

# Normal and pathogenic variation of *RFC1* repeat expansions: implications for clinical diagnosis

Natalia Dominik,<sup>1,†</sup> Stefania Magri,<sup>2,†</sup> Riccardo Currò,<sup>1,3</sup> Elena Abati,<sup>1,4</sup> Stefano Facchini,<sup>1,5</sup> Marinella Corbetta,<sup>2</sup> Hannah MacPherson,<sup>1</sup> Daniela Di Bella,<sup>2</sup> Elisa Sarto,<sup>2</sup> Igor Stevanovski,<sup>6,7</sup> Sanjog R. Chintalaphani,<sup>7</sup> Fulya Akcimen,<sup>8</sup> Arianna Manini,<sup>1,4,9</sup> Elisa Vegezzi,<sup>5</sup> Iliara Quartesan,<sup>3</sup> Kylie-Ann Montgomery,<sup>1</sup> Valentina Pirota,<sup>10,11,12</sup> Emmanuele Crespan,<sup>10</sup> Cecilia Perini,<sup>10</sup> Glenda Paola Grupelli,<sup>10</sup> Pedro J. Tomaselli,<sup>13</sup> Wilson Marques,<sup>13</sup> Genomics England Research Consortium, Joseph Shaw,<sup>1</sup> James Polke,<sup>1</sup> Ettore Salsano,<sup>14</sup> Silvia Fenu,<sup>14</sup> Davide Pareyson,<sup>14</sup> Chiara Pisciotta,<sup>14</sup> George K. Tofaris,<sup>15</sup> Andrea H. Nemeth,<sup>15,16</sup> John Ealing,<sup>17</sup> Aleksandar Radunovic,<sup>18</sup> Seamus Kearney,<sup>19</sup> Kishore R. Kumar,<sup>20,21,22</sup> Steve Vucic,<sup>22,23</sup> Marina Kennerson,<sup>21,24,25</sup> Mary M. Reilly,<sup>1</sup> Henry Houlden,<sup>1</sup> Ira Deveson,<sup>6</sup> Arianna Tucci,<sup>1</sup> Franco Taroni<sup>2</sup> and Andrea Cortese<sup>1,3</sup>

**†These authors contributed equally to this work.**

## Abstract

Cerebellar Ataxia, Neuropathy and Vestibular Areflexia Syndrome (CANVAS) is an autosomal recessive neurodegenerative disease, usually caused by biallelic AAGGG repeat expansions in *RFC1*. In this study, we leveraged whole genome sequencing (WGS) data from nearly 10,000 individuals recruited within the Genomics England sequencing project to investigate the normal and pathogenic variation of the *RFC1* repeat.

We identified three novel repeat motifs, AGGGC (n=6 from 5 families), AAGGC (n=2 from 1 family), AGAGG (n=1), associated with CANVAS in the homozygous or compound heterozygous state with the common pathogenic AAGGG expansion. While AAAAG, AAAGGG and AAGAG expansions appear to be benign, here we show a pathogenic role for large AAAGG repeat configuration expansions (n=5). Long read sequencing was used to fully characterise the entire repeat sequence and revealed a pure AGGGC expansion in six patients, whereas the other patients presented complex motifs with AAGGG or AAAGG interruptions. All pathogenic motifs seem to have arisen from a common haplotype and are predicted to form highly stable G quadruplexes, which have been previously demonstrated to affect gene transcription in other conditions.

1 The assessment of these novel configurations is warranted in CANVAS patients with negative  
2 or inconclusive genetic testing. Particular attention should be paid to carriers of compound  
3 AAGGG/AAAGG expansions, since the AAAGG motif when very large (>500 repeats) or in  
4 the presence of AAGGG interruptions.

5 Accurate sizing and full sequencing of the satellite repeat with long read is recommended in  
6 clinically selected cases, in order to achieve an accurate molecular diagnosis and counsel  
7 patients and their families.

8

9 **Author affiliations:**

10 1 Department of Neuromuscular Diseases, University College London, WC1N 3BG London,  
11 UK

12 2 Unit of Medical Genetics and Neurogenetics, Fondazione IRCCS Istituto Neurologico Carlo  
13 Besta, 20133, Milan, Italy

14 3 Department of Brain and Behavioral Sciences, University of Pavia, 27100, Pavia, Italy

15 4 Department of Pathophysiology and Transplantation, University of Milan, 20122, Milan,  
16 Italy;

17 5 IRCCS Mondino Foundation, 27100, Pavia, Italy

18 6 Genomics Pillar, Garvan Institute of Medical Research, Sydney, 2010, Australia.

19 7 Centre for Population Genomics, Garvan Institute of Medical Research and Murdoch  
20 Children's Research Institute, 2010, Australia

21 8 Laboratory of Neurogenetics, National Institute on Aging, National Institutes of Health,  
22 Bethesda, MD, 2292, USA

23 9 Department of Neurology and Laboratory of Neuroscience, IRCCS Istituto Auxologico  
24 Italiano, 20145 Milan, Italy

25 10 Institute of Molecular Genetics IGM-CNR "Luigi Luca Cavalli-Sforza", 27100 Pavia, Italy

26 11 Department of Chemistry, University of Pavia, 27100 Pavia, Italy

27 12 G4-INTERACT, USERN, Pavia, Italy

28 13 Department of Neurology, School of Medicine of Ribeirão Preto, University of São Paulo,  
29 2650 Ribeirão Preto, Brazil

1 14 Clinic of Central and Peripheral Degenerative Neuropathies Unit, IRCCS Foundation, C.  
2 Besta Neurological Institute, 20126 Milan, Italy

3 15 Nuffield Department of Clinical Neurosciences, University of Oxford, OX3 9DU UK

4 16 Oxford Centre for Genomic Medicine, Oxford University Hospitals NHS Foundation Trust,  
5 OX3 7HE, UK

6 17 Salford Royal NHS Foundation Trust Greater Manchester Neuroscience Centre, Manchester  
7 Centre for Clinical Neurosciences Salford, Greater Manchester, M6 8HD, UK

8 18 Barts MND Centre, Royal London Hospital, Whitechapel, London, UK, E1 1BB

9 19 Department of Neurology, Royal Victoria Hospital, Belfast, BT12 6BA, UK

10 20 Kinghorn Centre for Clinical Genomics, Garvan Institute of Medical Research,  
11 Darlinghurst, NSW, 2010, Australia

12 21 Molecular Medicine Laboratory, Concord Hospital, Concord, NSW, 2139, Australia

13 22 Concord Clinical School, Faculty of Medicine and Health, University of Sydney, Sydney,  
14 NSW, Australia

15 23 Brain and Nerve Research Centre, Concord Hospital, Sydney, NSW, 2139, Australia

16 24 Northcott Neuroscience Laboratory, ANZAC Research Institute SLHD, Sydney, NSW,  
17 2050, Australia

18 25 School of Medical Sciences, Faculty of Medicine and Health, University of Sydney, Sydney,  
19 NSW, 2050, Australia

20

21 Correspondence to: Andrea Cortese, MD PhD

22 Department of Neuromuscular Diseases, UCL Queen Square Institute of Neurology, Queen  
23 Square, WC1N 3BG, London, UK

24 E-mail: andrea.cortese@ucl.ac.uk

25

26 **Running title:** Normal and pathogenic variation of *RFC1*

27 **Keywords:** *RFC1*; CANVAS; ataxia; neuropathy; repeat expansions; long-read sequencing

1 **Abbreviations:** CANVAS = Cerebellar Ataxia, Neuropathy and Vestibular Areflexia  
2 Syndrome; CCS = Circular Consensus maps; DIG = digoxigenin granules solution; DLS =  
3 Direct Label and Stain; EHDN = ExpansionHunterDeNovo; GEL = Genomics England; HMW  
4 = high molecular weight; NHS = National Health Service; nt = nucleotides; OGM = Optical  
5 genome mapping; PacBio = Pacific Biosciences; QGRS = Quadruplex forming G-Rich  
6 Sequences; RFC1 = Replication Factor Complex subunit 1; STR = short-tandem repeat; WGS  
7 =whole genome sequencing

8

## 9 **Introduction**

10 Cerebellar Ataxia, Neuropathy and Vestibular Areflexia Syndrome (CANVAS) is an  
11 autosomal recessive neurodegenerative disease characterized by adult onset and slowly  
12 progressive ataxia caused by the concurrent impairment of sensory neurons, the vestibular  
13 system and the cerebellum. In most cases, the disease is caused by biallelic AAGGG repeat  
14 expansions in the second intron of the Replication Factor Complex subunit 1 (*RF1*) gene.<sup>1-19</sup>  
15 Moreover, additional pathogenic (AAAGG)<sub>10-25</sub>(AAGGG)<sub>n</sub> and ACAGG configurations were  
16 identified in people from Oceania and East Asia suggesting the possibility of genetic  
17 heterogeneity, at the repeat locus, underlying this condition.<sup>20-23</sup>

18 In this study we leveraged whole genome sequencing (WGS) data from the 100,000 Genomes  
19 Project to investigate the normal and pathogenic variation of the *RF1* repeat as well as  
20 identifying additional pathogenic motifs causing CANVAS, that were further analysed by  
21 targeted long read sequencing.

22 We identified three novel pathogenic repeat configurations, AAGGC, AGGGC, and AGAGG,  
23 either in the homozygous or compound heterozygous state with AAGGG repeat, showing a  
24 similar or larger size compared to the common AAGGG expansion. In addition, pathogenic  
25 uninterrupted or interrupted AAAGG expansions were identified, with significantly larger  
26 sizes compared to the more frequent non-pathogenic AAAGG repeat.

27

## 28 **Materials and methods**

### 29 **Whole Genome Sequencing data analysis**

1 The 100,000 Genomes Project, run by Genomics England (GEL), was established to sequence  
 2 whole genomes of patients of the National Health Service (NHS) of the United Kingdom,  
 3 affected by rare diseases and cancer.<sup>24</sup> In this study, we leveraged GEL WGS data and screened  
 4 for the presence of pentanucleotide expansions in *RFC1* in 893 samples from patients  
 5 diagnosed with ataxia and 8107 controls, all aged 30 years or older. Repeat expansions were  
 6 detected using ExpansionHunterDeNovo (EHDN) v0.9.0. We considered all motifs composed  
 7 of 5 or 6 nucleotides occurring at the *RFC1* locus. Repeat motifs which were present, in the  
 8 homozygous or compound heterozygous state with the AAGGG expansion in ataxia cases, but  
 9 absent or significantly less frequent in controls, were considered to be possibly pathogenic and  
 10 were further assessed as described below.

11 Structural variants were detected using Manta25 as described in [https://re-](https://re-docs.genomicsengland.co.uk/genomic_data/)  
 12 [docs.genomicsengland.co.uk/genomic\\_data/](https://re-docs.genomicsengland.co.uk/genomic_data/).

13 Predicted genetic ancestries for samples in GEL are based on PCA analysis, using as reference  
 14 populations the five macro-ethnicities of 1000 Genomes project (European, African, South  
 15 Asian, East Asian, American). Samples where none of the components reaches 95% are  
 16 classified as Mixed.

17

## 18 **RP-PCR**

19 Samples identified to carry novel pathogenic repeat motifs by EHDN were tested by RP-PCR.  
 20 In addition, we screened a cohort of 540 samples with genetically confirmed *RFC1* CANVAS,  
 21 as defined by the presence of a positive RP-PCR for the AAGGG expansion and the absence  
 22 of an amplifiable PCR product from the flanking PCR, looking for expansions of different  
 23 repeat motifs on the second allele. RP-PCR for AAAAG, AAAGG and AAGGG expansions  
 24 were performed as previously described.<sup>1</sup> Moreover, the following primers were used:  
 25 AGGGC-Rv: 5'-CAGGAAACAGCTATGACCAACAGAGCAAGACTCTGT  
 26 TTCAAAAAGGGCAGGGCAGGGCAGGGCA-3'; AAGGC -Rv; 5'-AAGGC:  
 27 CAGGAAACAGCTATGACCAACAGAGCAAGACTCTGTTTCAAAA GGCAAGGCA  
 28 AGGCAA-3'; or AGAGG-Rv: 5'-  
 29 CAGGAAACAGCTATGACCAACAGAGCAAGACTCTGTTTCAAAAAGGAGAGGAG  
 30 AGGAGAGGAGA-3', depending on the configuration tested. The PCR conditions for  
 31 AGGGC and AAGGC were modified to 30s denaturation per cycle as opposed to 10s for all  
 32 the other configurations.

1

## 2 **Southern Blotting**

3 Briefly, 5µg of high molecular weight (HMW) DNA was enzymatically digested with EcoRI  
4 for 3 h and size fractionated on a 1.2% agarose gel for 15 hours. The gel was washed in  
5 depurination, denaturing and neutralising solutions for 45 minutes each, after which the blot  
6 was assembled to transfer DNA from the gel onto a positively charged membrane using an  
7 upward transfer method for 15 h. The DNA was UV crosslinked to the membrane and  
8 hybridised with a mixture of salmon sperm and *RFC1* probe in digoxigenin granules solution  
9 (DIG) (Roche) overnight. The membrane was then washed, blocked and anti-DIG antibody  
10 was added after which detection buffer and a chemiluminescent CDP-STAR substrate (Roche)  
11 was used to visualize hybridization fragments.

12

## 13 **Targeted RFC1 long read sequencing**

14 We performed long read sequencing to establish the precise repeat sequence in patients  
15 carrying novel likely pathogenic expansion of *RFC1*. Given the technical hurdle in sequencing  
16 large repeat expansions, samples were sequenced on different platforms, including Oxford  
17 Nanopore and Pacific Biosciences (PacBio) and target enrichment was performed with either  
18 clustered regularly interspaced short palindromic repeats CRISPR/CRISPR-associated protein-  
19 9 nuclease (Cas9) system or ReadUntil programmable selective sequencing.

20 Samples were extracted from blood using the Qiagen MagAttract HMW DNA Kit and quality  
21 was checked using readouts from a Thermo Scientific NanoDrop. For CRISPR/Cas9 targeted  
22 sequencing, fragment lengths were assessed using the Agilent Femto Pulse Genomic DNA  
23 165kb kit and only samples with the majority of fragments over 25kb were used., Libraries  
24 were prepared from 5µg of input DNA for each sample for both the Pacific Biosciences No-  
25 Amp Targeted Sequencing Utilizing the CRISPR-Cas9 System protocol (Version 09) and the  
26 Oxford Nanopore Ligation sequencing gDNA Cas9 enrichment (SQK-LSK109) protocol  
27 (Version: ENR\_9084\_v109\_revT\_04Dec2018). Libraries were sequenced on the Oxford  
28 Nanopore PromethION or MinION platforms or the Pacific Biosciences (PacBio) Sequel IIe,  
29 respectively. For the Oxford Nanopore Ligation sequencing gDNA Cas9 enrichment, we used  
30 four CRISPR-Cas9 guides from Nakamura et al. <sup>22</sup>, which were RFC1-F1: 5'-  
31 GACAGTAACTGTACCACAATGGG-3', RFC1-R1: 5'-

1 CTATATTCGTGGAAGGAACTATCTTGG-3', RFC1-F2: 5'-  
 2 ACACTCTTTGAAGGAATAACAGG-3' and RFC1-R2: 5'-TGAGGTATGAAT  
 3 CATCCTGAGGG-3', except for cases IV-1, XI-1, XII-1, for which only the two guides RFC1-  
 4 F2 and RFC1-R2 were used. The guides RFC1-F3: 5'-GAAACTAAATAGAACCAGCC-3'  
 5 RFC1-R3: 5'-GACTATGGCTTACCTGAGTG-3', which were designed in house, were used  
 6 for Pacific Biosciences No-Amp Targeted Sequencing and up to 10 samples were multiplexed  
 7 using PacBio barcoded adapters. Libraries loaded onto the PromethION and MinION were run  
 8 for 72 hours with standard loading protocols. Sequel IIe libraries were run for a movie time of  
 9 30 hours with an immobilisation time of 4 hours. All libraries were loaded neat.  
 10 Programmable targeted sequencing was performed as described previously.<sup>26</sup> HMW DNA was  
 11 sheared to ~20kb fragment size using Covaris G-tubes. Sequencing libraries were prepared  
 12 from ~3-5ug of HMW DNA, using native library prep kit SQK-LSK110, according to the  
 13 manufacturer's instructions. Each library was loaded onto a FLO-MIN106D (R9.4.1) flow cell  
 14 and run on an ONT MinION device with live target selection/rejection executed by the  
 15 *ReadFish* software package.<sup>27</sup> Detailed descriptions of software and hardware configurations  
 16 used for ReadFish experiments are provided in a recent publication that demonstrates the  
 17 suitability of this approach for profiling tandem repeats.<sup>26</sup> The target used in this study was the  
 18 *RFC1* gene locus  $\pm 50$  kb. Samples were run for a maximum duration of 72 h, with nuclease  
 19 flushes and library reloading performed at approximately 24 and 48 h timepoints for targeted  
 20 sequencing runs, to maximise sequencing yield.

21

## 22 **Bioinformatic analysis**

23 Alignment to the hg38 reference of Nanopore reads, PacBio CCS and PacBio subreads was  
 24 done using minimap2<sup>28</sup> with additional options “-r 10000 -g 20000 -E 4,0”. For PacBio  
 25 sequences, the recommended step of generating Circular Consensus maps (CCS) from subreads  
 26 was not always possible because of low depth of the sequencing data. The only CCS we could  
 27 obtain was for the AAGGG allele of Case V-1. After alignment, we used PacBio scripts  
 28 [<https://github.com/PacificBiosciences/apps-scripts>] to extract the repeat region  
 29 (extractRegions.py) and obtain waterfall plots (waterfall.py) for the following motifs:  
 30 AAGGG, AGAGG, AGGGC, AAGGC and AAAGG.

31 For programmable targeted sequencing, raw ONT sequencing data was converted to BLOW5  
 32 format using *slow5tools* (v0.3.0)<sup>29</sup> then base-called using *Guppy* (v6). Resulting FASTQ files

1 were aligned to the *hg38* reference genome using *minimap2* (v2.14-r883). The short-tandem  
2 repeat (STR) site within *RFC1* locus was genotyped using a process validated in our recent  
3 manuscript.<sup>27</sup> This method involves local haplotype-aware assembly of ONT reads spanning a  
4 given STR site and annotation of STR size, motif and other summary statistics using Tandem  
5 Repeats Finder (4.09), followed by manual inspection and motif counting.

## 7 **Haplotype analysis**

8 We used SHAPEITv4<sup>30</sup> with default parameters to phase a 2Mb region (chr4:38020000-  
9 40550000) encompassing the *RFC1* gene. To maximise available haplotype information, the  
10 entire Rare Diseases panel in Genomics England (78195 samples from patients affected by rare  
11 diseases) were jointly phased. Input data format was an aggregate VCF file with a total of  
12 551795 variants.

13 The estimation of the haplotype age is based on the online application *Genetic Mutation Age*  
14 *Estimator* (<https://shiny.wehi.edu.au/rafehi.h/mutation-dating/>)<sup>31</sup>. The method requires as  
15 input a list of ancestral segments for sampled individuals. We used the five individuals with  
16 pathogenic expansions (see Figure 3): AAGGG hom, ACAGG hom, case VII-1, case I-1, case  
17 III-3.

## 19 **Optical genome mapping (OGM)**

20 Patients for whom whole blood was available were subjected to BioNano optical genome  
21 mapping to gather additional information on the precise size of the expanded repeat. Ultra  
22 HMW genomic DNA was isolated as described by the Bionano Prep SP Frozen Human Blood  
23 DNA Isolation Protocol v2. Homogeneous ultra HMW DNA was labelled using the Bionano  
24 Prep Direct Label and Stain (DLS) Protocol provided with the kit and the homogeneous  
25 labelled DNA was loaded onto a Saphyr chip. Optical mapping was performed at theoretical  
26 coverage of 400x. Molecule files (.bnx) were aligned to hg38 with Bionano Solve script  
27 “align\_bnx\_to\_cmap.py” from Bionano Solve v3.6 ([https://bionano.com/software-](https://bionano.com/software-downloads/)  
28 [downloads/](https://bionano.com/software-downloads/)) using standard parameters. For each sample, molecules overlapping both markers  
29 flanking the repeat expansion were extracted (Marker IDs: 7723 and 7724). Intermarker  
30 distances were analysed by decomposing into two Gaussian components and using the gaussian



1 mean as allele size and repeat expansion size was calculated as the difference between the  
2 gaussian mean and the intermarker distance of a non-expanded allele (6858 bp).

3

#### 4 **G-quadruplexes analysis**

5 The propensity of the different repeat configurations in *RFC1* to form G-quadruplexes (G4)<sup>32</sup>  
6 was predicted using the Quadruplex forming G-Rich Sequences (QGRS) Mapper<sup>33</sup> and G4-  
7 Hunter software<sup>34</sup>, through which the likelihood to form a stable G4 is rated in terms of G-  
8 score values. Putative G4s were identified considering the following restriction parameters: for  
9 QGRS, maximum sequence length of 30 nucleotides (nt), minimum G-tetrads number in a G4  
10 of 2, loops lengths in the range of 0–36 nt, G-Score values > 15; For G4-Hunter threshold of  
11 1.5 with a window size of 20 nt.

#### 12 **Data availability**

13 Anonymized data are available from the corresponding author.

14

## 15 **Results**

### 16 **Identification of novel pathogenic repeat motifs in RFC1 in the 100,000** 17 **Genome project**

18 Out of 893 cases diagnosed with adult-onset ataxia (over the age of 30 years) recruited as part  
19 of the 100,000 Genome project, 124 cases had at least one AAGGG repeat expansion and 48  
20 had biallelic AAGGG repeat expansions, thus confirming a diagnosis of CANVAS and disease  
21 spectrum.

22 To identify additional likely pathogenic repeat motifs in *RFC1* we specifically looked for rare  
23 repeat configurations present in patients diagnosed with adult-onset ataxia (over the age of 30  
24 years) or in compound heterozygous state with the known pathogenic AAGGG repeat  
25 expansion, but absent or significantly less frequent in controls under the same conditions  
26 (**Table 1**).

27 We identified three cases carrying repeat expansions AAGGC (Case I-1), AGGGC (case II-  
28 1), or AGAGG (case VII-1) repeat motifs which were absent in non-neurological controls.  
29 AAGGC was present in the homozygous state, while AGGGC and AGAGG were in the

1 compound heterozygous state with the AAGGG expansion. One additional case of self-  
2 reported Asian ancestry carried the previously reported rare pathogenic ACAGG repeat  
3 expansion in the homozygous state.

4 AAAAG, AAAGGG, AAGAG expansions were found with similar frequency in patients and  
5 controls, supporting their non-pathogenic significance, while there was a higher percentage of  
6 compound heterozygous AAGGG/AAAGG carriers in ataxia cases ( $p=0.05$ ).

7 All predicted genetic ancestries for individuals carrying rare homozygous or compound  
8 heterozygous expansions in RFC1 are reported in **Supplementary Table 2**. Patients carrying  
9 AAGGC (Case I-1) and AGGGC (case II-1) expansion were of predicted South Asian and  
10 mixed ethnicity, respectively; ACAGG expansion carrier was confirmed to be East Asian based  
11 on the predicted genetic ancestry, while other repeat configurations were mostly identified in  
12 individuals of European or mixed ethnicity.

13 We did not identify any loss-of-function variant or structural variant in the RFC1 gene in  
14 individuals carrying heterozygous AAGGG repeat expansions.

15 The presence of AGGGC, AAGGC or AGAGG repeat expansions was confirmed by RP-PCR  
16 in all three cases, and the AAGGC repeat segregated with the disease in family I as it was also  
17 present in the affected sister I-2. (**Figure 1A**)

18 Additionally, one case with isolated cerebellar ataxia carried the AAGGG expansion along  
19 with an ACGGG repeat, which was also absent in controls. However, Sanger sequencing  
20 showed that the ACGGG expansion was only 50 repeats, which is considerably below the lower  
21 limit of pathogenicity for the pathogenic AAGGG motifs of 250 repeats and was therefore  
22 considered likely non-pathogenic in this case. Notably, the patient exhibited isolated cerebellar  
23 ataxia but no neuropathy, which is unusual in RFC1 disease.

24 We next screened by RP-PCR an internal cohort of 540 DNA samples from cases with sensory  
25 neuropathy, ataxia or CANVAS and identified five additional cases carrying an AGGGC  
26 expansion (cases III-1, IV-1, V-1, V-2 and VI-1) and three cases carrying AAAGG expansions  
27 on the second allele (cases X-1, XI-1, XII-1) (**Table 2**). We did not identify additional AGAGG  
28 or AAGGC repeat expansion carriers. All cases were of self-reported Caucasian ethnicity.

29 Based on Southern blotting, OGM or long-read sequencing (**Figure 1B-C**), when available, we  
30 observed that the repeat size of the rare AGGGC, AAGGC and AGAGG was >600 repeats in  
31 all cases (mean $\pm$ SD, 892 $\pm$ 247) (**Figure 2A**).

1 Also, enough DNA for Southern blotting was available from 5 patients with CANVAS  
 2 spectrum as defined by the presence of sensory neuropathy and at least one of the additional  
 3 features of the full syndrome (cerebellar dysfunction, vestibular areflexia, cough) (cases VI-X)  
 4 and 8 controls carrying compound heterozygous AAGGG/AAAGG expansions (**Figure 2B**).

5 In CANVAS patients, the AAAGG expansions were always  $\geq 600$  repeats (mean $\pm$ SD,  
 6  $979\pm 257$ ), and were significantly larger compared to AAAGG expansions ( $238\pm 142$  repeat  
 7 units) found in controls ( $p < 0.0001$ ), suggesting that, although the AAAGG repeat is usually  
 8 small and non-pathogenic, as shown in **Figure 2A**, larger AAAGG repeat expansion occur and  
 9 may have a pathogenic role.

10

### 11 **Long read sequencing confirms the sequence of the expanded repeats**

12 To gain further insight into the exact sequence of the novel pathogenic motifs, we performed  
 13 targeted long-read sequencing (**Figure 1D and Supplementary Table 1**).

14 We confirmed the presence of uninterrupted AGGGC<sub>1240</sub> in case II-1 and AGGGC<sub>3200</sub> in case  
 15 III-1. Moreover, long-read sequencing was able to accurately define the exact repeat  
 16 composition of AGAGG and AAGGC expansions, which revealed the presence of mixed  
 17 repeat motifs (AAGGC)<sub>900</sub>(AAGGG)<sub>940</sub> and (AGAGG)<sub>470</sub>(AAAGG)<sub>470</sub>, in case I-1 and VII-1,  
 18 respectively. Long-read sequencing was also performed in five cases carrying large AAAGG  
 19 expansion and showed the presence of uninterrupted AAAGG motifs in three (Cases X-1, XI-  
 20 1 and XII-1), with sizes of 980, 800 and 600 repeat units, respectively, while two probands  
 21 (case VIII-1 and IX-1) carried complex (AAAGG)<sub>610</sub>(AAGGG)<sub>390</sub>  
 22 and (AAAGG)<sub>700</sub>(AAGGG)<sub>200</sub> repeats.

23

### 24 **All pathogenic repeat configurations share an ancestral haplotype**

25 Subsequently, we looked at the inferred haplotypes associated with the novel pathogenic repeat  
 26 motifs. A region of 66kb (**Figure 3**, between markers B and C, chr4:39302305-39366034,  
 27 hg38) is shared among all pathogenic alleles. It is worth noting that a larger region of 207 kb  
 28 (between markers A and C), and which contains *WDR19* and *RFC1* genes, is shared among all  
 29 pathogenic alleles but one (Case III-1), where the haplotype becomes the same as the wild type  
 30 allele. This suggests a probable more recent recombination event at marker B for Case III-1.  
 31 The larger shared region identified in carriers of the novel pathogenic configurations, as well

1 as in AAGGG and AAAGG carriers, supports the existence of an ancestral haplotype which  
2 gave rise to these expanded alleles. Notably, non-pathogenic AAAAG<sup>(9-11)</sup> and expanded  
3 AAAAG repeat originated from a different haplotype.

4 We estimated that the ancestral haplotype which gave rise to different pathogenic repeat  
5 configurations in RFC1 likely dates back to 56,100 years ago (CI at 95% 27,680-115,580 years)

## 6 **Clinical features of patients carrying novel pathogenic repeat configurations** 7 **in *RFC1***

8 We found 14 patients from 12 families carrying novel pathogenic *RFC1* repeat configurations.  
9 Demographic and clinical characteristics of patients are summarized in **Table 2**. All patients  
10 were Europeans, apart from patient I-1 and I-2 who were from India and patient X-1 who was  
11 from Australia. Mean age of onset was  $51.5 \pm 13.7$  (24-73) years, and mean disease duration  
12 at examination was  $17.2 \text{ years} \pm 8.7$  (3-34) years. Six patients had isolated sensory neuropathy,  
13 which was associated with cough in four of them, one patient had sensory neuropathy and  
14 vestibular dysfunction, while seven cases had full CANVAS. Additional features were  
15 observed in some cases including early onset and rapid progression (case I-1), cognitive  
16 impairment (III-1, VI-1), muscle cramps (I-1, II-1, III-1 and IV-1), and REM sleep behaviour  
17 disorder with positive dopamine transporter scan (DatScan) (IX-1). Autonomic dysfunction  
18 was observed in six cases and in two of them (II-1, III-1), who both carried AGGGC expansion,  
19 was severe and led to syncopal episodes. Detailed descriptions of the clinical features are  
20 provided in **supplementary note**.

21

## 22 **Pathogenic configurations in *RFC1* are predicted to form G quadruplexes**

23 As repetitive G-rich sequences are known to form G quadruplexes (G4)<sup>32,35,36</sup>, a secondary  
24 DNA structure which act as transcriptional regulator by impeding transcription factor binding  
25 to duplex-DNA or stalling the progression of RNA polymerase, we set out to evaluate the  
26 propensity of the different repeat configurations in *RFC1* to form G4.

27 All pathogenic repeat configurations showed high G4 scores, in the range observed for the  
28 well-known G4-forming regions of the *cMYC*<sup>37</sup> and *HRAS1*<sup>38</sup> genes, as predicted by QGRS-  
29 Mapper and G4Hunter in contrast to the non-pathogenic AAAAG (**Table 3**).

30

## 1 Discussion

2 We leveraged WGS data from nearly 10,000 individuals recruited to the Genomics England  
3 sequencing project to investigate the normal and pathogenic variation of the RFC1 repeat. We  
4 identified three novel repeat configurations associated with CANVAS and disease spectrum,  
5 including AGGGC, AAGGC and AGAGG. Notably, we also showed a pathogenic role for  
6 large uninterrupted or interrupted AAAGG expansions, while AAAAG, AAGAG and  
7 AAAGGG expansions are likely always benign (**Figure 4**).

8 Most pathogenic repeat expansion were found in individuals of Caucasian ancestry, however  
9 ACAGG seems to be common in East Asians, while AAGGC was identified in a family of  
10 South Asian ancestry. Interestingly most pathogenic repeats seem to have arisen from a shared  
11 region 207 Kb, supporting their origin from a common ancestor who lived ~50,000 years ago.  
12 *Rafehi et al.*<sup>2</sup> had previously identified a larger ancestral haplotype in Australian patients  
13 affected by CANVAS of 360 Kb and estimated that the most recent common ancestor lived  
14 approximately 25,880 (CI: 14,080–48,020) years ago.<sup>2</sup> In our study, the inclusion of additional  
15 pathogenic repeat configurations and multiple ethnicities allowed the identification of a smaller  
16 core haplotype and has extended further back in time the origin of the common ancestor  
17 carrying a pathogenic repeat in *RF1*. It is reasonable to believe that occurrence of subsequent  
18 A-G transition and A-G or G-C transversion in the polyA tail of the AluSx3 element on the  
19 ancestral haplotype has favoured the further expansion of GC rich motifs over the millennia.  
20 Since the most significant recent wave out of Africa is estimated to have taken place about  
21 70,000–50,000 years ago, we can speculate that the repeat containing haplotype spread with  
22 the migration of early modern humans from Africa through the Near East and to the rest of the  
23 world.

24 Patients showed clinical features undistinguishable from those of patients carrying biallelic  
25 AAGGG expansion. In some cases, however, the disease appeared to be more severe due to  
26 symptomatic dysautonomia, early cerebellar involvement or disabling gait disturbance.

27 The identification of these motifs has direct clinical implications. Given their frequency, RP-  
28 PCR for AAAGG and AGGGC should be considered in all cases. Particular attention should  
29 be paid to carriers of compound AAGGG/AAAGG expansions and accurate sizing and full  
30 sequencing of the satellite through long read is recommended to establish its possible  
31 pathogenicity. In addition, depending on availability, Southern blotting, genome optical  
32 mapping, or long read sequencing, are warranted in patients with a suggestive clinical

1 phenotype but inconclusive screening, such as in cases with absence of PCR amplifiable  
2 product on flanking PCR but negative RP-PCR for AAGGG expansion).

3 The findings of this study highlight the genetic complexity of *RFC1*-related disease and lend  
4 support to the hypothesis that the size and GC content of the pathogenic repeat is more  
5 important than the exact repeat motif. Consistently, all pathogenic repeat configurations are  
6 rich in G content and are predicted to form highly stable G quadruplexes, which have been  
7 previously demonstrated to affect gene transcription in other pathogenic conditions.<sup>35,36</sup>

8 Both Nanopore or Pacbio sequencing platforms and either the targeted CRISPR/Cas9 or  
9 adaptive selection approach were used to increase accuracy of the sequencing of *RFC1* repeat  
10 locus. Despite several attempts and similarly to other large satellites, long read sequencing of  
11 the *RFC1* repeat remains challenging and, depending on the specific configurations, size, and  
12 DNA quality, only few reads were available for analysis in some cases. Notably, an uneven  
13 coverage at the *RFC1* locus across samples was also observed in a recent study looking at  
14 *RFC1* repeat composition through Nanopore sequencing. The authors attributed the variability  
15 to variable degrees of DNA fragmentation depending on the delay between blood sampling and  
16 DNA extraction. Hopefully, constant advancements in long read sequencing platforms and  
17 decrease in cost (currently ~1,000 USD per sample) will soon translate into increased  
18 accessibility and higher accuracy of this technology.

19 In conclusion, the study expanded the genetic heterogeneity underlying *RFC1* CANVAS and  
20 disease spectrum, and identified three additional pathogenic AAGGC, AGGGC and AGAGG  
21 repeat motifs. We also demonstrated a pathogenic role for large uninterrupted or interrupted  
22 AAAGG expansions, thereby highlighting the importance of sizing and, if possible, full  
23 sequencing of the *RFC1* satellite expansion in clinically selected cases, in order to correctly  
24 diagnose and counsel patients and their families.

25

## 26 **Appendix 1**

### 27 **Genomics England Research Consortium**

28 J. C. Ambrose, P. Arumugam, E. L. Baple, M. Bleda, F. Boardman-Pretty, J. M. Boissiere, C.  
29 R. Boustred, H. Brittain, M. J. Caulfield, G. C. Chan, C. E. H. Craig, L. C. Daugherty, A. de  
30 Burca, A. Devereau, G. Elgar, R. E. Foulger, T. Fowler, P. Furió-Tarí, E. Gustavsson, J. M.  
31 Hackett, D. Halai, A. Hamblin, S. Henderson, J. E. Holman, T. J. P. Hubbard, K. Ibáñez, R.

1 Jackson, L. J. Jones, D. Kasperaviciute, M. Kayikci, L. Lahnstein, K. Lawson, S. E. A. Leigh,  
2 I. U. S. Leong, F. J. Lopez, F. Maleady-Crowe, J. Mason, E. M. McDonagh, L. Moutsianas,  
3 M. Mueller, N. Murugaesu, A. C. Need, C. A. Odhams, C. Patch, D. Perez-Gil, D.  
4 Polychronopoulos, J. Pullinger, T. Rahim, A. Rendon, P. Riesgo-Ferreiro, T. Rogers, M. Ryten,  
5 B. Rugginini, K. Savage, K. Sawant, R. H. Scott, A. Siddiq, A. Sieghart, D. Smedley, K. R.  
6 Smith, A. Sosinsky, W. Spooner, H. E. Stevens, A. Stuckey, R. Sultana, E. R. A. Thomas, S.  
7 R. Thompson, C. Tregidgo, A. Tucci, E. Walsh, S. A. Watters, M. J. Welland, E. Williams, K.  
8 Witkowska, S. M. Wood, M. Zarowiecki.

9 Further details are available in the Supplementary material.

10

11

## 12 **Funding**

13 This work was supported by the Medical Research Council (MR/T001712/1), Fondazione  
14 Cariplo (grant n. 2019-1836), the Inherited Neuropathy Consortium, Fondazione Regionale per  
15 la Ricerca Biomedica (Regione Lombardia, project ID 1751723), National Ministry of Health  
16 (Ricerca Corrente 2021-2022) and Italian Ministry for Universities and Research (MUR,  
17 20229MMHXP) to awarded to A.C. This work has also been supported by a Medical Research  
18 Future Fund Genomics Health Futures Mission grant (APP2007681) awarded to M.L.K. and  
19 S.V. and grant CP 20/2018 from the Fondazione Regionale per la Ricerca Biomedica to F.T.  
20 F.A. was supported by NIH Intramural Research Program, the US National Institute on Aging.  
21 E.A. was partially supported by the Telethon Foundation and by the Rotary Club Milano Ovest.

22

## 23 **Competing interests**

24 The authors report no competing interests.

25

## 26 **Supplementary material**

27 Supplementary material is available at *Brain* online

28

## 1 **References**

- 2 1. Cortese A, Simone R, Sullivan R, et al. Biallelic expansion of an intronic repeat in  
3 RFC1 is a common cause of late-onset ataxia. *Nat Genet.* 2019;51(4):649–58.
- 4 2. Rafehi H, Szmulewicz DJ, Bennett MF, et al. Bioinformatics-Based Identification of  
5 Expanded Repeats: A Non-reference Intronic Pentamer Expansion in RFC1 Causes CANVAS.  
6 *Am J Hum Genet.* 2019 Jun 12;
- 7 3. Cortese A, Reilly MM, Houlden H. RFC1 CANVAS / Spectrum Disorder. In: Adam  
8 MP, Ardinger HH, Pagon RA, Wallace SE, Bean LJ, Stephens K, et al., editors. *GeneReviews®*  
9 [Internet]. Seattle (WA): University of Washington, Seattle; 1993 [cited 2021 Jan 5]. Available  
10 from: <http://www.ncbi.nlm.nih.gov/books/NBK564656/>
- 11 4. Cortese A, Curro' R, Vegezzi E, Yau WY, Houlden H, Reilly MM. Cerebellar ataxia,  
12 neuropathy and vestibular areflexia syndrome (CANVAS): genetic and clinical aspects. *Pract*  
13 *Neurol.* 2022 Feb;22(1):14–8.
- 14 5. Cortese A, Tozza S, Yau WY, et al. Cerebellar ataxia, neuropathy, vestibular areflexia  
15 syndrome due to RFC1 repeat expansion. *Brain.* 2020 Feb 1;143(2):480–90.
- 16 6. Currò R, Salvalaggio A, Tozza S, et al. RFC1 expansions are a common cause of  
17 idiopathic sensory neuropathy. *Brain.* 2021 Jun 22;144(5):1542–50.
- 18 7. Kumar KR, Cortese A, Tomlinson SE, et al. RFC1 expansions can mimic hereditary  
19 sensory neuropathy with cough and Sjögren syndrome. *Brain.* 2020 Oct 1;143(10):e82.
- 20 8. Ronco R, Perini C, Currò R, et al. Truncating Variants in RFC1 in Cerebellar Ataxia,  
21 Neuropathy, and Vestibular Areflexia Syndrome. *Neurology.* 2023 Jan 31;100(5):e543–54.
- 22 9. Benkirane M, Da Cunha D, Marelli C, et al. RFC1 nonsense and frameshift variants  
23 cause CANVAS: clues for an unsolved pathophysiology. *Brain.* 2022 Nov 21;145(11):3770–  
24 5.
- 25 10. Huin V, Coarelli G, Guemy C, et al. Motor neuron pathology in CANVAS due to RFC1  
26 expansions. *Brain.* 2021 Dec 20;awab449.
- 27 11. Traschütz A, Cortese A, Reich S, et al. Natural History, Phenotypic Spectrum, and  
28 Discriminative Features of Multisystemic RFC1 Disease. *Neurology.* 2021 Mar  
29 2;96(9):e1369–82.



- 1 12. About Syriani D, Wong D, Andani S, et al. Prevalence of RFC1-mediated  
2 spinocerebellar ataxia in a North American ataxia cohort. *Neurol Genet.* 2020 Jun;6(3):e440.
- 3 13. Beijer D, Dohrn MF, De Winter J, et al. RFC1 repeat expansions: A recurrent cause of  
4 sensory and autonomic neuropathy with cough and ataxia. *Eur J Neurol.* 2022 Jul;29(7):2156–  
5 61.
- 6 14. Gisatulin M, Dobricic V, Zühlke C, et al. Clinical spectrum of the pentanucleotide  
7 repeat expansion in the RFC1 gene in ataxia syndromes. *Neurology.* 2020 Sep 1;
- 8 15. Tagliapietra M, Cardellini D, Ferrarini M, et al. RFC1 AAGGG repeat expansion  
9 masquerading as Chronic Idiopathic Axonal Polyneuropathy. *J Neurol.* 2021  
10 Nov;268(11):4280–90.
- 11 16. Montaut S, Diedhiou N, Fahrer P, et al. Biallelic RFC1-expansion in a French  
12 multicentric sporadic ataxia cohort. *J Neurol.* 2021 Sep;268(9):3337–43.
- 13 17. Van Daele SH, Vermeer S, Van Eesbeeck A, et al. Diagnostic yield of testing for RFC1  
14 repeat expansions in patients with unexplained adult-onset cerebellar ataxia. *J Neurol*  
15 *Neurosurg Psychiatry.* 2020 Jul 30;
- 16 18. Ghorbani F, de Boer-Bergsma J, Verschuuren-Bemelmans CC, et al. Prevalence of  
17 intronic repeat expansions in RFC1 in Dutch patients with CANVAS and adult-onset ataxia. *J*  
18 *Neurol.* 2022 Nov;269(11):6086–93.
- 19 19. Erdmann H, Schöberl F, Giurgiu M, et al. Parallel in-depth analysis of repeat  
20 expansions in ataxia patients by long-read sequencing. *Brain.* 2022 Oct 13;awac377.
- 21 20. Beecroft SJ, Cortese A, Sullivan R, et al. A Māori specific RFC1 pathogenic repeat  
22 configuration in CANVAS, likely due to a founder allele. *Brain.* 2020 Sep 1;143(9):2673–80.
- 23 21. Scriba CK, Beecroft SJ, Clayton JS, et al. A novel RFC1 repeat motif (ACAGG) in two  
24 Asia-Pacific CANVAS families. *Brain.* 2020 Oct 1;143(10):2904–10.
- 25 22. Nakamura H, Doi H, Mitsuhashi S, et al. Long-read sequencing identifies the  
26 pathogenic nucleotide repeat expansion in RFC1 in a Japanese case of CANVAS. *J Hum Genet.*  
27 2020 May;65(5):475–80.
- 28 23. Miyatake S, Yoshida K, Koshimizu E, et al. Repeat conformation heterogeneity in  
29 cerebellar ataxia, neuropathy, vestibular areflexia syndrome. *Brain.* 2022 Apr 29;145(3):1139–  
30 50.

- 1 24. 100,000 Genomes Project Pilot Investigators, Smedley D, Smith KR, Martin A, et al.  
2 100,000 Genomes Pilot on Rare-Disease Diagnosis in Health Care - Preliminary Report. *N*  
3 *Engl J Med*. 2021 Nov 11;385(20):1868–80.
- 4 25. Chen X, Schulz-Trieglaff O, Shaw R, et al. Manta: rapid detection of structural variants  
5 and indels for germline and cancer sequencing applications. *Bioinformatics*. 2016 Apr  
6 15;32(8):1220–2.
- 7 26. Stevanovski I, Chintalaphani SR, Gamaarachchi H, et al. Comprehensive genetic  
8 diagnosis of tandem repeat expansion disorders with programmable targeted nanopore  
9 sequencing. *Sci Adv*. 2022 Mar 4;8(9):eabm5386.
- 10 27. Payne A, Holmes N, Clarke T, Munro R, Debebe BJ, Loose M. Readfish enables  
11 targeted nanopore sequencing of gigabase-sized genomes. *Nat Biotechnol*. 2021  
12 Apr;39(4):442–50.
- 13 28. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018  
14 Sep 15;34(18):3094–100.
- 15 29. Gamaarachchi H, Samarakoon H, Jenner SP, et al. Fast nanopore sequencing data  
16 analysis with SLOW5. *Nat Biotechnol*. 2022 Jul;40(7):1026–9.
- 17 30. Delaneau O, Zagury JF, Robinson MR, Marchini JL, Dermitzakis ET. Accurate,  
18 scalable and integrative haplotype estimation. *Nat Commun*. 2019 Nov 28;10(1):5436.
- 19 31. Gandolfo LC, Bahlo M, Speed TP. Dating Rare Mutations from Small Samples with  
20 Dense Marker Data. *Genetics*. 2014 Aug 1;197(4):1315–27.
- 21 32. Frasson I, Pirota V, Richter SN, Doria F. Multimeric G-quadruplexes: A review on their  
22 biological roles and targeting. *Int J Biol Macromol*. 2022 Apr 15;204:89–102.
- 23 33. Kikin O, D'Antonio L, Bagga PS. QGRS Mapper: a web-based server for predicting  
24 G-quadruplexes in nucleotide sequences. *Nucleic Acids Res*. 2006 Jul 1;34(Web Server  
25 issue):W676-682.
- 26 34. Bedrat A, Lacroix L, Mergny JL. Re-evaluation of G-quadruplex propensity with  
27 G4Hunter. *Nucleic Acids Res*. 2016 Feb 29;44(4):1746–59.
- 28 35. Varshney D, Spiegel J, Zyner K, Tannahill D, Balasubramanian S. The regulation and  
29 functions of DNA and RNA G-quadruplexes. *Nat Rev Mol Cell Biol*. 2020 Aug;21(8):459–74.

- 1 36. Wang E, Thombre R, Shah Y, Latanich R, Wang J. G-Quadruplexes as pathogenic  
2 drivers in neurodegenerative disorders. *Nucleic Acids Res.* 2021 May 21;49(9):4816–30.
- 3 37. Dickerhoff J, Dai J, Yang D. Structural recognition of the MYC promoter G-quadruplex  
4 by a quinoline derivative: insights into molecular targeting of parallel G-quadruplexes. *Nucleic*  
5 *Acids Research.* 2021 Jun 4;49(10):5905–15.
- 6 38. Cogo S, Shchekotikhin AE, Xodo LE. HRAS is silenced by two neighboring G-  
7 quadruplexes and activated by MAZ, a zinc-finger transcription factor with DNA unfolding  
8 property. *Nucleic Acids Res.* 2014 Sep 1;42(13):8379–88.
- 9  
10

ACCEPTED MANUSCRIPT

## 1 **Figure legends**

2 **Figure 1 Long read sequencing defines the precise sequence of the novel pathogenic RFC1**  
3 **motifs.** Pedigrees with arrow and P indicating proband (**panel A**), RP-PCR plots and, where  
4 available, southern blot images and optical genome mapping plots (**panel B**), long-read  
5 sequencing results (**panel C**) of representative patients with AAGGC, AGGGC, AGAGG and  
6 AAAGG expansions (cases I-1, III-1, VII-1, XII-1). In case III-1 only partial reads, which did  
7 not span the entire RFC1 repeat locus, could be obtained from the AAGGG allele.

8

9 **Figure 2 RFC1 repeat expansion size.** A) Comparison of repeat sizes of alleles carrying  
10 AAGGG, AAAGG, AAGGC, AGGGC and AGAGG expansion from this and previous  
11 studies.<sup>1,5,6</sup> The dotted lines refer to the smallest pathogenic of 250 AAGGG repeat identified  
12 so far B) Comparison of the AAAGG repeat sizes in the compound heterozygous state with the  
13 AAGGG expansion, in patients with CANVAS and disease spectrum versus controls.

14

15 **Figure 3 A shared ancestral haplotype in patients with pathogenic RFC1 motifs.** Graphical  
16 representation of the haplotypes associated with AAGGG, ACAGG and novel pathogenic  
17 repeat motifs identified in this study. For each SNP, the reference allele is represented in blue,  
18 while the alternative allele is represented in yellow. The repeat expansion locus is marked with  
19 a red line (R). There is a shared region (B-C, -rs2066782-rs6851075, chr4:39302305-39366034,  
20 hg38) of 66kb for all novel configurations. A larger region of 207kb (A-C, rs148316325-  
21 rs6851075, chr4:39158847-39366034, hg38), which is flanked by two recombination hotspots  
22 (arrows), is also shared among all but one allele for Case III-1 suggesting a recombination  
23 event at B (rs2066782) in this family. The shared haplotype lies in a region of low-  
24 recombination rate (HapMap data) and is delimited by small peaks at A and C. A smaller  
25 increase in recombination rate is also visible at B.

26

27 **Figure 4 Normal and pathogenic significance of repeat expansion motifs at RFC1 locus.**

28

29

30

1 **Table 1 Normal and pathogenic variation of the RFC1 repeat locus on the 100,000 Genome Project**

	<b>Hereditary ataxia (N=893)</b>	<b>Non neurological controls (N=8107)</b>	<b>P values</b>
<b>Rare homozygous (&lt;1%) repeat expansions present in ataxia cases and absent in controls</b>			
ACAGG (hom)	1 (0.01%)	0 (0%)	-
<b>AAGGC (hom)</b>	1 (0.01%)	0 (0%)	-
<b>Repeat expansion found in compound heterozygous state with AAGGG expansions (allele 1/allele2)</b>			
AAGGG/AAAAG	21 (2.3%)	248 (3%)	Ns
AAGGG/AAAGGG	5 (0.6%)	32 (0.4%)	Ns
AAGGG/AAGAG	3 (0.3%)	16 (0.2%)	Ns
AAGGG/ <b>AAAGG</b>	10 (1.1%)	47 (0.6%)	0.05
AAGGG/ACGGG*	1 (0.01%)	0 (0%)	-
AAGGG/ <b>AGAGG</b>	1 (0.01%)	0 (0%)	-
AAGGG/ <b>AGGGC</b>	1 (0.01%)	0 (0%)	-

2 Ns=not significant. \*Small (ACGGG)<sub>50</sub> expansion in typical non-pathogenic range (10-220). Novel pathogenic repeat motifs identified in this  
3 study are bold highlighted.  
4

1

**Table 2 Clinical and demographic features of patients carrying novel pathogenic repeat configurations in RFC1**

	RFC1 genotype	Sex	Ethnicity	Phenotype	Age of onset	Disease duration (years)	Chronic cough	Cerebellar syndrome	Sensory neuropathy	Bilateral vestibular areflexia	Dysautonomia	Use of walking aids (age)	Additional features
<b>AAGGC expansion</b>													
Caselle-I	Allele 1: (AAGGG) <sub>510</sub> (AAGGC) <sub>880</sub> Allele 2: (AAGGG) <sub>940</sub> (AAGGC) <sub>900</sub>	F	Caucasian (India n)	CANVAS	24	17	Yes	Yes	Yes	Yes	No	Stick (36)	Cramps, pyramidal signs
Caselle-2	Allele 1: (AAGGG) <sub>n</sub> (AAGGC) <sub>n</sub> Allele 2: (AAGGG) <sub>n</sub> (AAGGC) <sub>n</sub>	F	Caucasian (India n)	Sensory neuropathy + cough	34	8	Yes	N/A	Yes	N/A	N/A	No	/
<b>AGGGC expansion</b>													
Caselle-III	Allele 1: (AGGGC) <sub>1240</sub> Allele 2: (AAGGG) <sub>930</sub>	M	Mixed (Lebanese)	Sensory neuropathy + vestibular dysfunction	53	11	Yes	No	Yes	Yes	Yes	No	Cramps
Caselle-III	Allele 1: (AGGGC) <sub>3200</sub> Allele 2: (AAGGG) <sub>1000</sub>	M	Caucasian (British)	CANVAS	71	12	Yes	Yes	Yes	N/A	Yes	Wheelchair (81)	Cramps, cognitive/behavioural abnormalities after age 80
Caselle-IV	Allele 1: (AGGGC) <sub>1875</sub> Allele 2: (AAGGG) <sub>500</sub>	M	Caucasian (Italian)	CANVAS	41	34	No	Yes	Yes	Yes	Yes	Wheelchair (72)	Cramps
Caselle-V	Allele 1: (AGGGC) <sub>n</sub> Allele 2: (AAGGG) <sub>n</sub>	F	Caucasian (Italian)	Sensory neuropathy + cough	60	13	Yes	No	Yes	No	No	No	/
Caselle-V	Allele 1: (AGGGC) <sub>n</sub> Allele 2: (AAGGG) <sub>n</sub>	F	Caucasian (Italian)	Sensory neuropathy	40	20	No	No	Yes	No	No	No	/
Caselle-VI	Allele 1: (AGGGC) <sub>n</sub> Allele 2: (AAGGG) <sub>n</sub>	F	Caucasian (Italian)	Sensory gangliopathy + cough	62	23	Yes	No	Yes	N/A	Yes	No	Voice and hand tremor, urinary incontinence
<b>AGAGG expansion</b>													
Caselle-VI	Allele 1: (AAAGG) <sub>470</sub> (AGAGG) <sub>470</sub> Allele 2: (AAGGG) <sub>1140</sub>	F	Caucasian (British)	CANVAS	50	24	Yes	Yes	Yes	Yes	No	Walker (69), wheelchair (74)	/
<b>AAAGG expansion</b>													
Caselle-VI	Allele 1: (AAAGG) <sub>610</sub> (AAGGG) <sub>390</sub> Allele 2: (AAGGG) <sub>110</sub>	M	Caucasian (British)	CANVAS	55	20	Yes	Yes	Yes	N/A	Yes	Walker and wheelchair (74)	Cognitive impairment since age 72

	0													
C as e IX -I	Allele 1: (AAGGG) <sub>700</sub> (AAGGG) <sub>200</sub> / Allele 2: (AAGGG) <sub>117</sub>	M	Cauc asian (Britis h)	CANV AS	45	31	Yes	Yes	Yes	Yes	Yes	Walker (75)	RBD, positive DatScan	
C as e X- I	Allele 1: (AAAGG) <sub>980</sub> Allele 2: (AAGGG) <sub>101</sub>	M	Cauc asian (Aust ralian )	CANV AS	58	15	Yes	Yes	Yes	Yes	N/A	N/A	/	
C as e XI -I	Allele 1: (AAAGG) <sub>800</sub> Allele 2: (AAGGG) <sub>500</sub>	F	Cauc asian (Italia n)	Sensor y ganglio nopath y + cough	73	3	Yes	No	Yes	No	No	Stick (77)	/	
C as e XI I- I	Allele 1: (AAAGG) <sub>600</sub> Allele 2: (AAGGG) <sub>390</sub>	M	Cauc asian (Italia n)	Sensor y ganglio nopath y + cough	56	10	Yes	No	Yes	No	No	No	/	

1  
2

ACCEPTED MANUSCRIPT

1  
2**Table 3 Pathogenic RFC I motifs are predicted to form G quadruplexes**

Gene - Analyzed Sequences	QGRS-Mapper Score	G4Hunter Score
<i>RFC1</i> – (AGGGC) <sub>10</sub>	42	1.83
<i>RFC1</i> – (AAGGG) <sub>10</sub>	42	2.00
<i>RFC1</i> – (AAGGC) <sub>10</sub>	21	1.82
<i>RFC1</i> – (AAAGG) <sub>10</sub>	21	0.94
<i>RFC1</i> – (AGAGG) <sub>10</sub>	21	1.12
<i>RFC1</i> – (AAAAG) <sub>10</sub>	No putative G4 identified	
<i>c-MYC</i> – TGGGGAGGTGGGGAGGGTGGGGAAGG	41	2.59
<i>HRAS-1</i> – TCGGGTTGCGGCGCAGGCACGGGCG	41	1.19

G-score values comparison between repeat configurations found in *RFC1* and well-known G4s-forming sequences.

3  
4  
5  
6  
7

ACCEPTED MANUSCRIPT



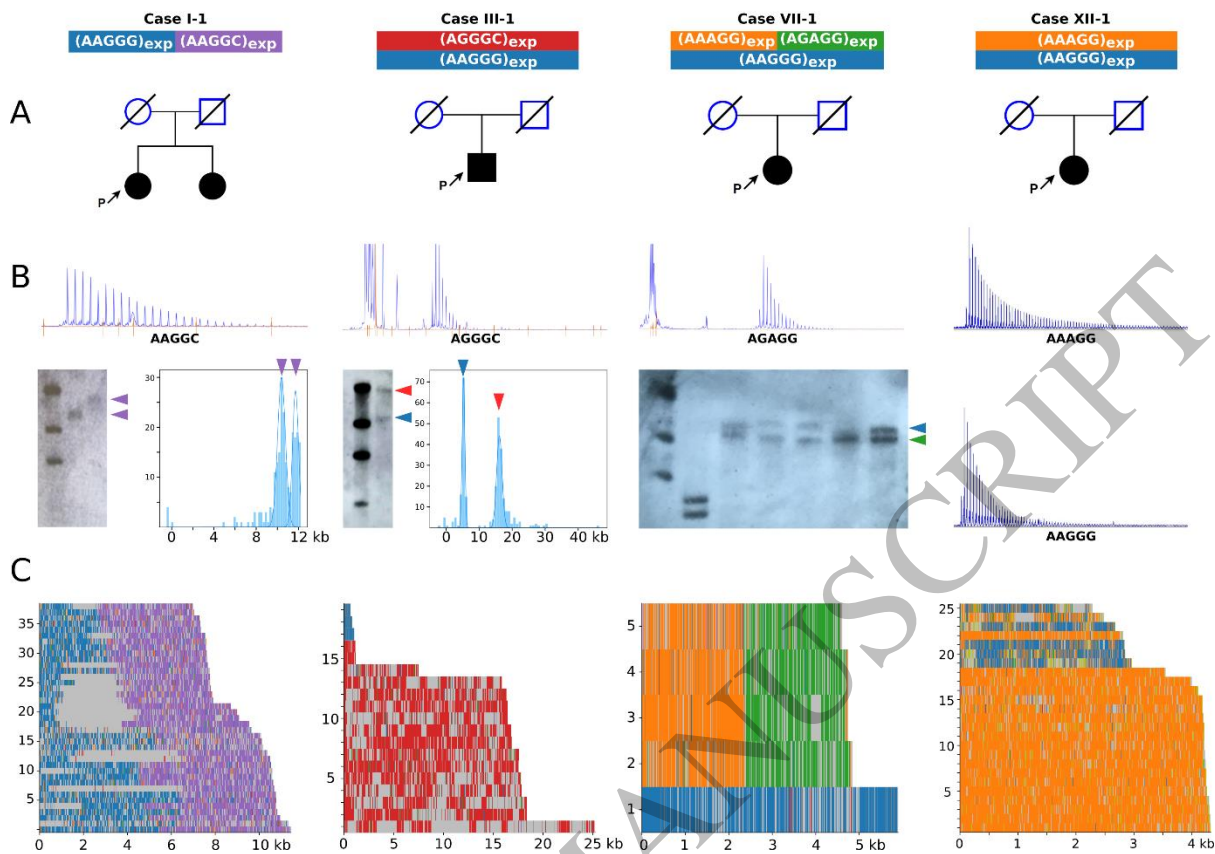


Figure 1  
543x381 mm (x DPI)

1  
2  
3  
4

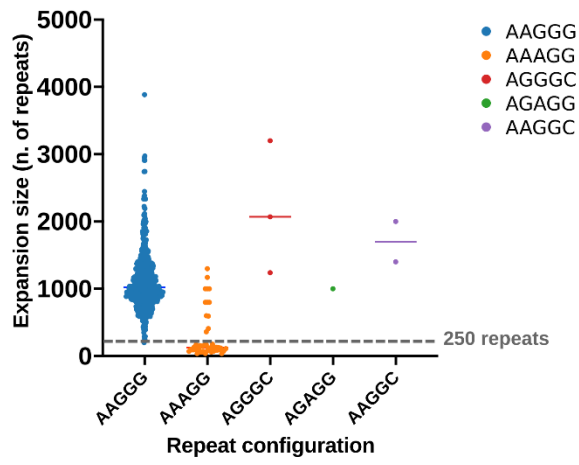
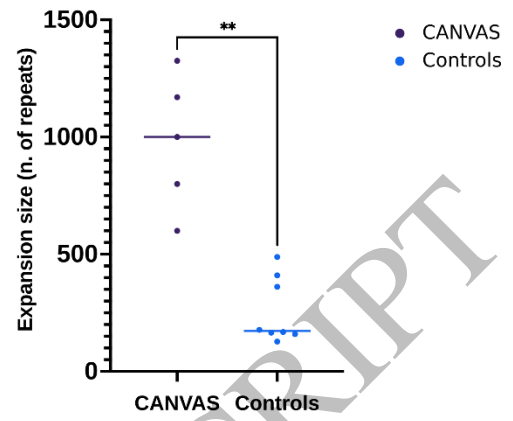
**A****B**

Figure 2  
543x248 mm (x DPI)

1  
2  
3  
4

ACCEPTED MANUSCRIPT

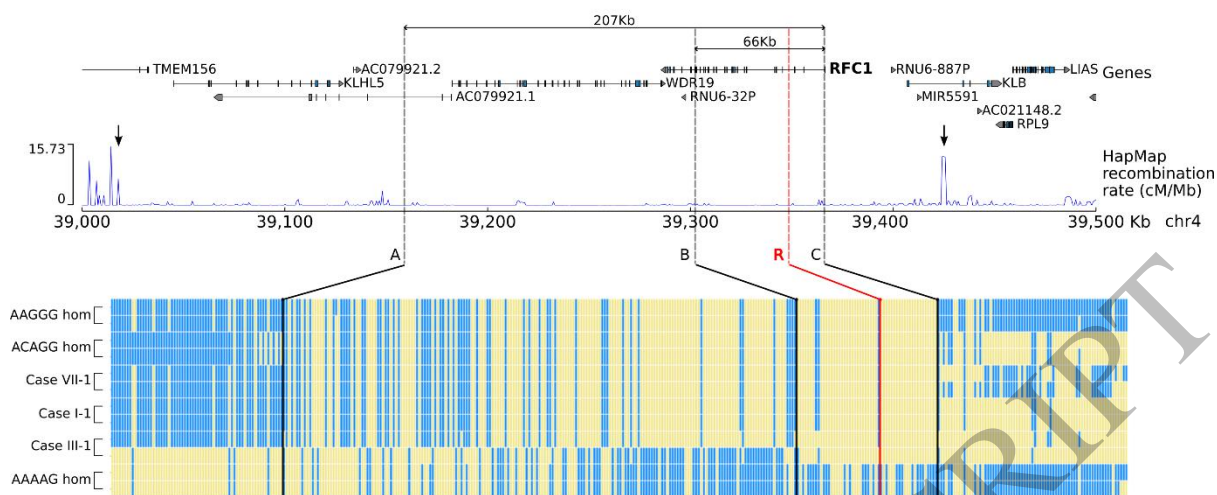


Figure 3  
543x223 mm (x DPI)

1  
2  
3  
4

ACCEPTED MANUSCRIPT

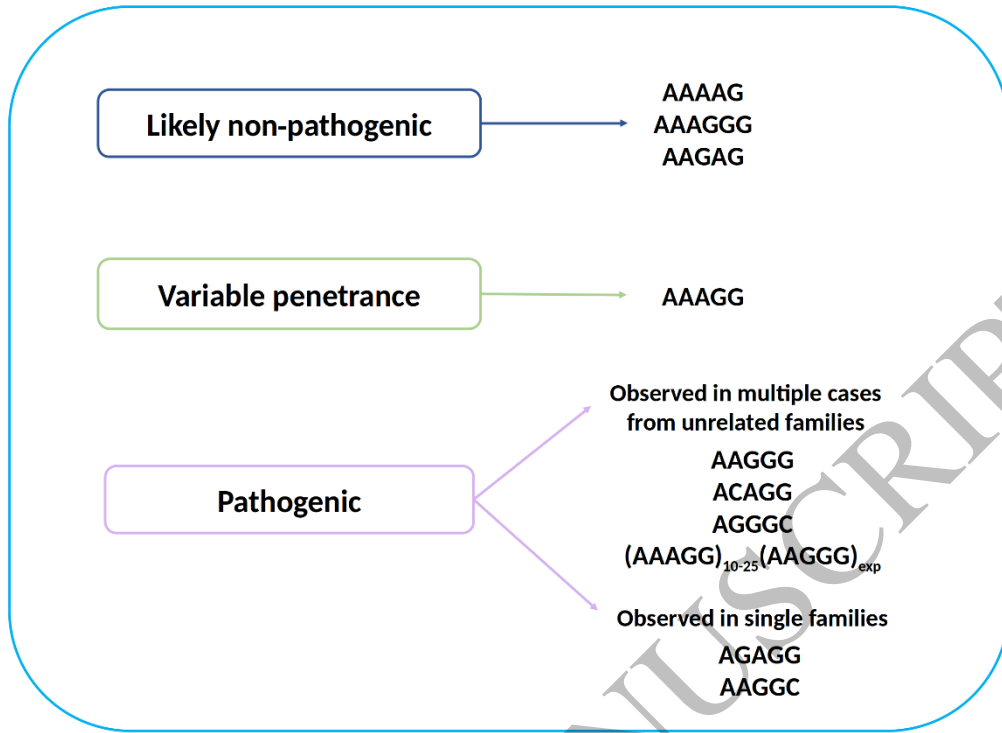


Figure 4  
559x347 mm (x DPI)

1  
2  
3

ACCEPTED MANUSCRIPT