

## Supplemental Data

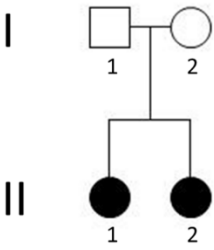
### Bioinformatics-Based Identification of Expanded

### Repeats: A Non-reference Intronic Pentamer

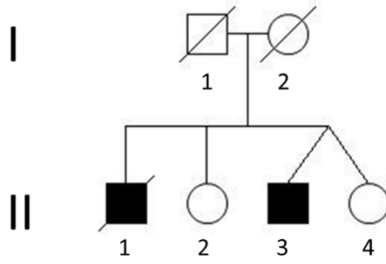
### Expansion in *RFC1* Causes CANVAS

Haloom Rafehi, David J. Szmulewicz, Mark F. Bennett, Nara L.M. Sobreira, Kate Pope, Katherine R. Smith, Greta Gillies, Peter Diakumis, Egor Dolzhenko, Michael A. Eberle, María García Barcina, David P. Breen, Andrew M. Chancellor, Phillip D. Cremer, Martin B. Delatycki, Brent L. Fogel, Anna Hackett, G. Michael Halmagyi, Solange Kapetanovic, Anthony Lang, Stuart Mossman, Weiyi Mu, Peter Patrikios, Susan L. Perlman, Ian Rosemergy, Elsdon Storey, Shaun R.D. Watson, Michael A. Wilson, David S. Zee, David Valle, David J. Amor, Melanie Bahlo, and Paul J. Lockhart

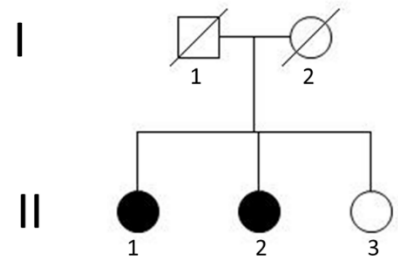
CANVAS1



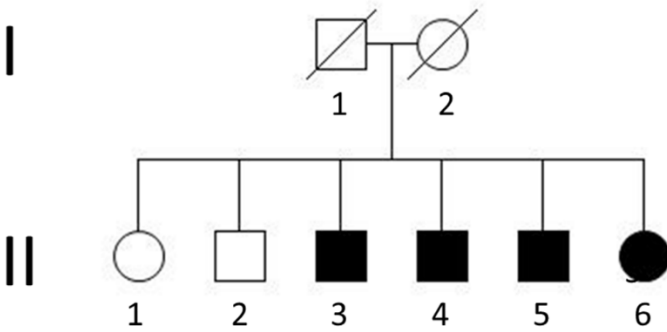
CANVAS2



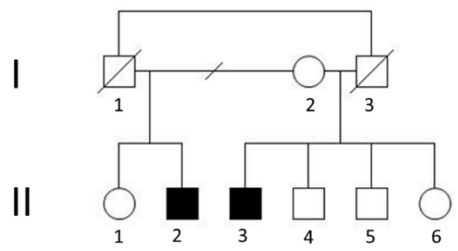
CANVAS3



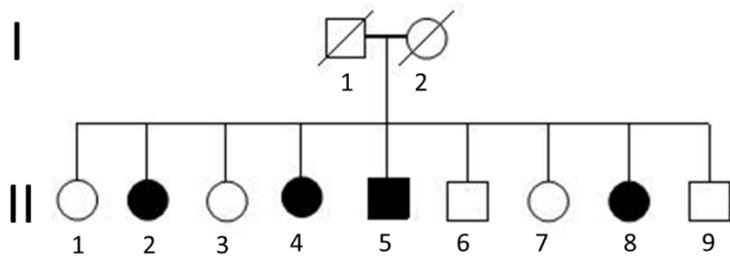
CANVAS4



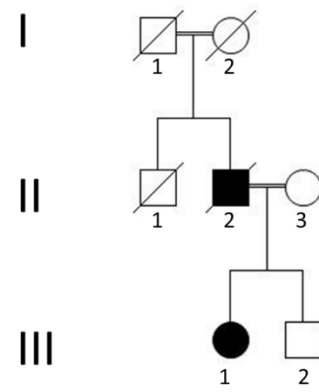
CANVAS17



CANVAS9



CANVAS21

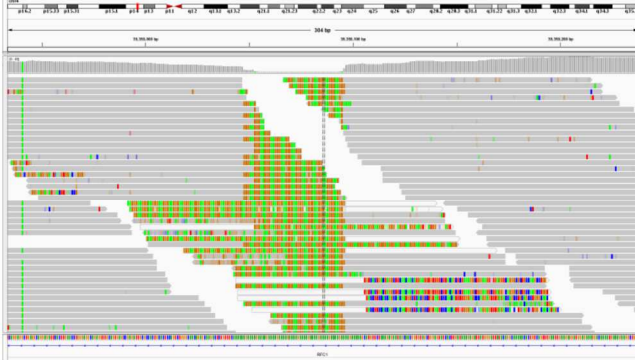


**Figure S1: Pedigree structure of CANVAS families utilized for linkage studies (CANVAS1, 2, 3, 4 and 9) or demonstrating multigenerational inheritance.**

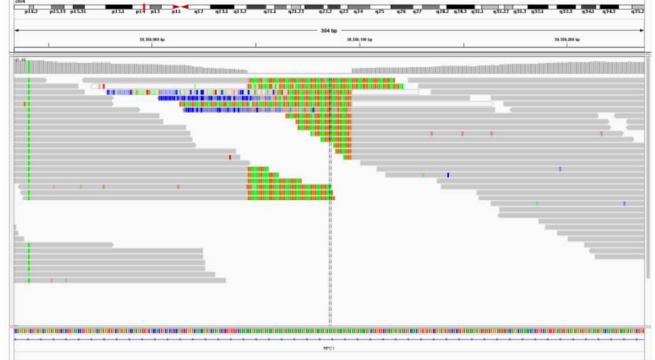
PCR-based WGS,  
~60x coverage

PCR-free WGS,  
~30x coverage

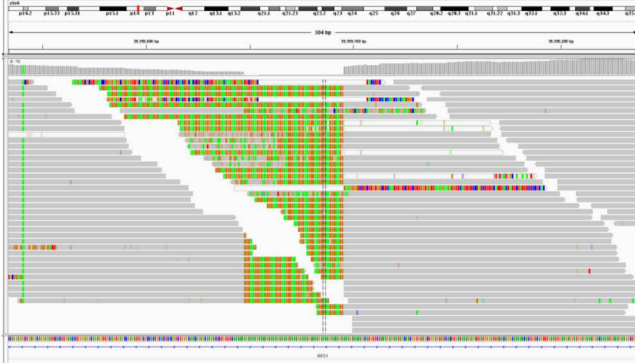
CANVAS1



CANVAS2



CANVAS9



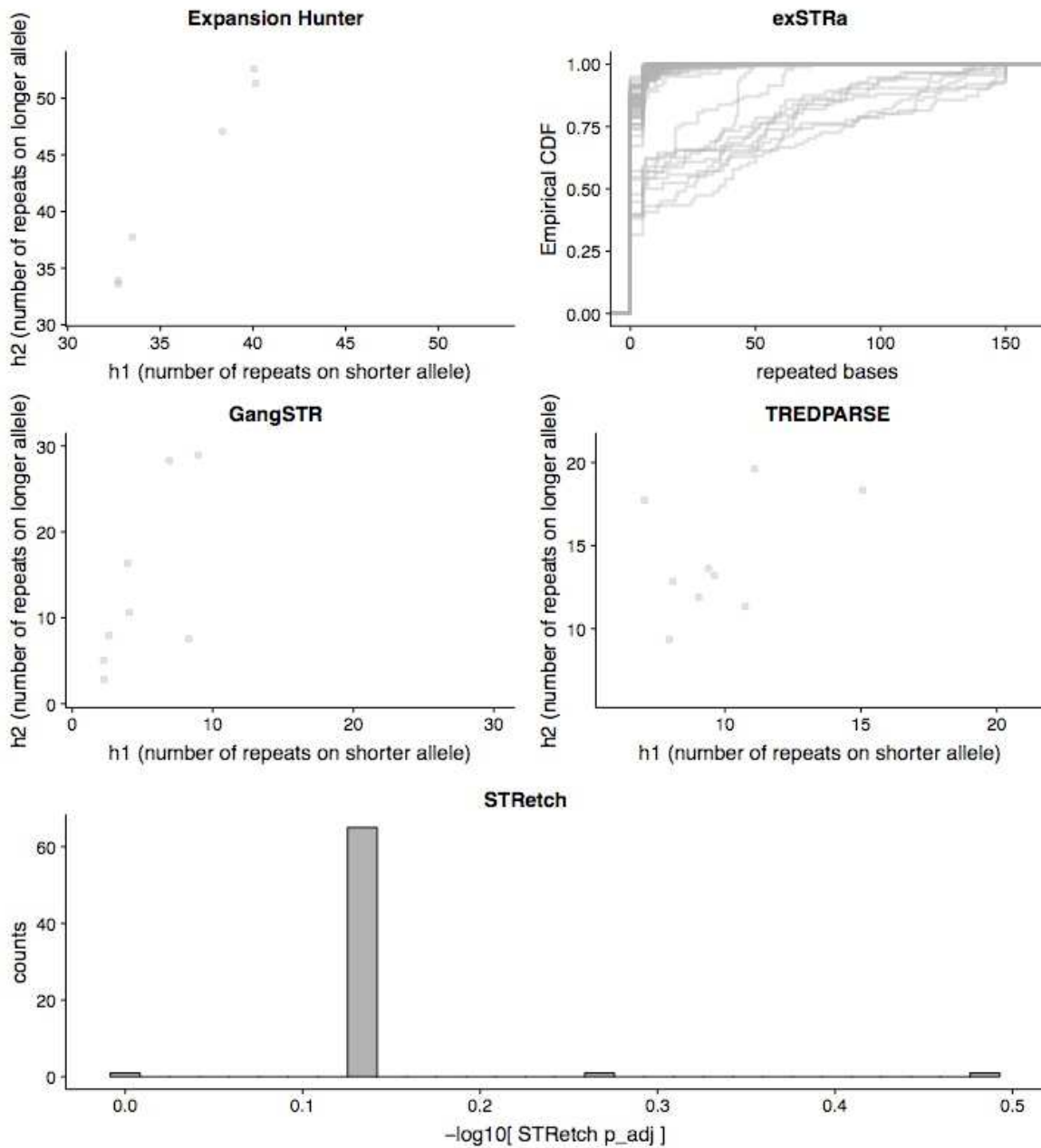
CANVAS8



**Figure S2: IGV snapshots of the (AAGGG)<sub>n</sub> locus in *RFC1*.**

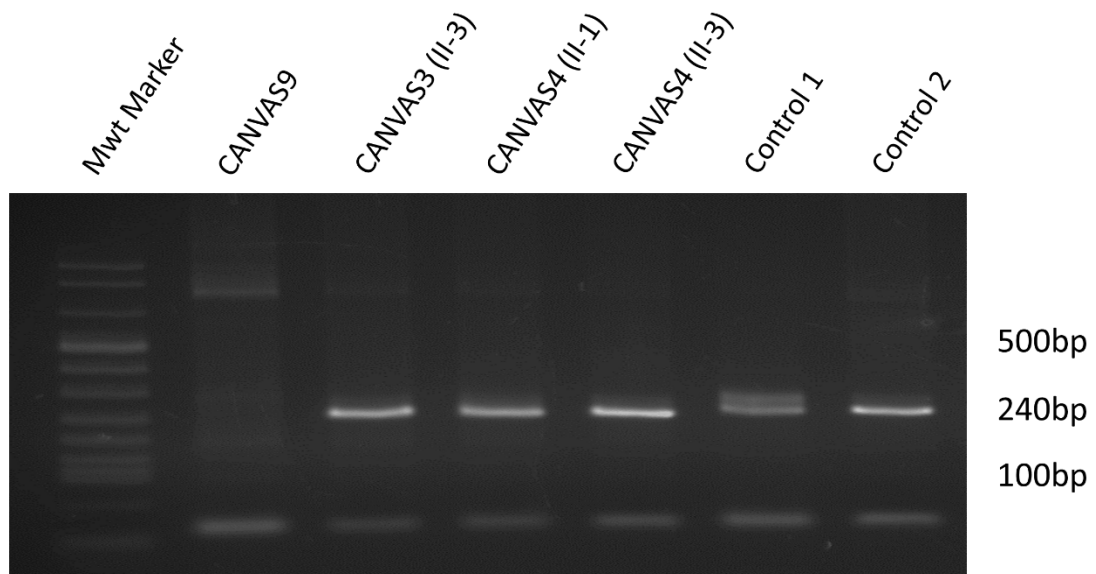
Illustrated are affected individuals from homozygous carriers of AAGGG (CANVAS1, 2, 9) and an individual who carries the AAGGG allele in addition to an expanded AAAGG allele (CANVAS8).



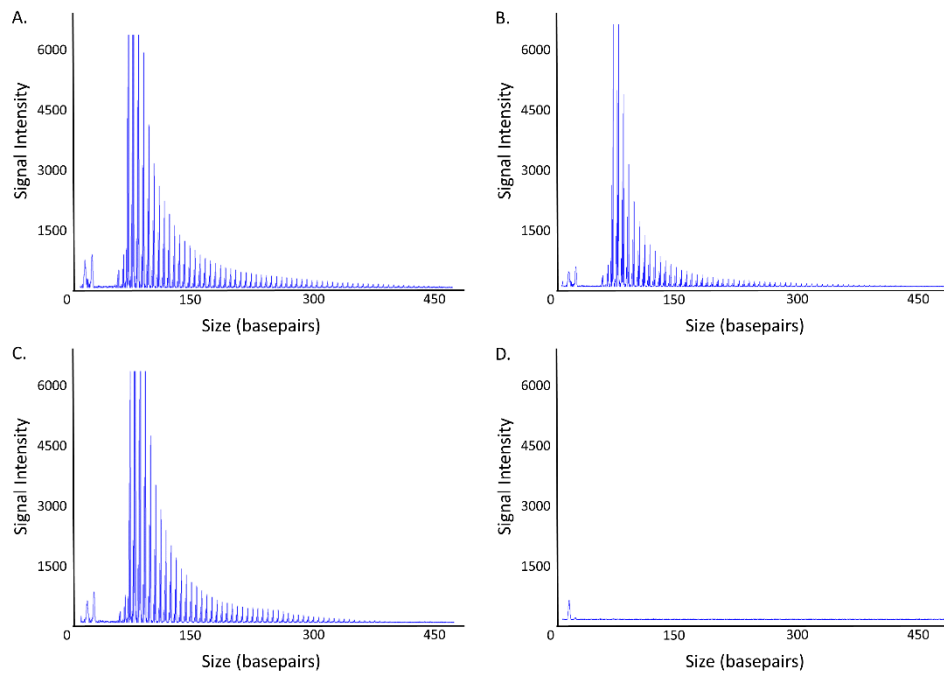


**Figure S4: Computational validation of the (AAGGG)<sub>n</sub> STR.**

The WGS from 69 unrelated non-CANVAS individuals (Coriell dataset) were analysed at the coordinate's chr4:39350045-39350095, using the tools exSTRa, EH, GangSTR, TREDPARSE and STRetch.

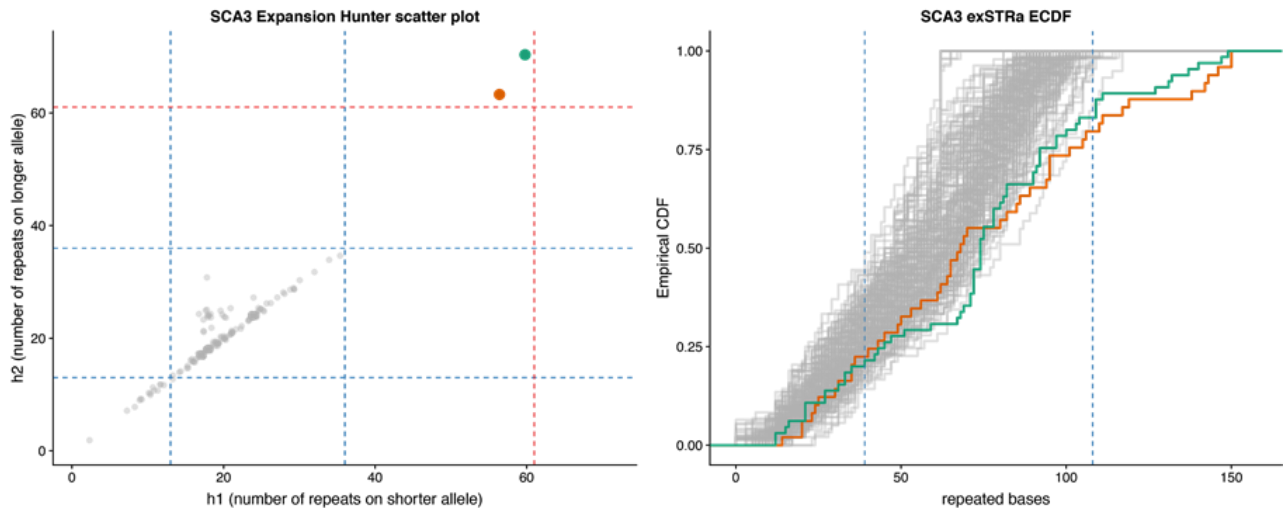


**Figure S5: PCR analysis of unaffected individuals from *RFC1* linked CANVAS families.** Genomic DNA from three unaffected family members from CANVAS3 and CANVAS4 was analysed for the presence of a non-expanded *RFC1* allele. All three individuals carried at least one non-expanded allele, as indicated by the ~250bp PCR product. CANVAS9 is homozygous for the expanded allele and indicates the pattern expected in the absence of the reference allele. Pedigree structure and affected status are illustrated in Figure S1.



**Figure S6: Repeat-primed PCR analysis of unaffected individuals with (AAGGG)<sup>exp</sup> RE.**

Representative images of the repeat-primed PCR for the (AAGGG)<sup>exp</sup> RE demonstrating a saw-toothed product with 5 base pair repeat unit size, amplified from gDNA of control individuals with heterozygous (A, B) or homozygous (C) alleles encoding the (AAGGG)<sup>exp</sup> RE. No product was observed for the no gDNA template negative control (D).



**Figure S7: Identification of a pathogenic SCA3 expansion in CANVAS13.**

The WGS data for SCA3 (green dot) and an individual with a confirmed SCA3 RE (orange dot) was analysed using ExpansionHunter and exSTRa. This analysis identified a heterozygous, pathogenic SCA3 expansion in CANVAS13.



Primer name	Position (hg19)	Sequence
CANVAS_RFC1_3F	chr4: 39350172-39350192	5'-ACTGACAGTGTTTTTGCCTGT
CANVAS_RFC1_3R	chr4:39349940-39349959	5'-GGCTGAGGCAGGAGATTCAC
TPP_CANVAS_FAM_2F	chr4:39350172-39350192	FAM-5'-ACTGACAGTGTTTTTGCCTGT
5R_TPP_M13R_CANVAS_RE_R	NA	5'-CAGGAAACAGCTATGACC_AAGGGAAGGGAAGGGAAGGGAAGGG
TPP_M13R	NA	5'-CAGGAAACAGCTATGACC

**Table S1: Primer sequences for analysis of RFC1.**

The gene reference sequences utilized were NC\_000004.11 and NM\_002913.4 (*RFC1*).

ENSEMBL Gene ID	Gene name	Gene biotype	OMIM gene ID	OMIM disease phenotype (phenotype, OMIM phenotype ID, inheritance pattern)
ENSG00000121895	TMEM156	protein_coding	NA	NA
ENSG00000109790	KLHL5	protein_coding	608064	NA
ENSG00000249207	RP11-360F5.1	antisense	NA	NA
ENSG00000249685	RP11-360F5.3	lincRNA	NA	NA
ENSG00000157796	WDR19	protein_coding	608151	?Cranioectodermal dysplasia 4, 614378, AR; ?Short-rib thoracic dysplasia 5 with or without polydactyly, 614376, AR;Nephronophthisis 13, 614377, AR; Senior-Loken syndrome 8, 616307, AR
ENSG00000035928	RFC1	protein_coding	102579	NA
ENSG00000206675	RNU6-32P	snRNA	NA	NA
ENSG00000222592	RNU6-887P	snRNA	NA	NA
ENSG00000134962	KLB	protein_coding	611135	NA
ENSG00000264621	MIR5591	miRNA	NA	NA
ENSG00000238797	Y_RNA	misc_RNA	NA	NA
ENSG00000163682	RPL9	protein_coding	603686	NA
ENSG00000121897	LIAS	protein_coding	607031	Hyperglycinemia, lactic acidosis, and seizures, 614462, AR
ENSG00000224097	RP11-472B18.1	pseudogene	NA	NA
ENSG00000109814	UGDH	protein_coding	603370	NA
ENSG00000249348	UGDH-AS1	antisense	NA	NA
ENSG00000163683	SMIM14	protein_coding	NA	NA
ENSG00000252796	RNU7-11P	snRNA	NA	NA
ENSG00000255458	RP11-539G18.2	lincRNA	NA	NA
ENSG00000078140	UBE2K	protein_coding	602846	NA
ENSG00000252975	Y_RNA	misc_RNA	NA	NA
ENSG00000249019	RP11-539G18.1	pseudogene	NA	NA
ENSG00000243260	RN7SL558P	misc_RNA	NA	NA
ENSG00000180610	ZBTB12P1	pseudogene	NA	NA

ENSG00000121892	PDS5A	protein_coding	613200	NA
ENSG00000271278	TCEB1P33	pseudogene	NA	NA
ENSG00000252970	RNA5SP159	rRNA	NA	NA
ENSG00000250568	RP11-333E13.2	pseudogene	NA	NA
ENSG00000231707	PABPC1P1	pseudogene	NA	NA
ENSG00000249064	KRT18P25	pseudogene	NA	NA
ENSG00000205794	RP11-333E13.4	pseudogene	NA	NA
ENSG00000078177	N4BP2	protein_coding	NA	NA
ENSG00000200455	RNU6-1112P	snRNA	NA	NA
ENSG00000201863	SNORA51	snoRNA	NA	NA
ENSG00000248977	RP11-395I6.1	pseudogene	NA	NA
ENSG00000260296	RP11-395I6.3	sense_overlapping	NA	NA
ENSG00000168421	RHOH	protein_coding	602037	{?Epidermodysplasia verruciformis, susceptibility to, 4}, 618307, AR
ENSG00000250338	RP11-395I6.2	lincRNA	NA	NA
ENSG00000249241	AC195454.1	lincRNA	NA	NA
ENSG00000174343	CHRNA9	protein_coding	605116	NA
ENSG00000239010	RNU7-74P	snRNA	NA	NA
ENSG00000250893	RP11-588L15.2	antisense	NA	NA

**Table S2. Genes and known disease associations within the CANVAS linkage region**

region	repeat.	pval	adj	Func.refGene	Gene.refGene
chr4:39350095-39350508	AAGGG	0.0027	1	intronic	RFC1
chr5:68499727-68500146	AATATATATATATAG	0.0027	1	intronic	CENPH
chr7:157344383-157344946	ACAGCCACCACCCACCCC	0.0027	1	intronic	PTPRN2
chr9:138739792-138740277	AAGGGGAGGGGAGTGGGGGG	0.0027	1	intronic	CAMSAP1
chr11:100495523-100496852	AATATGTGTATATATGT	0.0027	1	intergenic	CNTN5; LOC100128386
chr11:127149089-127150146	AAG	0.0027	1	ncRNA_intronic	LOC101929497
chr19:43901669-43901885	AAATATATATTATATAT	0.0027	1	intronic	TEX101
chr20:43066951-43067388	AAAAATATAATATAT	0.0027	1	intergenic	HNF4A; LINC01430
chr3:175416044-175416658	ACACATACATATAT	0.0028	1	intronic	NAALADL2
chr5:4286195-4286524	AAGCTATATATATATAGTG	0.0028	1	intergenic	IRX1; LINC02114
chr5:74613999-74614513	AAAG	0.0028	1	intergenic	ANKRD31; HMGCR
chr6:165352545-165353525	AAAG	0.0028	1	intergenic	MEAT6; C6orf118
chr7:4930864-4931703	AGAT	0.0028	1	intergenic	RADIL; MMD2
chr7:157846063-157847123	ACCCAGAGACGCAGAG	0.0028	1	intronic	PTPRN2
chr10:134789654-134790481	AATACATTCCACGTGTATC	0.0028	1	ncRNA_exonic	LINC01168
chr11:2182294-2182840	ACACCCCTGTCCCC	0.0028	1	UTR5	INS; INS-IGF2
chr11:120746805-120747444	ACC	0.0028	1	ncRNA_intronic	LOC101929227
chr11:134177695-134179032	ACC	0.0028	1	intronic	GLB1L3
chr13:40788773-40788913	AATATAT	0.0028	1	ncRNA_intronic	LINC00548

**Table S3: Expansion Hunter de novo results - screening for RE in two CANVAS individuals compared to 31 controls (all PCR based WGS,  $p < 0.01$ ) ranked by p-value and ordered by chromosome**

