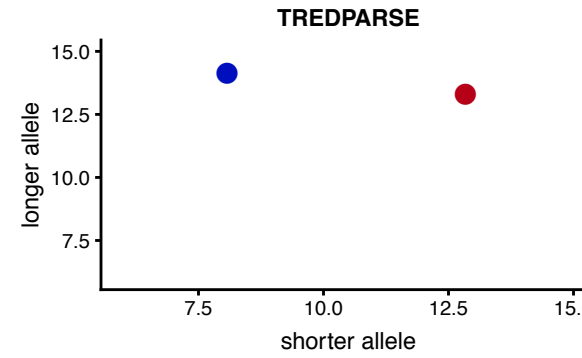
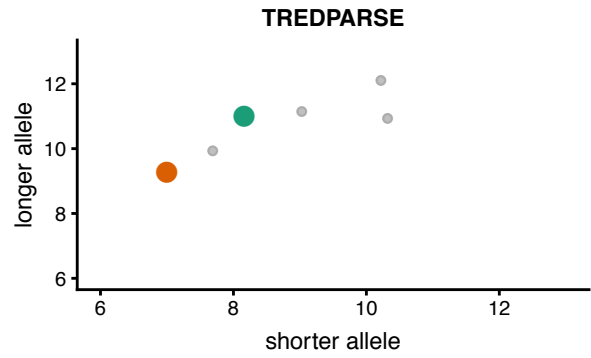
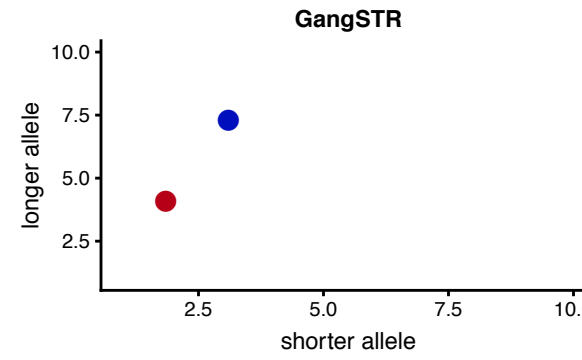
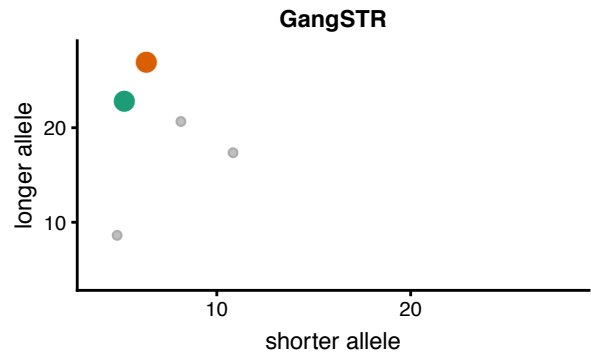
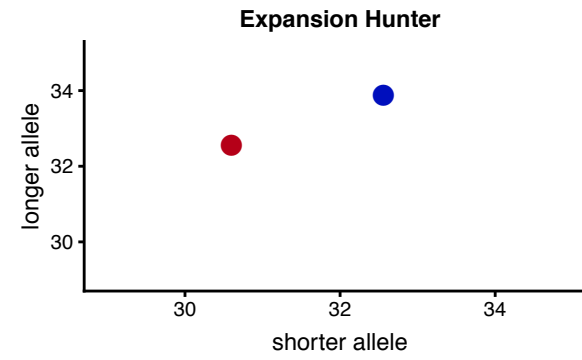
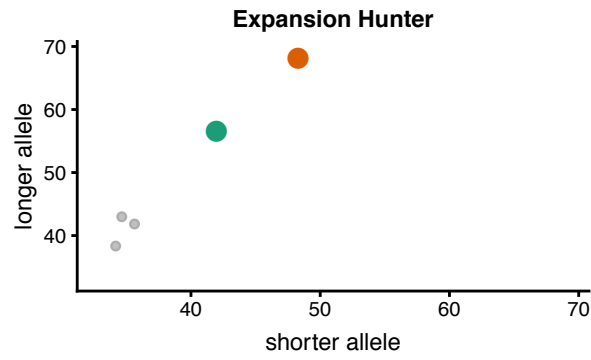
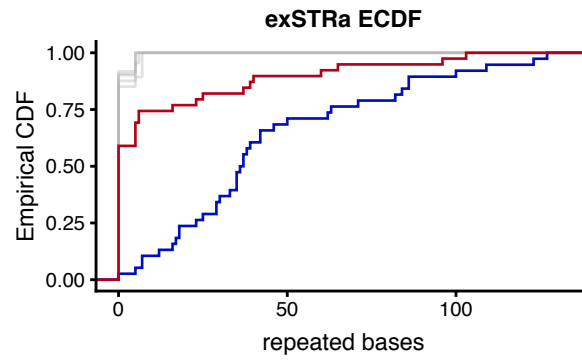
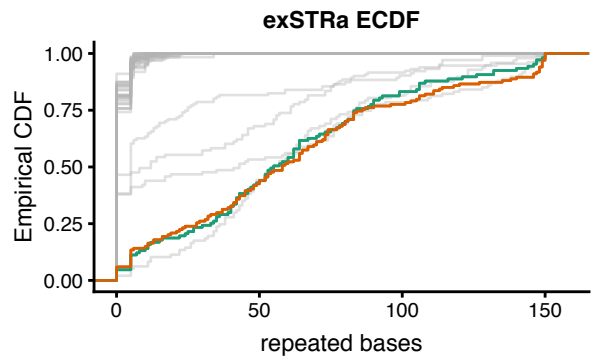
## D





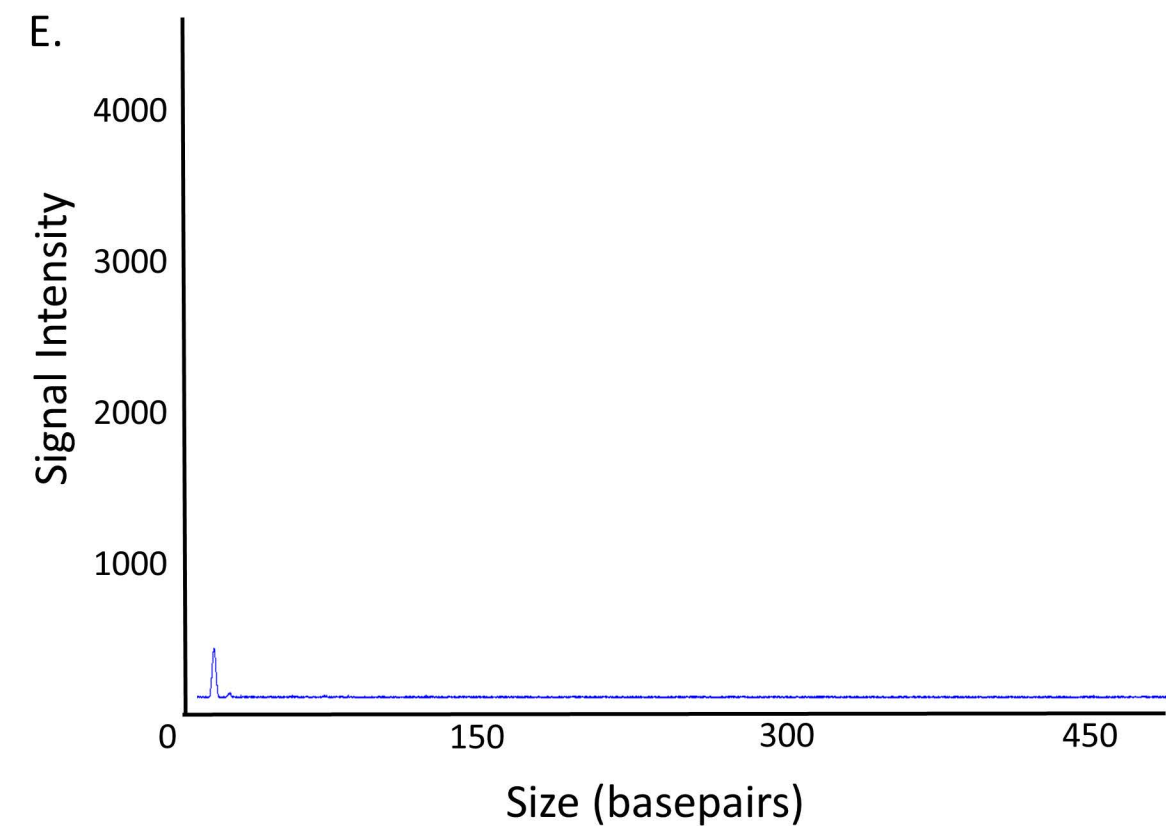
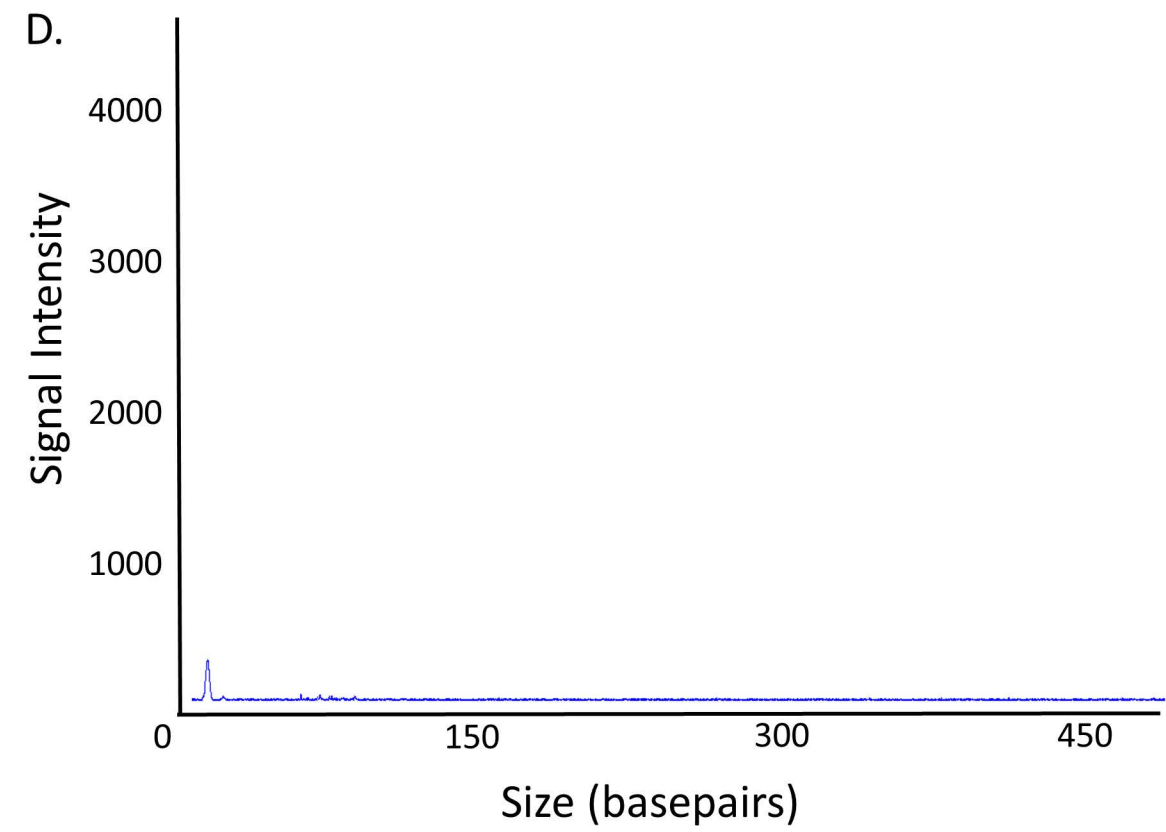
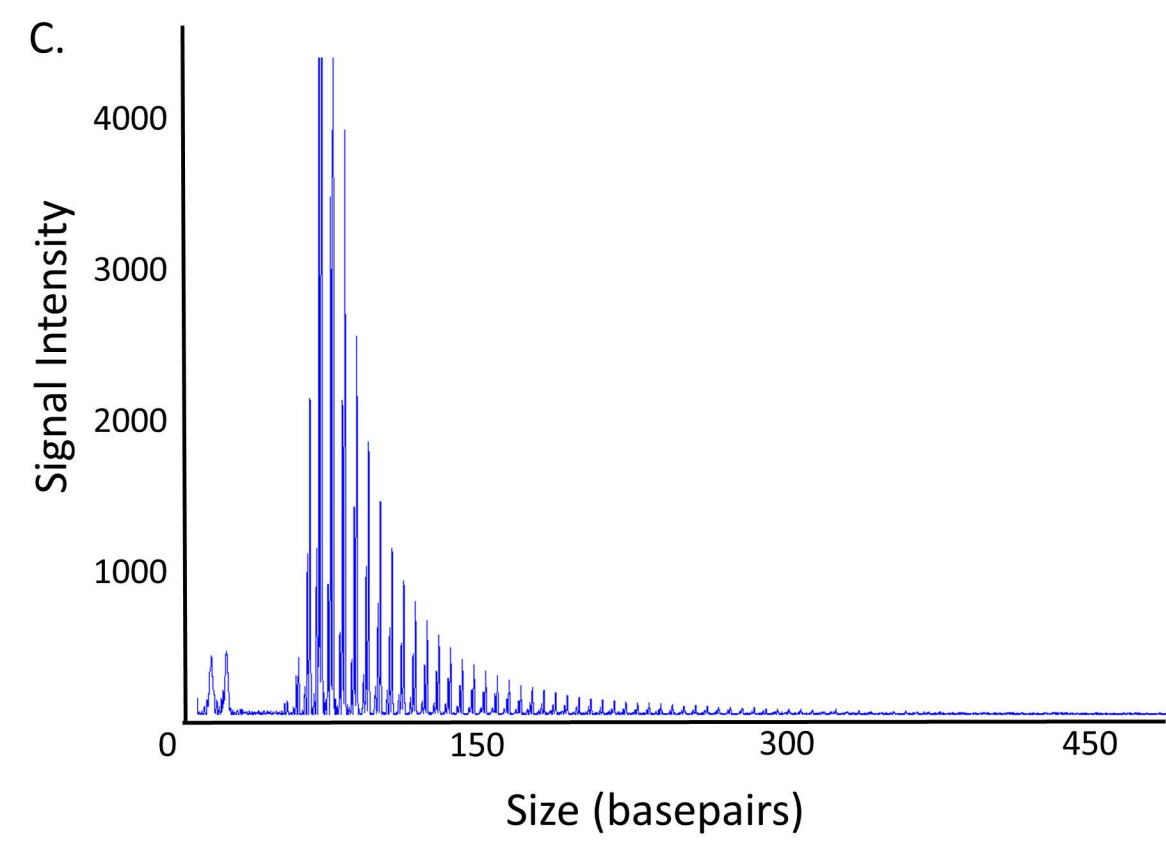
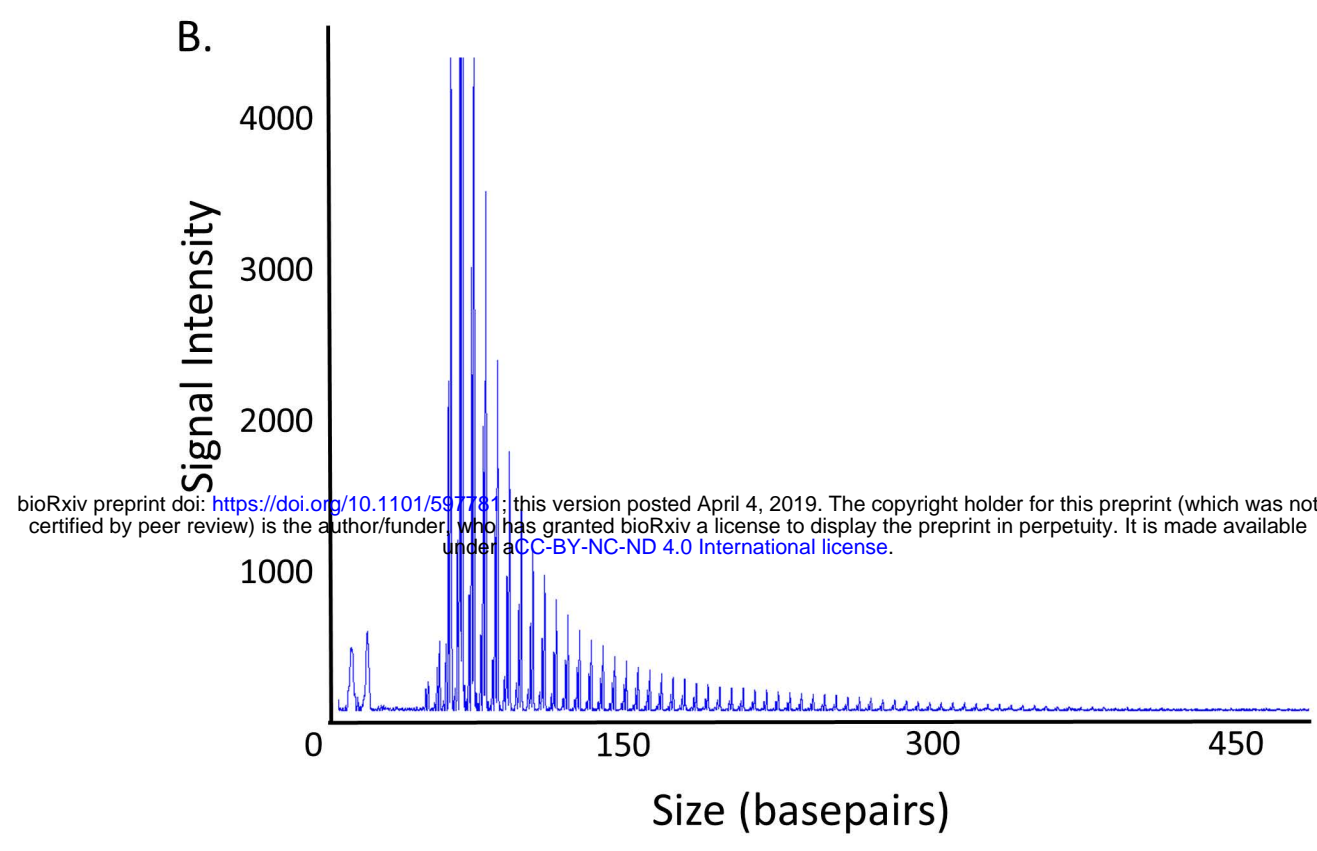
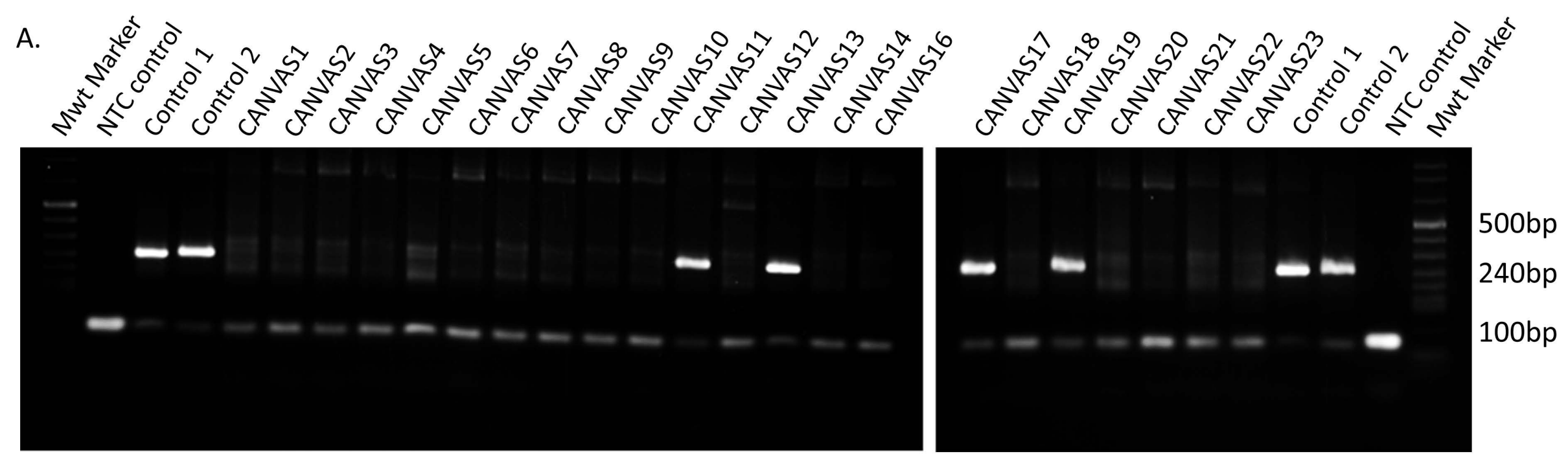
## PCR-based WGS

## PCR-free WGS

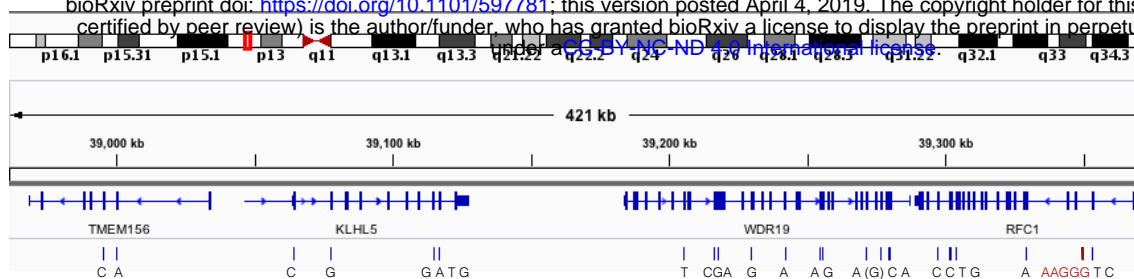


● CANVAS9  
● CANVAS1  
● non-CANVAS

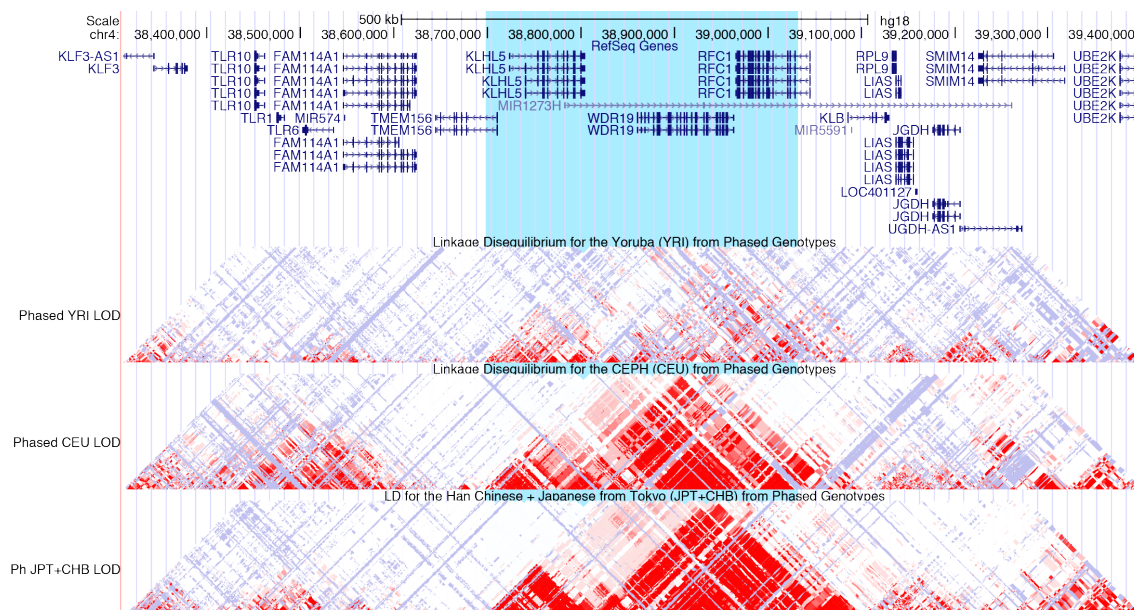
● CANVAS2  
● CANVAS8  
● non-CANVAS



A



B



C

