

1 **Detecting tandem repeat expansions in cohorts sequenced with short-read**  
2 **sequencing data**

3

4 Rick M Tankard<sup>1,2,3</sup>, Mark F Bennett<sup>1,2</sup>, Peter Degorski,<sup>1,2</sup> Martin B Delatycki,<sup>4,5,6</sup>  
5 Paul J Lockhart,<sup>4,6†</sup> Melanie Bahlo<sup>1,2\* †</sup>

6

7 <sup>1</sup> Population Health and Immunity Division, The Walter and Eliza Hall Institute of  
8 Medical Research, Parkville 3052, VIC, Australia

9 <sup>2</sup> Department of Medical Biology, The University of Melbourne, Melbourne 3010,  
10 VIC, Australia

11 <sup>3</sup> Mathematics and Statistics, Murdoch University, Murdoch 6150, WA, Australia

12 <sup>4</sup> Bruce Lefroy Centre for Genetic Health Research, Murdoch Children's Research  
13 Institute, Royal Children's Hospital, Parkville, Victoria, Australia.

14 <sup>5</sup> Victorian Clinical Genetics Services, Parkville, Victoria, Australia.

15 <sup>6</sup> Department of Paediatrics, University of Melbourne, Parkville, Victoria, Australia.

16

17

18 † Joint senior authors

19 \* Corresponding Author: Melanie Bahlo, [bahlo@wehi.edu.au](mailto:bahlo@wehi.edu.au), twitter handle  
20 @MelanieBahlo

21

22 **Keywords:** next-generation sequencing; repeat expansion disorders; short tandem  
23 repeats; whole exome sequencing; whole genome sequencing

24

25 **Abstract**

26 Repeat expansions cause over 30, predominantly neurogenetic, inherited  
27 disorders. These can present with overlapping clinical phenotypes, making molecular  
28 diagnosis challenging. Single gene or small panel PCR-based methods are employed  
29 to identify the precise genetic cause, but can be slow and costly, and often yield no  
30 result. Genomic analysis via whole exome and whole genome sequencing (WES and  
31 WGS) is being increasingly performed to diagnose genetic disorders. However, until  
32 recently analysis protocols could not identify repeat expansions in these datasets.

33 A new method, called exSTRa (**expanded Short Tandem Repeat algorithm**)  
34 for the identification of repeat expansions using either WES or WGS was developed  
35 and performance of exSTRa was assessed in a simulation study. In addition, four  
36 retrospective cohorts of individuals with eleven different known repeat expansion  
37 disorders were analysed with the new method. Results were assessed by comparing to  
38 known disease status. Performance was also compared to three other analysis methods  
39 (ExpansionHunter, STRetch and TREDPARSE), which were developed specifically  
40 for WGS data. Expansions in the STR loci assessed were successfully identified in  
41 WES and WGS datasets by all four methods, with high specificity and sensitivity,  
42 excepting the FRAXA STR where expansions were unlikely to be detected. Overall  
43 exSTRa demonstrated more robust/superior performance for WES data in comparison  
44 to the other three methods. exSTRa can be applied to existing WES or WGS data to  
45 identify likely repeat expansions and can be used to investigate any STR of interest,  
46 by specifying location and repeat motif. We demonstrate that methods such as  
47 exSTRa can be effectively utilized as a screening tool to interrogate WES data  
48 generated with PCR-based library preparations and WGS data generated using either  
49 PCR-based or PCR-free library protocols, for repeat expansions which can then be

50 followed up with specific diagnostic tests. exSTRa is available via GitHub  
51 (<https://github.com/bahlolab/exSTRa>).

52

### 53 **Introduction**

54 Thousands of short tandem repeats (STRs), also called microsatellites, are scattered  
55 throughout the human genome. STRs vary in size but are commonly defined as  
56 having a repeat motif 2-6 base pairs (bps) in size. They are underrepresented in the  
57 coding regions of the human genome<sup>1</sup>, despite the vast majority being population  
58 polymorphisms with no, or very little, phenotypic consequence. STRs were used as  
59 genetic markers for linkage mapping for human studies for many years, and continue  
60 to be used, but primarily for non-human studies. A subset of STRs can however cause  
61 disease. Pathogenic STRs have either one or two alleles, depending on the genetic  
62 model, that exceed some threshold for biological tolerance. These diseases are known  
63 as repeat expansion disorders. The abnormal STR allele(s), may affect gene  
64 expression levels, cause premature truncation of the protein or result in aberrant  
65 protein folding.<sup>2</sup> Repeat expansions at different STR loci share biological  
66 consequences. Common disease mechanisms mediated by repeat expansion disorders  
67 include Repeat-associated non-AUG translation and MBNL spliceosome interference,  
68 for example caused by CUG expansions in Myotonic Dystrophy Type 1 (DM1).  
69 These mechanisms are reviewed in Hannan<sup>3</sup>.

70

71 Repeat expansions cause ~30 inherited germline human disorders, predominantly  
72 neurogenetic diseases most often presenting with ataxia as a clinical feature. The size  
73 of pathogenic allele varies from ~60 repeats observed in the gene encoding the  
74 Calcium Voltage-Gated Channel Subunit Alpha1 A (*CACNA1A*) to several thousand

75 repeats observed in the gene encoding the Calcium Voltage-Gated Channel Subunit  
76 Alpha1 A (*C9orf72*) (Table 1). Remarkably 12 repeat expansions have now been  
77 identified as causing dominant forms of spinocerebellar ataxias. Other disorders  
78 caused by repeat expansions include fragile X syndrome (OMIM #300624, a repeat in  
79 the 5'UTR of *FMRI*), Huntington Disease (OMIM #606438, a repeat in exon 1 of  
80 *HTT*), myotonic dystrophy (OMIM #602668, repeats in *DMPK* and *ZNF9*), fronto-  
81 temporal dementia and amyotrophic lateral sclerosis 1 (OMIM #105550, a 6-mer  
82 repeat in *C9orf72*) and Unverricht-Lundborg disease, a severe myoclonic epilepsy  
83 (OMIM #254800, in *CSTB*). The genetic mode of inheritance encompasses autosomal  
84 dominant (e.g. SCA1, OMIM #164400) and recessive (e.g. Friedreich ataxia, OMIM  
85 #229300), as well as X-linked recessive (e.g. fragile X syndrome, OMIM #300624).  
86 Novel pathogenic alleles underlying repeat expansion disorders continue to be  
87 discovered, with the two most recently described STRs being pentamer repeats<sup>4,5</sup>. A  
88 selected list of repeat expansion disorders are shown in Table 1.

89

90 Many repeat expansion disorders show anticipation; a phenomenon whereby younger  
91 generations are affected by earlier age of onset. Anticipation is usually caused by an  
92 increase in repeat size between generations. When anticipation is observed it indicates  
93 that a search for repeat expansions as the cause of disease is warranted.

94

95 Friedreich ataxia is the most common of the recessive repeat expansion disorders,  
96 with a disease prevalence of 3 to 4/100,000 but with a carrier frequency of 1/100.<sup>6</sup>  
97 Fragile X syndrome is the most common cause of inherited intellectual disability and  
98 affects ~1/5000 individuals.<sup>7; 8</sup> Hence these diseases as a whole contribute  
99 significantly to the overall Mendelian disease burden in human populations.

100

101 Diagnostic identification of repeat expansions can be time consuming and costly.  
102 Current medical diagnosis consists of precise PCR or Southern blot assay, which  
103 require diagnostic laboratories that have refined these assays for each different repeat  
104 expansion. The clinician has to determine which repeat expansions are most likely to  
105 be relevant and submit the patient's DNA to appropriate laboratories. This can be  
106 difficult, given the phenotypic overlap between the different STRs, the potential  
107 heterogeneity in the symptoms and the variation in penetrance and age of onset,  
108 which is also dependent on the size of the allele and effect of modifier genes.<sup>9;10</sup> In  
109 addition, up to 50% of individuals with a diagnosis of ataxia may be due to other  
110 mutation types, such as single nucleotide variants (SNVs) and short  
111 insertion/deletions (indels).<sup>11</sup> Therefore, molecular diagnosis of these disorders often  
112 also requires conventional sequencing of candidate genes, either by Sanger, targeted  
113 panel or Next Generation Sequencing (NGS) methods.

114

115 Short-read NGS data, such as that generated by the Illumina sequencing platform, is  
116 currently predominant in both research and clinical diagnostic applications. Moreover,  
117 Whole Genome Sequencing (WGS) is now an affordable technology, gradually  
118 replacing whole exome sequencing (WES) for clinical genomics. Illumina's HiSeq X  
119 and NovaSeq platforms are currently the most commonly used platform for the  
120 generation of human WGS data and in particular clinical human genome sequencing  
121 with low error rates and well-documented, consistent, performance.

122

123 Illumina HiSeq X data reads are 150 bp in length and are designed so that the reads  
124 are transcribed facing each other, where the template DNA predominantly has a small

125 gap between the reads that is not sequenced. This gap can vary in size, but standard  
126 library preparation methodologies generate insert fragment lengths of ~350 bps,  
127 resulting in a gap of ~50bp.

128

129 Standard clinical diagnostic pipelines focus on the identification of SNVs and indels.  
130 Bioinformatic tools have been developed to genotype STRs, but are almost entirely  
131 confined to those STR alleles that are spanned by reads.<sup>12-16</sup> Pathogenic repeat  
132 expansions are usually significantly longer than the reads generated by short-read  
133 sequencing platforms such as Illumina, and may be longer than the library insert  
134 fragments lengths. Therefore, the short reads cannot span many pathogenic repeat  
135 expansion alleles, such as those that cause SCA2 (OMIM #183090), or SCA7 (OMIM  
136 #164500, Table 1). Furthermore some of these reads are not mapped, or poorly  
137 mapped, to the STR allele, due to sequencing bias and alignment issues such as: (i)  
138 the repetitive nature of the repeat itself where the expanded alleles require alignments  
139 of additional repetitive bases, (ii) multiple occurrences of the same repeat throughout  
140 the genome, leading to multi-mapping reads, and (iii) GC bias. Despite this, these data  
141 do still carry information about the expanded allele with a larger number of reads  
142 mapping to the STR for an expanded allele than expected, based on the reference STR  
143 allele lengths.

144

145 Several methods now describe the detection of repeat expansion in short read NGS  
146 data. These include ExpansionHunter<sup>17</sup>, STRetch<sup>18</sup> and TREDPARSE<sup>19</sup>, reviewed in  
147 Bahlo et al<sup>20</sup>. These methods are focused on detection of repeat expansions in whole  
148 genome sequencing data, with a preference for PCR-free library free protocols.  
149 ExpansionHunter and TREDPARSE determine whether an individual has an

150 expansion based on pre-determined thresholds, however TREDPARSE also has a  
151 likelihood ratio test with a likelihood framework that determines the genetic model  
152 and the likelihood of expansion. STRetch uses a genome reference augmented with  
153 decoy chromosomes, consisting of long stretches of all 1 to 6 bp repeat expansions to  
154 competitively attract long repeats. None of these methods have been assessed for  
155 performance in comparison to each other or in WES data.

156

157 Here we describe the development of the STR repeat expansion-calling algorithm,  
158 exSTRa (**expanded STR algorithm**), which detects expanded repeat expansion  
159 allele(s) at repeat expansion loci, specified by the user, in cohorts of sequenced  
160 individuals. We demonstrate the utility of the method with twelve different verified  
161 repeat expansion disorders. exSTRa is designed to be applied to cohorts of individuals  
162 without requiring a set of controls. This is because exSTRa is designed as an outlier  
163 detection test, where the majority of individuals (>85%) are assumed to have normal  
164 length alleles at a particular repeat expansion locus. This assumption is robust for the  
165 majority of disease cohorts, even spinocerebellar ataxias. exSTRa also generates  
166 unique empirical cumulative distribution function (ECDF) plots of individual's repeat  
167 motif distributions, plotted for all individuals in a cohort, facilitating QC for batch  
168 effects and validity of assumptions. We demonstrate for the first time, that repeat  
169 expansion detection is possible with WES data and further demonstrate on additional  
170 STR loci, that PCR-based library preparation WGS, whilst inferior to PCR-free  
171 library preparation WGS data, can be used to confidently interrogate most known  
172 STR loci. This will enable researchers to interrogate the thousands of existing NGS  
173 datasets for repeat expansions at known repeat loci or any other loci they wish to  
174 investigate.

175

## 176 **Methods**

### 177 **Study cohorts and next-generation sequencing data generation**

178 Individuals with already diagnosed repeat expansion disorders were recruited for this  
179 study. The repeat expansion status was verified via standard diagnostic STR-specific  
180 PCR-based assays. Individuals affected by neurogenetics disorders not due to known  
181 repeat expansions were recruited as controls. These individuals were not tested for  
182 any of the known repeat expansion loci with standard methods as none of them are  
183 affected by symptoms that are typical of expansion disorders such as ataxia. All  
184 individuals were recruited at the Murdoch Children's Research Institute, and provided  
185 written informed consent (Human Research Ethics Committee #28097, #25043 and  
186 #22073).

187

188 Four cohorts underwent different types of NGS, with some individuals being  
189 sequenced multiple times. Individuals were sequenced with either: (i) WES with the  
190 Agilent V5+UTR capture platform (4 repeat expansion patients, with 4 different  
191 expansion disorders, 58 controls), (ii) WGS with the TruSeq Nano protocol, which  
192 includes a PCR step to increase sequencing material (17 repeat expansion patients,  
193 with 8 different expansion disorders, 16 controls), or (iii) WGS with the PCR-free  
194 cohort consisting 118 individuals (52 females and 66 males). Samples in this cohort  
195 were either affected with the repeat expansion disorder, or carriers, for one of:  
196 FRAXA (15 expanded, 19 intermediate), FRDA (25), DM1 (17), HD (13), SCA1 (3),  
197 DRPLA (2), SBMA (1) and SCA3 (1), or relatives with no known expansion (22),  
198 with all samples sourced from the Coriell resource. The WES cohort is designated as  
199 WES\_PCR. Two different cohorts were sequenced with protocol (ii). These are



200 designated as WGS\_PCR\_1 and WGS\_PCR\_2. The WGS cohort was designated as  
201 WGS\_PF. These cohorts are described in Table 2.

202

### 203 **Sequencing Data generation**

204 WGS data with PCR (WGS\_PCR1 and WGS\_PCR2) was generated by the Kinghorn  
205 Centre for Clinical Genomics, Garvan Institute of Medical Research, Sydney,  
206 Australia with HiSeq X Ten. The WES data (WES\_PCR) was generated by the  
207 Australian Genome Research Facility, Melbourne, Australia, and sequenced on a  
208 HiSeq 2500 sequencer. All WGS\_PF samples were sequenced on the Illumina HiSeq  
209 X sequencing platform at Illumina, La Jolla, California, USA. Further details can be  
210 found in Dolzhenko et al. All sequencing data was aligned to the hg19 human genome  
211 reference using the Bowtie 2 aligner<sup>21</sup> in local alignment mode.

212

### 213 **Definition of Repeat Expansion Loci**

214 Table 1 defines the chromosomal location, physical map location, disease, genetic  
215 disease model and repeat motif, normal and repeat expansion size for 24 repeat  
216 expansion loci, which cause neurological disorders. For the analyses in this paper we  
217 examined 21 of these STR loci, excluding the more recently discovered SCA37 and  
218 FAME1 loci, and the SCA31 locus, where the inserted repeat is not in the reference  
219 sequence. This focused the analysis on currently most likely tested expansion loci and  
220 in particular concentrating on the spinocerebellar ataxia repeat expansion loci.

221

### 222 **Data extraction for repeat expansions**

223 We developed a two-step analysis method, called exSTRa, detailed in the  
224 Supplemental data, to identify individuals likely to have a repeat expansion at a

225 particular STR locus. The analysis method extracts STR repeat content information  
226 for each read, stemming from a particular individual, which has been identified as  
227 mapping to one of the 21 STR loci. We designed a statistical test that captures the  
228 differences between an individual to be tested within a cohort of cases and controls.  
229 All N individuals within a cohort are examined in turn at each of the 21 known  
230 pathogenic repeat expansion loci by comparing each individual in turn to all N other  
231 individuals in each cohort. This generates 21xN test statistics per cohort. The  
232 empirical p-value of the test statistic was determined using a simulation method. All  
233 p-values over all STR loci for all individuals within each cohort were assessed for  
234 approximate uniform distribution with histograms and Quantile-Quantile (Q-Q) plots.  
235  
236 Raw data was visualized using empirical cumulative distribution functions (ECDFs),  
237 which display the distribution of the amount of STR repeat motif found in each read,  
238 ordered from smallest to largest content amount, as a step function. This allows  
239 comparison of the distributions, regardless of sequencing depth. Reads generated  
240 from expanded alleles have increased numbers of repeat motifs in their reads  
241 compared to reads stemming from normal alleles. This produces a shift of the read  
242 repeat motif distribution to the right for the individual with the repeat expansion, in  
243 comparison to reads from individuals with normal alleles.

244

### 245 **Simulation Study**

246 We conducted a simulation study using the next generation sequencing data  
247 simulation package ART<sup>22</sup>, which simulates NGS data with realistic error profiles  
248 based on supplied reference genomes. Alleles at STR loci were simulated using  
249 reference genomes where alleles (normal, intermediate, expanded) had been inserted

250 into the reference genome. STR loci such as HD do not have an intermediate range or  
251 only a very narrow range. We extensively searched the literature to determine  
252 pathogenic and non-pathogenic ranges of STR length alleles. We only used the  
253 ‘overall’ distribution, ignoring any ethnic specificity for these loci. We did not apply a  
254 stutter model in the simulations, as this was not feasible due to ARTs constraints. We  
255 simulated data for 20 STR loci (excluding FAME1, SCA31, SCA37 and SCA31), for  
256 200 controls, and ten normal, ten intermediate range and ten expanded individuals.  
257 These 30 individuals were tested for expansions. The STR genotype for the controls  
258 was randomly chosen based on the distributions of these as described in the literature  
259 (Supplemental Table S3). Ten normal, intermediate and expanded alleles were chosen  
260 based on uniform distances between alleles, covering the known normal, intermediate  
261 and expanded allele ranges as described in the literature (Supplementary Table S4);  
262 for autosomal dominant loci, the second allele was chosen randomly with the same  
263 method as the controls. For the recessive STR loci EPM1 and FRDA we sampled two  
264 expanded alleles for individuals with disease. To allow for STR loci assessment on  
265 the X chromosome (FRAXA, FRAXE and SBMA) we generated half of the samples  
266 as male and the other half as female, with males having a single X chromosome and  
267 hence a single STR allele. For the X chromosome STR loci we only investigated the  
268 male individuals. To investigate the effect of control sample size on detection with  
269 exSTRa we sub-sampled the control cohort at intervals of 50, with control cohort  
270 sizes ranging from 50 to 200 individuals. The ART command used to generate the  
271 simulated data was:

```
272 art_illumina -i ${file} -p -na -l 150 -f 50 -m 450 -s 50 -o $outfile/$base -l  
273 ${profiles}/HiSeqXPCRfreeL150R1.txt -2 ${profiles}/HiSeqXPCRfreeL150R2.txt
```

274

275 **Performance evaluation**

276 For exSTRa we called individuals as being normal or expanded based on the  
277 Bonferroni multiple testing corrected p-values derived from our empirical p-values.  
278 The number of Bonferroni corrections for the four cohorts was performed based on  
279 the 21 STRs tested per individual for the WGS cohorts and 13 for the WES cohort.  
280 Repeat expansion calls were compared to the known disease status. Performance of  
281 all four methods was evaluated by examining the number of true positives (TP), true  
282 negatives (TN), false positives (FP), false negatives (FN), sensitivity, which is defined  
283 as  $TP/(TP+FN)$  and specificity, which is defined as  $TN/(TN+FP)$  at each STR and  
284 then summarized across the STR loci, within cohorts.

285

286 **Comparison with ExpansionHunter, STRetch and TREDPARSE**

287 ExpansionHunter<sup>17</sup> estimates the repeat size using a parametric model but does not  
288 attempt to call repeat expansions in a probabilistic framework. ExpansionHunter was  
289 used to determine if alleles were larger than currently known smallest disease-causing  
290 repeat expansion alleles. STRetch was used to detect the presence of repeat expansion  
291 using its statistical test, which is also an outlier detection test. Bonferroni corrections  
292 were calculated as per the exSTRa analysis. TREDPARSE was used to both estimate  
293 the repeat size and to detect the presence of an expansion based on its likelihood  
294 model. Bonferroni corrections were applied in the same way as for exSTRa.

295

296 **Results**

297 **Simulation Study Results**

298 The simulation study of the 20 STR loci provide evidence of the validity and  
299 robustness of the exSTRa test statistic with respect to control cohort size, repeat

300 expansion size and known expansion status. Decreasing the control cohort in exSTRa  
301 showed that results were robust as the control sample size decreased (Supplementary  
302 Figure S5). exSTRa also showed consistent results when the size of the repeat  
303 expansion allele varied, with longer expansion alleles achieving smaller p-values  
304 (Supplementary Figure S4). Overall exSTRa p-values showed adequate Type 1 error,  
305 and good discriminatory ability between expansion and non-expansion individuals  
306 (Supplementary Figures S3 and S4). The ECDF plots, which are unique to exSTRa,  
307 show the effect of increasing expansion size in all STRs, with commensurate right  
308 shifts of the distributions. The ECDFs also allow heuristic determination of the  
309 genetic model, with larger shifts to the right for the recessive FRDA STR and the X-  
310 linked STRs (FRAXA, FRAXE and SBMA). Dominant loci only show the shift in  
311 ECDF for the upper half of the ECDF (Supplementary Figure 2). All STR loci  
312 performed well for repeat expansion detection in the simulation studies, including  
313 FRAXA and FRAXE. The simulated dataset is available to other researchers on  
314 request.

315

### 316 **Coverage and Alignment Results for study cohorts**

317 Full coverage and alignment results are in Supplemental Table S2 for three cohorts,  
318 but not WGS\_PF\_3, which is described in Dolzhenko et al<sup>17</sup>. The median coverage  
319 achieved was 44, 66, 82 and 46.3 for cohorts WES, WGS\_PCR\_1, WGS\_PCR\_2 and  
320 WGS\_PF\_3 respectively, with 1<sup>st</sup> and 3<sup>rd</sup> quartile coverage of (37,48.25), (49.5,71),  
321 (76.5,84) and (44.9,47.9) Genome-wide sample specific coverage variability, as  
322 measured by the median IQR of the mean coverage library size corrected samples,  
323 was very similar between all three WGS\_cohorts (WGS\_PCR\_1 median IQR = 8,

324 WGS\_PCR\_2 median IQR = 5.7, WGS\_PF\_3 median IQR = 8.3). In contrast the  
325 WES data showed substantial variability (median IQR = 22.3).

326

### 327 **STR loci sequencing coverage ability**

328 We examined the 21 STR loci for coverage in our four study cohorts. As expected  
329 WES\_PCR only achieved reasonable coverage for repeat expansion detection in a  
330 subset of the STR loci. However, this included many of the known repeat expansion  
331 STRs located in coding regions (8 out of 10) (Figure 1 and Supplemental Figure S1).  
332 SCA6 (OMIM #183086, *CACNA1A*) and SCA7 (*ATXN7*) are poorly covered. Despite  
333 the use of the Agilent SureSelect V5+UTR capture platform, which incorporates  
334 UTRs we achieved no, or very low coverage, for the known repeat expansion loci  
335 located in the UTR, such as FRAXA (OMIM #300624), FRAXE (OMIM #309548)  
336 and DM1 (OMIM #160900). DM1 and SCA7 are not captured by the Agilent  
337 enrichment platform (Supplemental Table S3), however both FRAXA and FRAXE  
338 are targeted and therefore should be captured. In general, WGS data outperformed  
339 WES over all STR loci, with one exception, SCA3 (OMIM #109150), located in the  
340 coding region of *ATXN3*. The reason for this is currently unknown.

341

### 342 **Visualizations of repeat motif distributions**

343 ECDF curves of selected loci are shown for each cohort to illustrate the data. Full  
344 results for all 21 loci, for all WGS cohorts, and 10 covered loci for WES cohort, are  
345 given in Supplemental Figures S6-S11. STR loci varied in their coverage with several  
346 loci consistently poorly captured. These were usually loci that are rich in GC content.  
347 Short read NGS data has a known GC bias with a GC content of 40-55% maximizing  
348 sequencing yield, depending on sequencing platform.<sup>23</sup> The shape of the ECDF is

349 affected by additional factors such as: the genetic model (dominant, recessive or X-  
350 linked) and capture efficiency (for WES).

351

352 The STR loci also showed differences in variability with regards to STR motif  
353 lengths. Some STR loci, such as SCA17 (OMIM #607136) and HDL2 (OMIM  
354 #606438), showed little variability in STR allele distributions, regardless of NGS  
355 platform in our cohorts. Identification of outliers is easier for these loci, with low  
356 background variability. Those repeat expansion disorders that are autosomal recessive  
357 or X-linked recessive (in males), also show much clearer outlier distributions (Figure  
358 2, top right panel). This is due to the outlier distribution deviating for either both  
359 alleles, or, in the case of the X-chromosome, and only males, just the one allele  
360 having to be examined (not performed in this analysis).

361

#### 362 **Statistical test results for exSTRa**

363 Test statistics were generated for all 21 loci for all N individuals for all four cohorts  
364 with exSTRa. Combined p-values over all STR loci for all individuals within each  
365 cohort showed approximate uniform distribution with histograms (Supplemental  
366 Figure S12) and Q-Q plots (Supplemental Figure S13), albeit with some inflation of  
367 p-values at both tails. Our study cohorts had very small numbers of control  
368 individuals for some of the cohorts.

369

#### 370 **Expansion call results**

371 Expansion call results are presented in summary form in Tables 3 and 4, and at the  
372 individual level in Supplemental table S4 and S5. For the cohorts WES\_PCR,  
373 WGS\_PCR\_1, WGS\_PCR\_2, WGS\_PCR\_2\_30X\_1, WGS\_PCR\_2\_30X\_2 and

374 WGS\_PF exSTRa achieved sensitivities of 1, 0.67, 0.81, 0.81 and 0.75 and 0.77  
375 respectively, for these cohorts (Table 4), with very high specificity (all cohorts  
376 >0.97). Sensitivity is poorly estimated due to the small number of true positives (TPs)  
377 in some cohorts, which leads to large variability. This is particularly the case for  
378 WES\_PCR (4 cases) and WGS\_PCR\_1 (3 cases). This has also resulted in highly  
379 variable results for the other methods. FRAXA was the STR most refractory to  
380 analysis, performing poorly regardless of sequencing platform and repeat expansion  
381 detection method. Excluding this locus in the evaluation of WGS\_PF increased the  
382 sensitivity from 0.77 to 0.84, but specificity remained unchanged at 0.97.

383

384 We divided the WGS\_PCR\_2 cohort data into two sub-cohorts, where each sample's  
385 data comes from a single flow cell lane that has ~30X coverage. This allowed an  
386 investigation of reproducibility, and assessment at the more standard 30X coverage.  
387 Results were highly reproducible between the two 30X replicates, with only one  
388 sample generating an alternative call between the two sequencing runs. We also  
389 observed very little change in performance between the 60X and 30X data with  
390 virtually identical sensitivity and specificity (Table 4).

391

### 392 **Comparison with other repeat expansion detection methods**

393 Across all cohorts (WES\_1, WGS\_PCR\_1, WGS\_PCR\_2, WGS\_PF) exSTRa called  
394 the most expansions (79 out of 100 known expansions) compared to ExpansionHunter  
395 75 expansions, STRetch 77 expansions, TREDPARSE-L 52 expansions and  
396 TREDPARSE-T 71 expansions, albeit with slightly different results in the REs  
397 identified. Excluding FRAXA exSTRa called 71 out of 82 (87%) expansions,  
398 ExpansionHunter 74 expansions, STRetch 77 expansions, TREDPARSE-L 51



399 expansions and TREDPARSE-T 71 expansions each. Notably, exSTRa was able to  
400 identify expanded repeats at all eleven STR expansions examined. STRetch was  
401 unable to identify the SCA6 expansions in any cohort (N=2 in WGS\_PCR\_2, N=1 in  
402 WGS\_PCR\_1 and N=1 in WES\_1). SCA6 is the shortest of all known repeat  
403 expansions. These shorter expansions fail to map preferentially to the decoy  
404 chromosome for the most part, leading to the inability to call this locus. This will also  
405 apply to other short repeat expansion alleles. However the other methods found most  
406 of the SCA6 expansions, regardless of sequencing platform. All four methods  
407 performed poorly when analyzing samples with an *FMRI* expansion (FRAXA). In the  
408 WGS\_PCR\_1 and 2 cohorts this is due to poor coverage at the FRAXA and FRAXE  
409 loci caused by GC bias issues (Supplemental Figure S1). Although there was a clear  
410 right shift of the exSTRa ECDF plots of both the full mutation and premutation *FMRI*  
411 samples (Figure 3 bottom left panel), this was not always statistically significant. The  
412 other methods similarly performed poorly with this expansion, often failing to detect  
413 it. However, ExpansionHunter and TREDPARSE-T and -R identified pre-mutation  
414 alleles for this locus ~75% of the time. exSTRa identified 5/15 FRAXA expansions,  
415 STRetch identified none and called three of these as SCA3 expansions instead.  
416 STRetch performed equal best with ExpansionHunter in the WGS\_PF cohort but was  
417 the best performer once FRAXA was ignored, finding all remaining repeat  
418 expansions, albeit with the highest false positive rate. TREDPARSE and STRetch  
419 both perform particularly well for large expansions where their use of “in-repeat  
420 reads”<sup>17; 20</sup>, or reads that map entirely to the repeat, is highly advantageous. exSTRa  
421 does not use this information and ExpansionHunter only uses it optionally, for large  
422 repeats. Remarkably all four methods call all 13 HD expansions correctly in the  
423 WGS\_PF\_3 cohort (Supplementary Table S5), suggesting highly robust detection of

424 HD expansions for WGS data. The four methods also unanimously identify the  
425 SBMA expansion and the two DRPLA expansions.

426

427 exSTRa was the equal best performing method for the WGS\_PCR cohorts with  
428 TREDPARSE, and performed best overall for the WES cohort. Overall all methods  
429 performed more poorly in the WES and WGS\_PCR cohorts in comparison to the  
430 WGS\_PF cohort. exSTRa performs well for small repeat expansions and for  
431 platforms where small read fragments have been preferentially selected (WES\_PCR,  
432 WGS\_PCR). Overall the results indicate that no single method is optimal over this  
433 breadth of sequencing library preparations and STR loci. These results suggest that a  
434 consensus call that makes use of all existing methods could be advantageous.  
435 Concordance with at least one other method will be useful to maximize detection of  
436 expansions, especially since specificity is high in all WGS cohorts, across all methods  
437 ( $\geq 0.97$ ). This drops to  $\geq 0.93$  for WES data. Using a rule whereby at least two  
438 expansion calls are required, with at least two calling methods showing concordant  
439 results to calculate a consensus call, leads to sensitivities of 1 for WES\_1, 1 for  
440 WGS\_PCR\_1, 0.81 for WGS\_PCR\_2 (1, if FRAXA is excluded), 0.77 for  
441 WGS\_PF\_3 and 0.94 for WGS\_PF\_3 (excluding FRAXA) (Supplementary Tables 4  
442 and 5, last columns).

443 Computational expense varied between the different repeat expansion tools. Running  
444 time for the WGS\_PF cohort comprising 118 samples using 8 CPUs, was  
445 approximately 0.5 hours for exSTRa with  $10^4$  permutations (12.6 hours for  $10^6$   
446 permutations), 0.6 hours for ExpansionHunter, 1.6 hours for TREDPARSE, and 2,300  
447 hours for STRetch. STRetch requires that data is realigned to its custom reference

448 genome, which comprises the majority of computation time and also creates  
449 additional data storage requirements.

450

## 451 **Discussion**

452 Genomic medicine, which uses genomic information about an individual as part of  
453 their clinical care, promises better patient outcomes and a more efficient health  
454 system through rapid diagnosis, early intervention, prevention and targeted therapy.<sup>24</sup>

455 <sup>25</sup> A single affordable front-line test that is able to comprehensively detect the genetic  
456 basis of human disease is the ultimate goal of diagnostics for genomic medicine and  
457 represents the logical way forward in an era of personalized medicine. Screening tests  
458 will play a major role in the implementation of preventative medicine.

459

460 Currently, the diagnostic pathway for suspected repeat expansion disorders utilizes  
461 single gene tests or small target panels, employing a condition-by-condition approach.  
462 This method is cost effective when the clinical diagnosis is straightforward. However,  
463 for some disorders, such as spinocerebellar ataxias, the ‘right’ test is not immediately  
464 obvious.<sup>26</sup> Many families remain unsolved, even after extensive genetic studies  
465 encompassing both gene sequencing and expansion repeat testing.<sup>26</sup> The  
466 implementation of a single NGS-based test that could identify causal point mutations,  
467 indels and expanded STRs is likely to be cost effective in this context. NGS-based  
468 tests will act as a screening tool, to identify putative expansions, which then need to  
469 be followed up with gold-standard methods such as Southern blot analysis or repeat-  
470 primed PCR. Pathogenicity will need to be determined by clinical geneticists once  
471 the precise make-up of the repeat is determined. SNVs and indels detected in NGS  
472 also have to be validated and clinically interpreted. Detecting repeat expansions using

473 NGS-based tests would include both increased diagnostic yield and a reduction in the  
474 diagnostic odyssey for many affected individuals.

475

476 Previously described methods such as hipSTR<sup>14</sup>, attempt to genotype STRs, i.e.  
477 estimate the allele sizes, which renders them ineffective when the repeat size exceeds  
478 the read length of the sequencing platform. To address this shortcoming several  
479 methods have now been developed that are designed to specifically call repeat  
480 expansions. By examining performance using >100 individuals known to have repeat  
481 expansions, spanning twelve different repeat expansion disorders, we show that  
482 exSTRa, does not require PCR-free library sequencing protocols, nor even WGS, to  
483 detect repeat expansions. We show that exSTRa delivers consistent, robust results in  
484 simulation studies.

485

486 exSTRa analysis can be run in a self contained cohort of modest size (>15  
487 individuals). It does not require any individuals that are known to be unaffected by  
488 repeat expansions because it makes use of expanded individuals as ‘controls’ for other  
489 loci by using all available data with its robust outlier detection method. exSTRa  
490 determines significance of the outlier test statistic by simulation from the cohort using  
491 a robust estimator. Hence, the default setting for exSTRa requires that not >15% of  
492 individuals in the cohort have the same repeat expansion. exSTRa has a trimming  
493 parameter which can be adjusted. Trimming too many observations leads to non-  
494 robust results. The default setting is 15%, but this can be increased up to 50% and can  
495 be assessed for performance with the ECDF plots. This was applied to the WGS\_PF  
496 cohort, which had large numbers of FRAXA (56/118, 47%) and FRDA individuals  
497 (25/118, 21%). Real disease cohorts, even ascertained from patients with diseases

498 such as spinocerebellar ataxia, which is known to be enriched for repeat expansions,  
499 are highly unlikely to reach >15% contributions from one particular repeat expansion,  
500 based on known frequencies of such expansions.

501

502 We show that exSTRa detected the most repeat expansions across all platforms and  
503 STR loci tested. It outperforms other methods at some loci, such as FRAXE, which is  
504 the highest frequency Mendelian cause of autism. exSTRa performs well in cohorts  
505 with sequencing data with more restrictions on size-fragments and greater PCR  
506 artifacts, such as WES and WGS with PCR-based library preparations. Other  
507 advantages are that it can be run with fewer requirements (no controls necessary, no  
508 size thresholds) and its graphical ECDF representation, which allows QC and fine-  
509 tuning of analysis. The exSTRa input file is easily amended to add further loci beyond  
510 the 21 investigated. These can be determined by making use of the Tandem Repeat  
511 Finder output in the UCSC genome browser. As part of the GitHub exSTRa archive  
512 we also supply an additional input file of STRs consisting of a genome wide list of  
513 STR loci that are specifically expressed in brain. This file can be amended by the user  
514 to target specific areas of the genome, such as regions identified in linkage analysis.  
515 In comparison, ExpansionHunter and TREDPARSE (for the threshold model)  
516 currently require knowledge of the pathogenic allele size, which will not be known  
517 for novel repeat expansion loci. STRetch investigates all STRs listed in its input file  
518 simultaneously and uses its novel decoy chromosome method, facilitating genome  
519 wide analysis. However this requires re-alignment to an augmented chromosome. We  
520 also found that the decoy chromosome method does not perform well with short  
521 expansions such as SCA6, since these shorter expanded alleles will preferentially find  
522 other sites in the genome, rather than the augmented genome (data not shown).

523 exSTRa does not attempt to call allele sizes, which TREDPARSE, ExpansionHunter  
524 and STRetch infer. However, gold standard validation with repeat-primed PCR or  
525 Southern blot still needs to occur prior to return of the genetic findings, and these  
526 methods size alleles more accurately than the NGS-based methods<sup>27</sup>.

527

528 We have not investigated the impact of different aligners in detail, but examination of  
529 ECDFs from the same cohort but aligned with BWA and Bowtie, the two most  
530 commonly used aligners, show highly concordant results. The ability to use existing  
531 alignments is a valuable time saving step for STR expansion analysis. exSTRa's  
532 ECDF plots inform researchers if re-alignment is necessary or not when batches from  
533 different cohorts are combined. Combining cohorts across sequencing platforms is not  
534 advisable because motif capture and hence distributions of motif sizes differ between  
535 platforms leading to batch effects.

536

537 Some expansion alleles show population heterogeneity in allele sizes, which could  
538 influence the inference of expansions with exSTRa, but will also affect other repeat  
539 expansion detection methods since they also implicitly assume homogeneity of repeat  
540 expansion distributions. One advantage of exSTRa in this context is that the ECDF  
541 method allows assessments of the results for such features. If appropriate, population  
542 heterogeneity/membership can be assessed with methods such as PLINK<sup>28</sup> or  
543 PEDDY<sup>29</sup>, allowing the identification and removal of population outliers or  
544 stratification of cohorts. Furthermore the exSTRa ECDF method allows assessments  
545 of the results for such features.

546

547 In the context of our results, exSTRa, and the other three methods appear to have  
548 potential as a population screening tool for carrier status. For example, all the  
549 methods should be able to identify carriers for Friedreich's ataxia, the most prevalent  
550 of the inherited ataxias, with a carrier frequency of ~1/100 with high sensitivity and  
551 specificity. More broadly, although the current version of exSTRa performed  
552 suboptimally for detection of *FMRI* expansions, we believe these limitations can be  
553 resolved with further refinements of exSTRa or similar detection methods. Fragile X  
554 syndrome (FXS) is the most common cause of inherited ID. Approximately 1/300  
555 individuals carry a premutation allele (55-200 repeats) which causes fragile X-  
556 associated tremor ataxia syndrome and fragile X primary ovarian insufficiency<sup>30</sup>.  
557 Currently, newborn/carrier screening is not performed for FXS. Historically, there  
558 was no medical advantage to early detection of FXS, although recent targeted  
559 treatments have shown potential benefits.<sup>31; 32</sup> There is now discussion regarding the  
560 clinical utility of screening *FMRI* for reproductive and personal healthcare.<sup>33</sup>

561

562 Given that the genetic basis of disease in many affected individuals currently remains  
563 unsolved, even after extensive genetic sequencing, we recommend the introduction of  
564 a protocol, such as exSTRa, into any standard sequencing analysis pipeline and that  
565 this be run both prospectively and retrospectively. This should identify missed repeat  
566 expansions in individuals that have only been tested for a subset of common repeat  
567 expansions, which is standard clinical practice, and will also expedite the diagnosis of  
568 individuals potentially suffering from a repeat expansion disorder. There are already  
569 >20 known repeat expansion loci, but more are likely awaiting discovery. In OMIM  
570 there are additional putative SCA loci, such as SCA25 (OMIM #608703, 2p21-p13),

571 with as yet unidentified genetic causes, but which are potentially due to novel  
572 pathogenic repeat expansions.

573

574 With large cohorts and further improvements in methodology, we believe methods  
575 such as exSTRa and future developments will facilitate the discovery of novel repeat  
576 expansion loci, which, in turn, will identify the etiology of neurodegenerative  
577 disorders in more affected individuals and families. exSTRa enables fast discovery of  
578 repeat expansions in next generation sequencing discovery cohorts including  
579 retrospective cohorts consisting mainly of WES data or WGS PCR-free library  
580 preparation data. An important new challenge lies in novel repeat expansions that are  
581 *de novo*<sup>4, 5</sup>, and not represented in the reference set of STRs that all four methods  
582 need to stipulate at which genomic locations to test. Addressing this current limitation  
583 of all RE detection algorithms will require refinement of existing/ the development of  
584 new bioinformatics tools.

585

586 The identification of a potentially pathogenic repeat expansion using detection  
587 methods such as exSTRa, should not replace the current diagnostic, locus-specific,  
588 PCR-based tests. Firstly, these will remain gold-standard, with higher sensitivity and  
589 specificity than the sequencing-based methods, and secondly, they give much more  
590 accurate estimates of the size of the expanded allele(s), and the makeup of the repeat,  
591 including whether there are interruptions, which has prognostic implications for the  
592 age of onset, disease progression and outcome.

593

594 We anticipate that there will be further improvements to all of the current methods  
595 that identify RE in NGS data. There are clearly sources of bias that affect certain loci



596 that are contributing to the poor performance at some of the STRs. For instance, we  
597 observed a GC bias for the repeat expansion alleles underlying FRAXA, FRAXE and  
598 FTDALS1, with far fewer reads able to capture these repeat expansions due to their  
599 extreme GC content. Notably FRAXA and FTDALS1 had substantially improved  
600 coverage with the PCR-free protocol.

601

602 Long read WGS will see further improvements in the detection of repeat expansion  
603 alleles, allowing capture of the entire expanded allele in a read fragment, but is  
604 currently not cost-effective, being almost 10 times more expensive than the prevailing  
605 Illumina HiSeq X sequencing platform. The development of methods such as exSTRa  
606 will lead to further improvements in patient care via clinical genomic sequencing.  
607 They will also facilitate the pending era of precision/preventative medicine, when  
608 screening tests will become much more prevalent. A universal single test will be cost  
609 and time effective in comparison to the array of existing tests currently required, to  
610 test for all known mutation types.

611

## 612 **Appendices**

### 613 **Alignment**

614 Alignment of each pair of FASTQ files was performed with Bowtie2<sup>21</sup> to the hg19  
615 human genome reference build in very sensitive local mode, with maximum insert  
616 sizes of 800 bp for WES samples and 1000 bp for WGS samples. BAM files were  
617 sorted and merged with the Novosort tool. Duplicate marking was performed with  
618 Picard. Local realignment and base score recalibration was performed with the GATK  
619 IndelAligner tool and the Base Quality Score Recalibration tool<sup>34</sup> to produce input  
620 ready BAM files.

621

622 **Software**

623 The first step of the analysis is performed with a Perl module, called  
624 Bio::STR::exSTRa, which carries out a heuristic procedure to extract repeat content.  
625 In summary, this procedure uses the data from the reference database for the 21 loci  
626 presented in Table 1 to identify all reads that map to each of the STR loci, for each  
627 individual to be examined. The number of repeat motifs contained by each read are  
628 determined by the heuristic procedure, which examines each read for the repeat units  
629 that that STR is known to contain. This allows for some mismatches due to impure  
630 repeats and sequencing errors. Additionally, this is more computationally efficient  
631 than determining the exact repeat start and end, and is more robust as determining the  
632 edge of the repeat can be difficult near the end of a read in the presence of  
633 mismatches.

634

635 **Bio::STR::exSTRa : A heuristic procedure to extract repeat units per read**

636 For simplicity, the following description of the data and analysis methods is only for a  
637 single locus. The algorithm is repeated independently at each locus.

638

639 Read information is extracted from a database of STR locations, such as 2–6bp repeat  
640 unit features generated using the Tandem Repeats Finder <sup>35</sup>, which is also available as  
641 the Simple Repeats track of UCSC Genome Browser. Information is extracted for one  
642 STR at a time, with the following algorithm repeated for each STR:

643

644 1. The method identifies ‘anchor’ reads that facilitates identifying reads within or  
645 overlapping the STR. To qualify as an anchor, the reads are required to map within

646 800 bp of the STR, with the anchor orientated towards the STR. An anchor may  
647 overlap the STR.

648

649 2. The anchor-mate mapping is checked. If the anchor-mate is mapped near the STR  
650 and is not overlapping or adjacent, then the read is discarded, while those reads  
651 overlapping the STR are taken forward to the next analysis step. Sometimes the read  
652 is unmapped, or mapped to another locus, which is then recovered for further  
653 interrogation in the next step.

654

655 3. Remaining anchor-mates have their sequence content matched for the presence of  
656 the repeat unit in the correct direction, allowing for the repeat to start at any base, or  
657 phase, of the repeat unit. For example, if the repeat unit is CAG, the method can also  
658 match AGC and GCA. The number of bases found to be part of the repeat unit is  
659 counted to derive a repeat-score for that read, that is designated at a given locus as  $x_{ij}$   
660 for sample  $i$  and read  $j$  (note that the maximum defined  $j$  depends on the sample). If  
661 both ends of a read-pair overlap within an STR, both reads undergo this procedure  
662 and each end is given a score that can be resolved during the statistical analysis of the  
663 data (the implementation in this paper did not investigate resolving these further, with  
664 both ends left in the analysis if any). An example of matching (lower case) a CAG on  
665 the opposite strand, thus matching CTG at any starting base, or phase, of the motif,  
666 i.e. CTG, TGC and GCT:

667

668 CGTTCAC**Cctg**GATGTGAACT**tctg**TC**tctg**ATAGGTCCCC**CctgctgctgctgctgctgctgTt**  
669 **gctgc**TTTT**gctgc**TGT**tctg**AAA

670

671 This 87 bp sequence has 48 bp marked (bold and lower case) as part of the repeat.

672

673 4. The method filters out reads where the score is lower than expected in random  
674 nucleotide sequences. While not precisely true, the assumption applied is that the four  
675 nucleotides are uniformly distributed and independent with respect to other positions.  
676 Short motifs are more likely to appear by chance. The method filters out scores where  
677  $x_{ij} < lk/4^k$ , where  $l$  is the read length and  $k$  is the motif length. 800 bp has been chosen  
678 to avoid discarding reads overlapping the STR, with the insert size of read pairs  
679 having median ~360 bp. Some protocols may need to analyse reads further than 800  
680 bp. This can be adjusted when calling the Perl module.

681

682 The output of this Perl module consists of a tab-delimited file consisting of a table  
683 where each row in the table is the repeat content of any read from a particular  
684 individual that has been identified as mapping to an STR locus that was to be  
685 investigated.

686

687 Note that these data do not represent the true size of the allele that the read has  
688 captured but where the method predicts an individual with repeat expansion allele at a  
689 particular STR locus to show an excess of reads and read content mapping to that  
690 STR.

691

692 **R package exSTRa : detecting outlier distributions of repeat content in reads**

693 Analysis methods for the second part of the analysis method are embedded in an R  
694 package, called exSTRa (expanded STR algorithm). The output data from step 1 can

695 be loaded and the data visualized. In particular visualizations of the data are  
696 performed with empirical cumulative distribution functions, or ECDFs.

697

698 The analysis of the samples is treated as an outlier detection problem. For the N  
699 individuals in the cohort the method compares each individual in turn to all others,  
700 including itself for robustness, for all STR loci that will be tested for repeat  
701 expansions. Since more reads with greater numbers of the repeat motif will be visible  
702 in an individual with a repeat expansion at a particular locus, the data at the repeat  
703 locus being interrogated is used in a statistical test of a difference of distribution in  
704 number of repeats that are observed for a particular individual in comparison to the  
705 set of controls. Individuals with an expanded repeat demonstrate a shift in the  
706 distribution in comparison to individuals with normal size alleles comprising their  
707 genotype for the STR locus being examined. To visualize the results, the output is  
708 plotted as empirical cumulative distribution functions (ECDFs) in R.

709

### 710 **Statistical Test**

711 We developed a statistical test to detect outlier samples in comparison to a  
712 background set of samples. These outlier samples are likely to be individuals  
713 harbouring repeat expansions. To apply this test the method utilizes an empirical  
714 quantile imputation procedure, implemented in the R function `quantile()`. This  
715 function calculates empirical quantiles for any desired probability, for example  
716 `probability = 0.5` generates the median observation in a dataset, but it is also capable  
717 of generating quantiles at probability points that have not been observed, by  
718 interpolating the probability distribution function based on the empirical observations.  
719 We make use of this function to firstly generate the same number of ‘observations’

720 for all samples to be tested, defined as  $M$ . In general,  $n$  is defined so that it is the  
721 largest number of observations for all of the samples, but other values could also be  
722 chosen, such as the median number of observations. The R function `quantile()` is  
723 applied to generate this dataset which consists of  $N$  samples, with  $M$   
724 observations/quantiles, leading to a dataset with  $N$  by  $M$  datapoints, or quantiles. This  
725 dataset is defined as  $Y=(y_{ij})$ , where  $y_{ij}$  is the repeat content of the  $j^{\text{th}}$  quantile from the  
726  $i^{\text{th}}$  individual.

727

728 The test statistic, which we call  $T_i$ , is defined as the average of multiple t-statistics  
729 generated at each quantile  $j$ , above a preset threshold  $0 \leq h < 1$ , which we usually  
730 define  $h = 0.5$ .

731

$$T_i = \frac{1}{D} \sum_{j:Pr(y_{ij}) \geq h}^M t_{ij}$$

$$D = |\{j : Pr(y_{ij}) \geq h\}|$$

732

733 Sixteen of the 21 STR repeat expansion loci to be examined have a dominant mode of  
734 inheritance, with only one copy of the expanded allele. This can be observed with the  
735 ECDF plots for the autosomal dominant STR loci, where deviations in the repeat  
736 composition of reads are only noticeable after the median quantile, when the y-axis  
737 (which is the probability) exceeds 0.5. Observations below this threshold are likely to  
738 carry no signal, and are thus would not contribute to any test statistic attempting to  
739 discriminate between expansions and normal sized alleles.

740

741 Each quantile test statistic,  $t_{ij}$ , is calculated similarly to a two-sample T-test like test  
742 statistic, but using a trimmed mean and variance, to robustly allow for the occurrence  
743 of more than one expansion in the background distribution, which is the case in the  
744 cohorts we tested but which will also likely be the case in other cohorts. The trimming  
745 percentage, or percentage of samples that are used is a parameter that can be set by  
746 the user in exSTRa, but the default is set at 0.15. Trimming is performed bilaterally,  
747 for both the lower and upper tails of the distributions, resulting in at least 30% of the  
748 samples being trimmed.

749

$$t_{ij} = \frac{y_{ij} - m_j}{S_j}$$

$$m_j = \frac{1}{n_j} \sum_{j:l_j \leq y_{ij} \leq u_j} y_{ij}$$

$$n_j = |\{j : l_j \leq y_{ij} \leq u_j\}|$$

$$S_j = s_j \sqrt{1 + \frac{1}{n_j}}$$

750

751

752 where  $l_j$  is the first observation included from the lower tail of the distribution after  
753 the trimmed observations and  $u_j$  the last observation included from the upper tail of  
754 the distribution, with all observations beyond this trimmed.  $s_j$  is the sample standard  
755 deviation of the trimmed samples.

756

757 We derive p-values for these test statistics using a simulation procedure.

758

759 Since the number of individuals in our simulations is not large and only test a single  
760 individual, standard permutation tests will not result in sufficient sampling of the

761 empirical distribution thus resulting in a very coarse grained empirical distribution.  
762 Instead we take advantage of the well-described empirical distributions of the samples  
763 by directly simulating from the background distribution, which represents the  
764 distribution of normal, or non-expanded alleles. We perform this using robust  
765 methods to ensure that samples with expanded alleles do not influence the simulation  
766 in the simulation study.

767

768 For simulation  $s$  we simulate  $M$  quantiles for  $N$  samples, by assuming that the  
769 distributions at each quantile follow large sample theory and are thus approximately  
770 normally distributed with mean  $m_j$  and standard deviation  $d_j$ , where  $j$  denotes the  
771 quantile. The method then tests this assumption by performing visual inspections of  
772 the distribution of quantiles after standardization with the R function `qqnorm()` and  
773 the approximation was reasonable.

774

775 The method then uses the median as our estimator for the mean, and the median  
776 absolute deviation (MAD) as our robust estimator for the standard deviation. Thus,

$$\hat{m}_j = \text{median}\{y_{.j}\}$$

$$\hat{d}_j = \frac{1}{(\Phi^{-1}(3/4))} \text{MAD}\{y_{.j}\}$$

$$\text{MAD}\{y_{.j}\} = \text{median}\{|y_{ij} - \text{median}\{y_{.j}\}|\}$$

777

778 Where  $\Phi^{-1}$ , and  $\Phi$  is the inverse of the cumulative distribution  
779 function of the standard normal distribution. The R function `mad()` incorporates the  
780 scaling factor that ensures consistency with the standard deviation when observations  
781 are normally distributed.

782



783 The method then uses the `rnorm()` function in R to randomly generate the N new  
784 observations for each quantile, using the STR locus and quantile specific estimators  
785 for the mean and standard deviation. The data is then sorted for each sample, as some  
786 of the new observations are no longer monotonically increasing as per definition of  
787 quantiles.

788

789 Finally, the test statistic  $T_s$  is calculated as defined above, but using the new data set  
790 generated from the simulation, where the first sample in the simulated data set is  
791 arbitrarily chosen to be the sample to be tested as an outlier. The method then repeat  
792 this for a desired number of simulations, say B, and then calculates the empirical p-  
793 value for our test statistic using standard methods, where:

794

$$p_{T_i} = \frac{\sum_{s=1}^B I([T_i > T_1^s]) + 1}{B + 1}$$

795

796 Here  $I(\cdot)$  is the indicator function.  $T_1^S$  is the test statistic for the dataset. The method  
797 calls individuals as expanded or not for each STR locus examined based on a  
798 Bonferroni corrected threshold at the 0.05 significance level, based on the number of  
799 STR tested for each sample.

800

801 Standard deviations for the empirical p-value estimator were also calculated as  
802 follows.

$$SD(\hat{p}) = \sqrt{\frac{1 + \sum_{i=1}^B x_i (1 - \frac{\sum_{i=1}^B x_i}{B+1})}{B}}$$
$$x_i = I([T_i > T_1^S])$$

803

804 **Calling expansions with ExpansionHunter, STRetch and TREDPARSE**

805 We performed analysis with ExpansionHunter (version 2.5.3), STRetch (GitHub  
806 commit 94d0516) and TREDPARSE (GitHub commit 83881b4), on the cohorts at the  
807 21 repeat expansion loci listed in Table 1. The input data was the same BAM files  
808 generated as described above. Only specification files (in JSON format) for the DM1,  
809 DRPLA, FRAXA, FRDA, FTDALS1, HD, SBMA, SCA1 and SCA3 loci were  
810 provided with ExpansionHunter. The JSON files for the remaining loci were obtained  
811 by personal communication with Egor Dolzhenko (Illumina, Inc. San Diego, CA,  
812 USA). For data aligned with bowtie2, the --min-anchor-mapq parameter was set to  
813 44, while for the original alignments of the Coriell samples this parameter was set to  
814 60. The --read-depth parameter was set the median coverage for each sample in the  
815 WES\_PCR cohort, otherwise this was computed by ExpansionHunter for the WGS  
816 samples. The list of STR loci provided with STRetch does not include FRDA, which  
817 was added manually. The EPM1 repeat motif is 12 bp and is not assessed using  
818 STRetch, which aligns to an augmented reference genome containing a decoy  
819 chromosome for each STR repeat motif up to 6 bp in size.

820

821 ExpansionHunter and TREDPARSE-T call allele lengths and genotypes. To call  
822 individuals as having expansions requires the user to define thresholds on allele sizes  
823 as to what constitutes an appropriate threshold. For FRAXA, we additionally tested  
824 using the premutation threshold (labelled FRAXA\_pre), in addition to testing for full  
825 expansions. To call an expansion, we used the same thresholds as Dolzhenko et al<sup>17</sup>  
826 (based on McMurray<sup>36</sup>) or the largest reported normal allele size at other loci. Other  
827 thresholds will change the sensitivity and specificity. TREDPARSE-L expansions  
828 calls were recorded for all samples labelled as “risk”. exSTRA p-values were  
829 Bonferroni corrected over the number of STRs tested. STRetch reports p-values

830 adjusted for multiple testing over all STRs genome wide, however unadjusted p-  
831 values were extracted and Bonferroni corrected over just the number of STRs tested.  
832 A threshold of  $p < 0.05$  was used for significance.

833

#### 834 **Supplemental Data**

835 Supplemental Data includes 13 figures and 5 tables.

836

#### 837 **Acknowledgements**

838 We would like to thank Egor Dolzhenko and Michael Eberle (Illumina Inc), who  
839 produced the STR specification files for ExpansionHunter and gave access to the  
840 EGA00001003562. We thank Leslie Burnett, Ben Lundie, Katie Ayres and Andrew  
841 Sinclair for access to control datasets. We thank Kate Pope and Greta Gillies for  
842 assistance with recruitment and sample preparation.

843

844 RT was supported by an Australian Postgraduate Award and funding from the Edith  
845 Moffat fund. PJJ was supported by NHMRC CDA2 (GNT1032364). MB was  
846 supported by NHMRC Program Grant (GNT1054618) and NHMRC Senior Research  
847 Fellowship (APP1102971). This work was supported by the Victorian Government's  
848 Operational Infrastructure Support Program and Australian Government National  
849 Health and Medical Research Council Independent Research Institute Infrastructure  
850 Support Scheme (NHMRC IRIISS).

851

#### 852 **Web Resources**

853 exSTRa <http://github.com/bahlolab/exSTRa>

854 ExpansionHunter <https://github.com/Illumina/ExpansionHunter>

855 TREDPARSE <https://github.com/humanlongevity/tredparse>

856 STRetch <https://github.com/Oshlack/STRetch>

857 Picard <http://broadinstitute.github.io/picard/>

858 Novosort <http://www.novocraft.com/products/novosort/>

859 OMIM <https://www.omim.org>

860 GATK IndelAligner <https://software.broadinstitute.org/gatk/>

861 Coriell <https://www.coriell.org/>

862

### 863 **Figure Legends**

864

865 **Figure 1** ECDF of repeat expansion composition of reads from the WES cohort,  
866 depicting four different known repeat expansion disorders captured by WES (HD,  
867 SCA2, SCA6 and SCA1). Sample rptWEHI3 (blue) is a known HD repeat expansion  
868 patient. The expanded allele size is not known. Sample rptWEHI1 (yellow) a known  
869 SCA2 repeat expansion of length 42 repeats, sample rptWEHI2 (red) a known SCA6,  
870 of length 22 repeats, and sample rptWEHI4 (green) a known SCA1 patient, of length  
871 52 repeats. The title at the top of each individual figure gives the locus being  
872 examined, the reference number of repeats in the hg19 human genome reference with  
873 the corresponding number of bps, and the smallest reported expanded allele in the  
874 literature (with the corresponding number of bps in brackets). The blue dashed  
875 vertical line in the plot denotes the largest known normal allele, the red dashed  
876 vertical line denotes the smallest known expanded allele.

877

878 **Figure 2** ECDFs of repeat expansion composition of reads from the WGS\_PCR\_2

879 cohort, depicting four different STR loci (top left = SCA1, length of the expanded

880 alleles are 52 and 45 repeats; top right = FRDA, length of the expanded alleles are  
881 320 and 788 repeats; bottom left = SCA7, length of the expanded allele is 39; bottom  
882 right = DM1, length of the expanded alleles are 173 and 83 repeats). Here coloured  
883 samples at each STR indicate those called by exSTRa as repeat expansions at the STR  
884 locus. The title at the top of each individual figure gives the locus being examined, the  
885 reference number of repeats in the hg19 human genome reference with the  
886 corresponding number of bps, and the smallest reported expanded allele in the  
887 literature (with the corresponding number of bps in brackets). The blue dashed  
888 vertical line in the plot denotes the largest known normal allele, the red dashed  
889 vertical line denotes the smallest known expanded allele.

890

891 **Figure 3.** ECDFs for four repeat expansion loci from WGS\_PF\_3 cohort .Top left,  
892 DM1; top right, FRDA; bottom left, FRAXA; bottom right, HD .The title at the top of  
893 each individual figure gives the locus being examined, the reference number of  
894 repeats in the hg19 human genome reference with the corresponding number of bps,  
895 and the smallest reported expanded allele in the literature (with the corresponding  
896 number of bps in brackets). The blue dashed vertical line in the plot denotes the  
897 largest known normal allele, the red dashed vertical line denotes the smallest known  
898 expanded allele.

899

### 900 **Table Legends**

901 **Table 1** Detailed STR loci information. TRF = Tandem Repeats Finder (Benson et al,  
902 1999). TRF match and TRF indel describe the purity of the repeat. AD = autosomal  
903 dominant, X = X-linked, AR = autosomal recessive.

904

905 **Table 2** Repeat type, genetic model, diseases, sample names and which cohorts  
906 samples appear in. Allele sizes are derived from standard laboratory tests for repeat  
907 expansions. Some individuals were not tested (Not sized) or the data was not  
908 available (not recorded). MOI, mode of inheritance; AD, autosomal dominant; X, X-  
909 linked, AR; autosomal recessive. Only the total number of controls are given denoted  
910 by (controls).

911 **Table 3** Repeat Expansion detection results for exSTRa, ExpansionHunter, STRetch  
912 and TREDPARSE over all four cohorts. TP, true positive; FN, false negative; FP,  
913 false positive; TN, true negative, Sensitivity,  $TP/(TP+FN)$ ; Specificity,  $TN/(FP+TN)$ ;  
914 NA, not applicable. WES cohort labeled with (\*) only assessed over eleven STR loci  
915 in the capture design. WGS\_PCR\_2 was also analysed split into two sub-cohorts, split  
916 by flow cell lane, and are designated as WGS\_PCR\_2\_30X\_1 and  
917 WGS\_PCR\_2\_30X\_2.

918

919

920

921

922

Disease	Symbol	OMIM	Inheritance	Gene	Cytogenetic Location	Type	Repeat Motif	Normal Range	Expansion Range	Strand	Start hg19	Reference Repeat Number	TRF Match (%)	TRF Indel (%)	Reference STR size (bp)
Huntington disease	HD	143100	AD	<i>HTT</i>	4p16.3	Coding	CAG	6-34	36-100+	+	3,076,604	21.3	96	0	64
Kennedy disease	SBMA	313200	X	<i>AR</i>	Xq12	Coding	CAG	9-35	38-62	+	66,765,159	33.3	86	9	103
Spinocerebellar ataxia 1	SCA1	164400	AD	<i>ATXN1</i>	6p23	Coding	CAG	6-38	39-82	-	16,327,865	30.3	95	0	91
Spinocerebellar ataxia 2	SCA2	183090	AD	<i>ATXN2</i>	12q24	Coding	CAG	15-24	32-200	-	112,036,754	23.3	97	0	70
Machado-Joseph disease	SCA3	109150	AD	<i>ATXN3</i>	14q32.1	Coding	CAG	13-36	61-84	-	92,537,355	14	84	0	42
Spinocerebellar ataxia 6	SCA6	183086	AD	<i>CACNA1A</i>	19p13	Coding	CAG	4-7	21-33	-	13,318,673	13.3	100	0	40
Spinocerebellar ataxia 7	SCA7	164500	AD	<i>ATXN7</i>	3p14.1	Coding	CAG	4-35	37-306	+	63,898,361	10.7	100	0	32
Spinocerebellar ataxia 17	SCA17	607136	AD	<i>TBP</i>	6q27	Coding	CAG	25-42	47-63	+	170,870,995	37	94	0	111
Dentatorubral-pallidolusian atrophy	DRPLA	125370	AD	<i>DRPLA/ATN1</i>	12p13.31	Coding	CAG	7-34	49-88	+	7,045,880	19.7	92	0	59
Huntington disease-like 2	HDL2	606438	AD	<i>JPH3</i>	16q24.3	Exon	CTG	7-28	66-78	+	87,637,889	15.3	95	4	47
Fragile-X site A	FRAXA	300624	X	<i>FMR1</i>	Xq27.3	5'UTR	CGG	6-54	200-1000+	+	146,993,555	25	90	5	75
Fragile-X site E	FRAXE	309548	X	<i>FMR2</i>	Xq28	5'UTR	CCG	4-39	200-900	+	147,582,159	15.3	100	0	46
Myotonic dystrophy 1	DM1	160900	AD	<i>DMPK</i>	19q13	3'UTR	CTG	5-37	50-10000	-	46,273,463	20.7	100	0	62
Friedreich ataxia	FRDA	229300	AR	<i>FXN</i>	9q13	Intron	GAA	6-32	200-1700	+	71,652,201	6.7	100	0	20
Myotonic dystrophy 2	DM2	602668	AD	<i>ZNF9/CNBP</i>	3q21.3	Intron	CCTG	10-26	75-11000	-	128,891,420	20.8	92	0	83

Frontotemporal dementia and/or amyotrophic lateral sclerosis 1	FTDALS1	105550	AD	<i>C9orf72</i>	9p21	Intron	GGGGCC	2-19	250-1600	-	27,573,483	10.8	74	8	62
Spinocerebellar ataxia 36	SCA36	614153	AD	<i>NOP56</i>	20p13	Intron	GGCCTG	3-8	1500-2500	+	2,633,379	7.2	97	0	43
Spinocerebellar ataxia 10	SCA10	603516	AD	<i>ATXN10</i>	22q13.31	Intron	ATTCT	10-20	500-4500	+	46,191,235	14	100	0	70
Myoclonic epilepsy of Unverricht and Lundborg	EPM1	254800	AR	<i>CSTB</i>	21q22.3	Promoter	CCCCGCC	2-3	40-80	-	45,196,324	3.1	100	0	37
Spinocerebellar ataxia 12	SCA12	604326	AD	<i>PPP2R2B</i>	5q32	Promoter	CAG	7-45	55-78	-	146,258,291	10.7	100	0	32
Spinocerebellar ataxia 8	SCA8	608768	AD	<i>ATXN8OS/ATXN8</i>	13q21	utRNA	CTG	16-34	74+	+	70,713,516	15.3	100	0	46
Spinocerebellar ataxia 31	SCA31	117210	AD	<i>BEANI/TK2</i>	16q21	Intron	TGGAA <sup>a</sup>	0	2.5-3.8kb <sup>b</sup>	+	66,524,302	0	N/A	N/A	N/A
Spinocerebellar ataxia 37	SCA37	615945	AD	<i>DABI</i>	1p32.3	Intron	ATTTC <sup>a</sup>	0	31-75	-	57,832,716 <sup>c</sup>	0	N/A	N/A	N/A
Familial adult myoclonic epilepsy 1 <sup>c</sup>	FAME1/ BAFME1	601068	AD	<i>SAMD12</i>	8q24	Intron	TTTCA <sup>a</sup>	0	440-3,680 <sup>f</sup>	-	119,379,055 <sup>d</sup>	0	N/A	N/A	N/A

923 **Table 1 Short tandem repeat loci information for STRs causing neurogenetic disorders. TRF, Tandem Repeats Finder (Benson et al, 1999). TRF match and TRF indel**  
924 **describe the purity of the repeat. AD, autosomal dominant; X, X-linked; AR, autosomal recessive; UTR, untranslated region. <sup>a</sup>These repeat expansions are novel insertions**  
925 **and thus not represented in the reference genome at their respective locations. <sup>b</sup>SCA31 is caused by the insertion of a complex repeat containing (TGGAA)<sub>n</sub>; hence the**  
926 **length is given in as the length of the expanded repeats in bps, instead of repeat number. <sup>c</sup>The SCA37 physical map location is given at the reference (ATTTT)<sub>n</sub> repeat,**  
927 **where affected individuals have the pathogenic (ATTTC)<sub>n</sub> inserted. <sup>d</sup>The FAME1 physical map location is given as the position of the reference (TTTTA)<sub>n</sub> repeat, at which**



928 affected individuals have (TTTCA)<sub>n</sub> inserted. <sup>e</sup>Ishiura et al. identified similar expansions associated with FAME6 and FAME7, in the genes TNRC6A and RAPGEF2  
929 respectively, but only in single families. These have not been listed. <sup>f</sup>The FAME1 repeat size is the estimated size of the combined expanded (TTTCA)<sub>n</sub> and the (TTTTA)<sub>n</sub>  
930 reference repeat.

931

932	Class	MOI	Diagnosis	Allele sizes	Gender	WES_PCR	WGS_PCR_1	WGS_PCR_2
933	PolyQ	AD	HD	Not recorded	male	rptWEHI3	HD-1	
934	PolyQ	AD	HD	17,39	female			WGSrpt_10
935	PolyQ	AD	HD	20,42	male			WGSrpt_12
936	PolyQ	AD	SCA1	36,52	female	rptWEHI4		WGSrpt_14
937	PolyQ	AD	SCA1	30,45	male			WGSrpt_16
938	PolyQ	AD	SCA2	21,42	female	rptWEHI1	SCA2-1	WGSrpt_18
939	PolyQ	AD	SCA2	23,39	male			WGSrpt_20
940	PolyQ	AD	SCA6	11,22	female	rptWEHI2	SCA6-1	WGSrpt_05
941	PolyQ	AD	SCA6	10,21	female			WGSrpt_07
942	PolyQ	AD	SCA7	13,39	female			WGSrpt_08
943	5'UTR	X	FRAXA	Not sized	male			WGSrpt_17
944	5'UTR	X	FRAXA	613-1680	male			WGSrpt_19
945	5'UTR	X	FRAXA (pre)	~100	female			WGSrpt_21
946	3'UTR	AD	DM1	8,173	female			WGSrpt_13
947	3'UTR	AD	DM1	13,83	male			WGSrpt_15
948	Intron	AR	FRDA	320,320	male			WGSrpt_09
949	Intron	AR	FRDA	788,788	male			WGSrpt_11
950			(controls)			58	14	2
951								

952 **Table 2. Repeat type, genetic model, diseases, sample names and which cohorts samples appear in. Allele sizes are derived from standard laboratory tests for repeat**  
 953 **expansions. Some individuals were not tested (Not sized) or the data was not available (not recorded). MOI = mode of inheritance (AD = autosomal dominant, X = X-**  
 954 **linked, AR = autosomal recessive). Only the total number of controls are given denoted by (controls).**

955	Class	MOI	Diagnosis	Expanded	Affected	Not Expanded
956	PolyQ	AD	HD	13	13	105
957	PolyQ	AD	SCA1	3	3	115
958	PolyQ	AD	SCA3	1	1	117
959	PolyQ	AD	DRPLA	2	2	116
960	PolyQ	AD	SBMA	1	1	117
961	5'UTR	X	FRAXA	16	16	102
962	5'UTR	X	FRAXA (pre)	33	21	85
963	3'UTR	AD	DM1	17	17	101
964	Intron	AR	FRDA	25	14	93
965			Total (FRAXA) <sup>a</sup>	78		40
966			Total (FRAXA pre)	95		23

967 **Table 3 WGS\_PF cohort. Cohort of 118 individuals sequenced with Illumina PCR-free library preparation. Only total number of samples are listed, rather than actual**  
 968 **samples. Details of samples are available in Dolzhenko et al 2017. <sup>a</sup>Total only includes FXS individuals, and no intermediate pre expansions.**

969

970

971	<b>Cohort</b>	<b>Cases</b>	<b>Controls<sup>^</sup></b>	<b>Method</b>	<b>TP</b>	<b>FN</b>	<b>TN</b>	<b>FP</b>	<b>Sensitivity</b>	<b>Specificity</b>
972	WES_PCR*	4	58	exSTRa	4	0	607	9	1	0.99
973				ExpansionHunter	2	2	616	0	0.5	1
974				STRetch <sup>a</sup>	3	1	613	3	0.75	1
975				TREDPARSE-T <sup>b</sup>	4	0	585	31	1	0.95
976				TREDPARSE-L <sup>b</sup>	4	0	574	42	1	0.93
977	WGS_PCR_1	3	14	exSTRa	2	1	343	11	0.67	0.97
978				ExpansionHunter	3	0	354	0	1	1
979				STRetch <sup>a</sup>	1	2	336	1	0.33	1
980				TREDPARSE-T <sup>b</sup>	3	0	354	0	1	1
981				TREDPARSE-L <sup>b</sup>	3	0	354	0	1	1
982	WGS_PCR_2	16	2	exSTRa	13	3	352	10	0.81	0.97
983				ExpansionHunter	8	8	362	0	0.5	1
984				STRetch <sup>a</sup>	11	5	338	6	0.69	0.98
985				TREDPARSE-T <sup>b</sup>	12	4	362	0	0.75	1
986				TREDPARSE-L <sup>b</sup>	11	5	362	0	0.69	1
987	WGS_PCR_2_30X_1	16	2	exSTRa	13	3	357	5	0.81	0.99
988				ExpansionHunter	8	8	362	0	0.5	1
989				STRetch <sup>a</sup>	11	5	340	4	0.69	0.99
990				TREDPARSE-T <sup>b</sup>	13	3	362	0	0.81	1
991				TREDPARSE-L <sup>b</sup>	9	7	362	0	0.56	1
992	WGS_PCR_2_30X_2	16	2	exSTRa	12	4	354	8	0.75	0.98
993				ExpansionHunter	8	8	362	0	0.5	1

994				STRetch <sup>a</sup>	11	5	336	8	0.69	0.98
995				TREDPARSE-T <sup>b</sup>	13	3	362	0	0.81	1
996				TREDPARSE-L <sup>b</sup>	10	6	362	0	0.62	1
997	WGS_PF	77	41	exSTRa	60	17	2330	71	0.78	0.97
998		77	41	ExpansionHunter <sup>c</sup>	62	15	2395	6	0.81	1
999		96	22	EH FRAXA_pre <sup>d</sup>	95	1	2374	8	0.99	1
1000		96	22	STRetch <sup>a</sup>	62	15	2207	76	0.81	0.97
1001		96	22	TREDPARSE-T <sup>b</sup>	52	25	2384	17	0.68	0.99
1002		96	22	TP-T FRAXA_pre <sup>d</sup>	72	24	2364	18	0.75	0.99
1003		66	52	TREDPARSE-L <sup>b</sup>	34	32	2396	16	0.52	0.99
1004		72	46	TP-L FRAXA_pre <sup>d</sup>	48	24	2383	23	0.67	0.99
1005	WGS_PF (no FRAXA)	62	56	exSTRa	52	10	2231	67	0.84	0.97
1006				ExpansionHunter <sup>c</sup>	61	1	2292	6	0.98	1
1007				STRetch <sup>a</sup>	62	0	2104	76	1	0.97
1008				TREDPARSE-T <sup>b</sup>	52	10	2281	17	0.84	0.99
1009		51	67	TREDPARSE-L <sup>b</sup>	34	17	2293	16	0.67	0.99

1010 **Table 4 Repeat expansion detection results for all four cohorts. ^Individuals designated as controls have no known repeat expansions. Individuals designated as cases**  
1011 **have one known repeat expansion, but are controls for all other loci tested. TP, true positive; FN, false negative, TN, true negative; FP, false positive; Sensitivity,**  
1012 **TP/(TP+FN); Specificity, TN/(FP+TN); WES cohort labeled with (\*) only assessed over ten STR loci in the capture design. WGS\_PCR\_2 was also analysed split into two sub-**  
1013 **cohorts, split by flow cell lane, and are designated as WGS\_PCR\_2\_30X\_1 and WGS\_PCR\_2\_30X\_2. <sup>a</sup>STRetch was Bonferroni corrected for the same number of tests as the**  
1014 **other methods, and not genome-wide corrected. <sup>b</sup>TREDPARSE results are given for the repeat expansion size threshold method (TREDPARSE-T) and for the likelihood**  
1015 **ratio test based method (TREDPARSE-L). For STR loci with recessive inheritance, samples with double expansions were designated as cases for TREDPARSE-L, which**

1016 takes into account the inheritance model. <sup>c</sup>For the WGS\_PF cohort the original ExpansionHunter results from Dolzhenko et al were used, which make use of reads aligned  
1017 with a different aligner. <sup>d</sup>For the WGS\_PF cohort, additional results were computed using the premutation threshold to test for FRAXA expansions with ExpansionHunter  
1018 (EH FRAXA\_pre), TREDPARSE-T (TP-T FRAXA\_pre) and TREDPARSE-L (TP-L FRAXA\_pre).

1019

1020 **References**

- 1021 1. Benson, G. (1999). Tandem repeats finder: a program to analyze DNA sequences.  
1022 Nucleic Acids Res 27, 573-580.
- 1023 2. Jones, L., Houlden, H., and Tabrizi, S.J. (2017). DNA repair in the trinucleotide  
1024 repeat disorders. The Lancet Neurology 16, 88-96.
- 1025 3. Hannan, A.J. (2018). Tandem repeats mediating genetic plasticity in health and  
1026 disease. Nat Rev Genet, 1-13.
- 1027 4. Seixas, A.I., Loureiro, J.R., Costa, C., Ordóñez-Ugalde, A., Marcelino, H.,  
1028 Oliveira, C.L., Loureiro, J.L., Dhingra, A., Brandão, E., Cruz, V.T., et al.  
1029 (2017). A Pentanucleotide ATTTTC Repeat Insertion in the Non-coding Region  
1030 of DAB1 , Mapping to SCA37 , Causes Spinocerebellar Ataxia. The American  
1031 Journal of Human Genetics 101, 87-103.
- 1032 5. Ishiura, H., Doi, K., Mitsui, J., Yoshimura, J., Matsukawa, M.K., Fujiyama, A.,  
1033 Toyoshima, Y., Kakita, A., Takahashi, H., Suzuki, Y., et al. (2018).  
1034 Expansions of intronic TTTC A and TTTTA repeats in benign adult familial  
1035 myoclonic epilepsy. Nat Genet, 1-14.
- 1036 6. Schulz, J.B., Boesch, S., Bürk, K., Durr, A., Giunti, P., Mariotti, C., Pousset, F.,  
1037 Schöls, L., Vankan, P., and Pandolfo, M. (2009). Diagnosis and treatment of  
1038 Friedreich ataxia: a European perspective. Nat Rev Neurol 5, 222-234.
- 1039 7. Seltzer, M.M., Baker, M.W., Hong, J., Maenner, M., Greenberg, J., and Mandel, D.  
1040 (2012). Prevalence of CGG expansions of the FMR1 gene in a US population-  
1041 based sample. Am J Med Genet B Neuropsychiatr Genet 159B, 589-597.
- 1042 8. Tassone, F. (2014). Newborn screening for fragile X syndrome. JAMA Neurol 71,  
1043 355-359.

- 1044 9. Genetic Modifiers of Huntington's Disease, C. (2015). Identification of Genetic  
1045 Factors that Modify Clinical Onset of Huntington's Disease. *Cell* 162, 516-  
1046 526.
- 1047 10. Bettencourt, C., Hensman Moss, D.J., Flower, M., Wiethoff, S., Brice, A., Goizet,  
1048 C., Stevanin, G., Koutsis, G., Karadima, G., Panas, M., et al. (2016). DNA  
1049 repair pathways underlie a common genetic mechanism modulating onset in  
1050 polyglutamine diseases. *Ann Neurol*.
- 1051 11. Németh, A.H., Kwasniewska, A.C., Lise, S., Parolin Schneckenberg, R., Becker,  
1052 E.B.E., Bera, K.D., Shanks, M.E., Gregory, L., Buck, D., Zameel Cader, M.,  
1053 et al. (2013). Next generation sequencing for molecular diagnosis of  
1054 neurological disorders using ataxias as a model. *Brain* 136, 3106-3118.
- 1055 12. Tae, H., Kim, D.-Y., McCormick, J., Settlage, R.E., and Garner, H.R. (2014).  
1056 Discretized Gaussian mixture for genotyping of microsatellite loci containing  
1057 homopolymer runs. *Bioinformatics* 30, 652-659.
- 1058 13. Gymrek, M., Golan, D., Rosset, S., and Erlich, Y. (2012). lobSTR: A short  
1059 tandem repeat profiler for personal genomes. *22*, 1154-1162.
- 1060 14. Willems, T., Zielinski, D., Yuan, J., Gordon, A., Gymrek, M., and Erlich, Y.  
1061 (2017). Genome-wide profiling of heritable and de novo STR variations. *Nat*  
1062 *Methods* 14, 590-592.
- 1063 15. Highnam, G., Franck, C., Martin, A., Stephens, C., Puthige, A., and Mittelman, D.  
1064 (2013). Accurate human microsatellite genotypes from high-throughput  
1065 resequencing data using informed error profiles. *Nucleic Acids Res* 41, e32.
- 1066 16. Cao, M.D., Tasker, E., Willadsen, K., Imelfort, M., Vishwanathan, S.,  
1067 Sureshkumar, S., Balasubramanian, S., and Boden, M. (2013). Inferring short

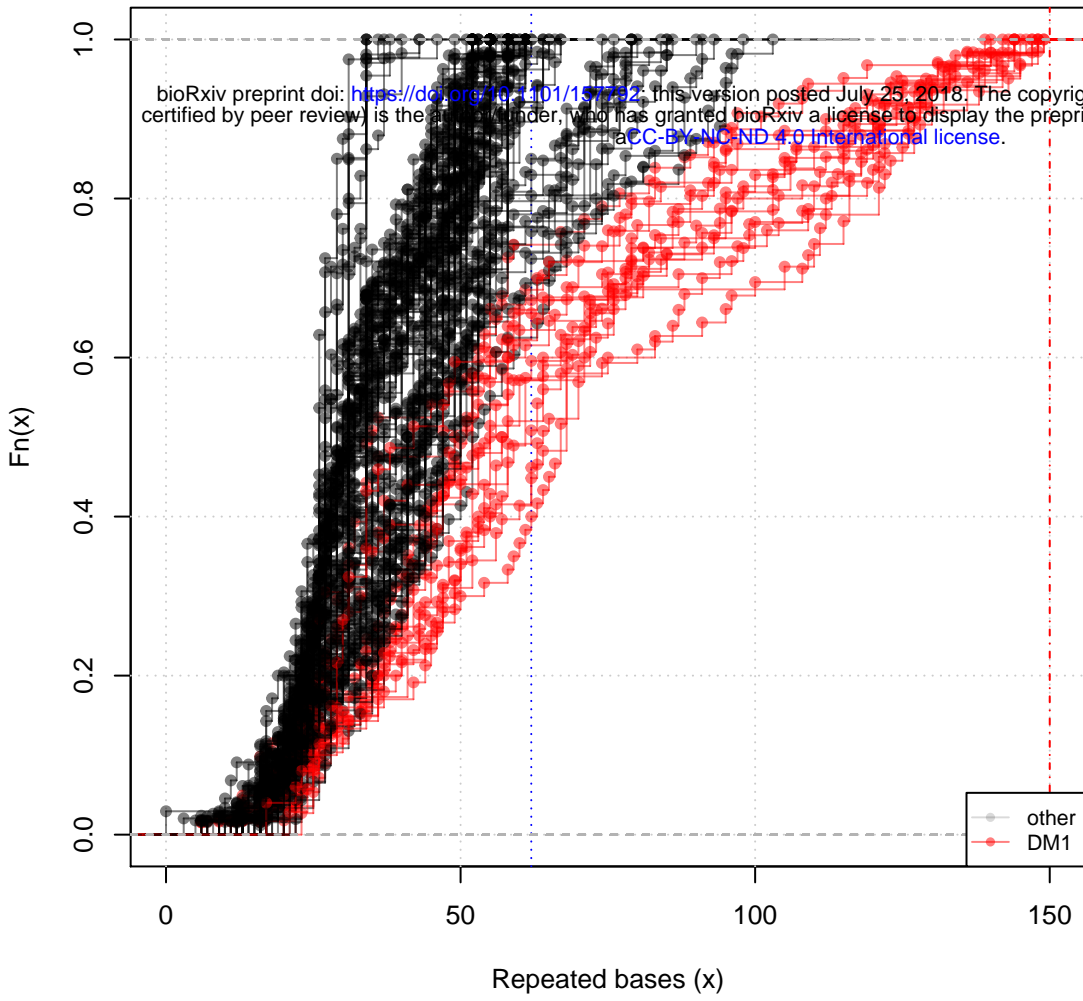


- 1068 tandem repeat variation from paired-end short reads. *Nucleic Acids Res* 42, p.  
1069 e16.
- 1070 17. Dolzhenko, E., van Vugt, J.J.F.A., Shaw, R.J., Bekritsky, M.A., van Blitterswijk,  
1071 M., Narzisi, G., Ajay, S.S., Rajan, V., Lajoie, B., Johnson, N.H., et al. (2017).  
1072 Detection of long repeat expansions from PCR-free whole-genome sequence  
1073 data. *Genome Res* 27, 1895-1903.
- 1074 18. Dashnow, H., Lek, M., Phipson, B., Halman, A., Davis, M., Lamont, P., Laing,  
1075 N., MacArthur, D., and Oshlack, A. (2017). STRetch: detecting and  
1076 discovering pathogenic short tandem repeats expansions. *bioRxiv*.
- 1077 19. Tang, H., Kirkness, E.F., Lippert, C., Biggs, W.H., Fabani, M., Guzman, E.,  
1078 Ramakrishnan, S., Lavrenko, V., Kakaradov, B., Hou, C., et al. (2017).  
1079 Profiling of Short-Tandem-Repeat Disease Alleles in 12,632 Human Whole  
1080 Genomes. *Am J Hum Genet* 101, 700-715.
- 1081 20. Bahlo, M., Bennett, M.F., Degorski, P., Tankard, R.M., Delatycki, M.B., and  
1082 Lockhart, P.J. (2018). Recent advances in the detection of repeat expansions  
1083 with short-read next-generation sequencing. *F1000Res* 7, 736.
- 1084 21. Langmead, B., Trapnell, C., Pop, M., Salzberg, S., Langmead, B., Trapnell, C.,  
1085 Pop, M., and Salzberg, S. (2009). Ultrafast and memory-efficient alignment of  
1086 short DNA sequences to the human genome. In *Genome Biol.* p R25.
- 1087 22. Huang, W., Li, L., Myers, J.R., and Marth, G.T. (2012). ART: a next-generation  
1088 sequencing read simulator. *Bioinformatics* 28, 593-594.
- 1089 23. Benjamini, Y., and Speed, T.P. (2012). Summarizing and correcting the GC  
1090 content bias in high-throughput sequencing. *Nucleic Acids Res* 40, e72.
- 1091 24. Rehm, H.L. (2017). Evolving health care through personal genomics. *Nat Rev*  
1092 *Genet* 18, 259-267.

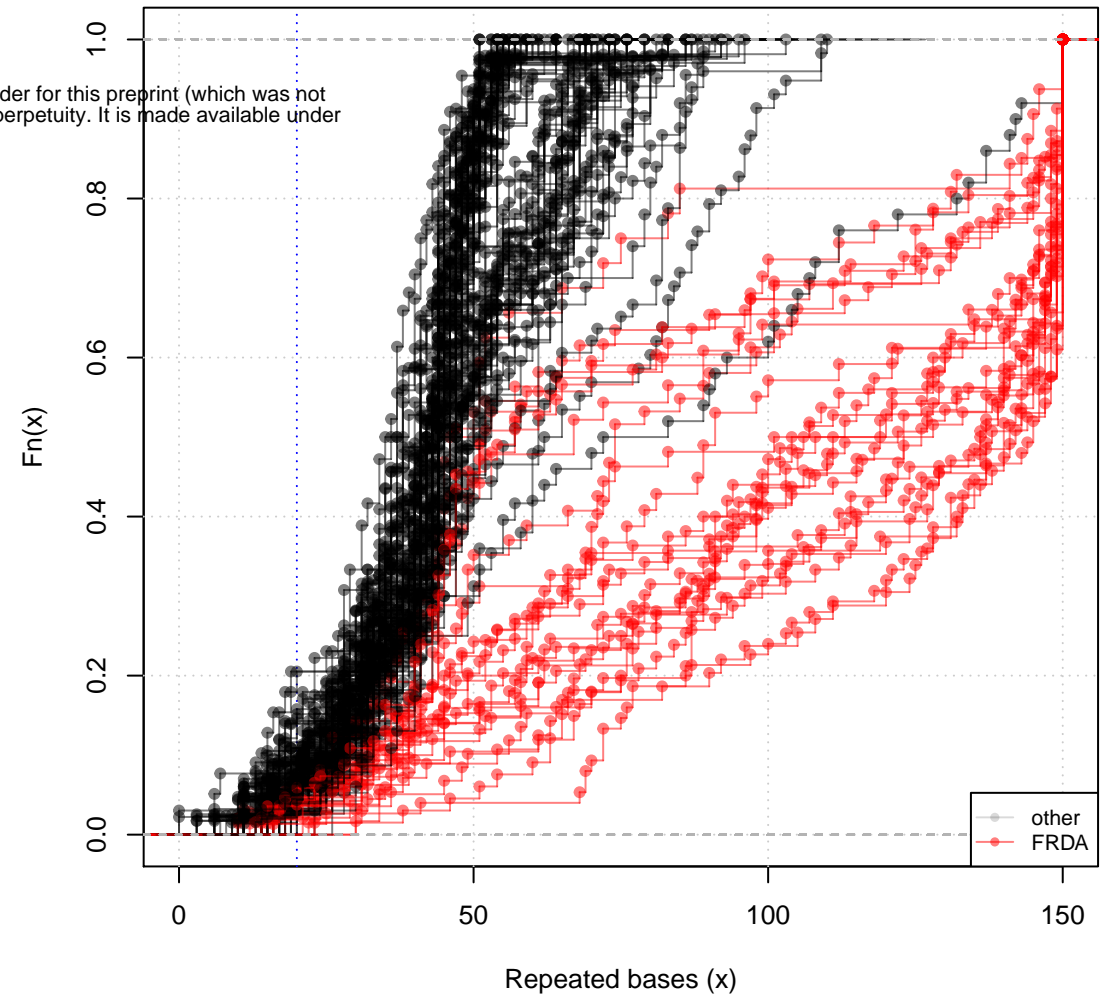
- 1093 25. Ashley, E.A. (2015). The precision medicine initiative: a new national effort.  
1094 JAMA 313, 2119-2120.
- 1095 26. Nemeth, A.H., Kwasniewska, A.C., Lise, S., Parolin Schneckenberg, R., Becker,  
1096 E.B.E., Bera, K.D., Shanks, M.E., Gregory, L., Buck, D., Zameel Cader, M.,  
1097 et al. (2013). Next generation sequencing for molecular diagnosis of  
1098 neurological disorders using ataxias as a model. In Brain. (
- 1099 27. Mousavi, N., Shleizer-Burko, S., and Gymrek, M. (2018). Profiling the genome-  
1100 wide landscape of tandem repeat expansions. bioRxiv.
- 1101 28. Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A.R., Bender, D.,  
1102 Maller, J., Sklar, P., de Bakker, P.I.W., Daly, M.J., et al. (2007). PLINK: a  
1103 tool set for whole-genome association and population-based linkage analyses.  
1104 Am J Hum Genet 81, 559-575.
- 1105 29. Pedersen, B.S., and Quinlan, A.R. (2017). Who's Who? Detecting and Resolving  
1106 Sample Anomalies in Human DNA Sequencing Studies with Peddy. Am J  
1107 Hum Genet 100, 406-413.
- 1108 30. Hunter, J., Rivero-Arias, O., Angelov, A., Kim, E., Fotheringham, I., and Leal, J.  
1109 (2014). Epidemiology of fragile X syndrome: a systematic review and meta-  
1110 analysis. Am J Med Genet 164A, 1648-1658.
- 1111 31. Ligsay, A., and Hagerman, R.J. (2016). Review of targeted treatments in fragile X  
1112 syndrome. Intractable Rare Dis Res 5, 158-167.
- 1113 32. Hagerman, R.J., and Polussa, J. (2015). Treatment of the psychiatric problems  
1114 associated with fragile X syndrome. Curr Opin Psychiatry 28, 107-112.
- 1115 33. Hagerman, R., and Hagerman, P. (2013). Advances in clinical and molecular  
1116 understanding of the FMR1 premutation and fragile X-associated  
1117 tremor/ataxia syndrome. The Lancet Neurology 12, 786-798.

- 1118 34. DePristo, M.A., Banks, E., Poplin, R., Garimella, K.V., Maguire, J.R., Hartl, C.,  
1119 Philippakis, A.A., del Angel, G., Rivas, M.A., Hanna, M., et al. (2011). A  
1120 framework for variation discovery and genotyping using next-generation DNA  
1121 sequencing data. *Nat Genet* 43, 491-498.
- 1122 35. Benson, G. (1999). Tandem repeats finder: a program to analyze DNA sequences.  
1123 In *Nucleic Acids Res.* pp 573-580.
- 1124 36. McMurray, C.T. (2010). Mechanisms of trinucleotide repeat instability during  
1125 human development. *Nat Rev Genet* 11, 786-799.
- 1126
- 1127

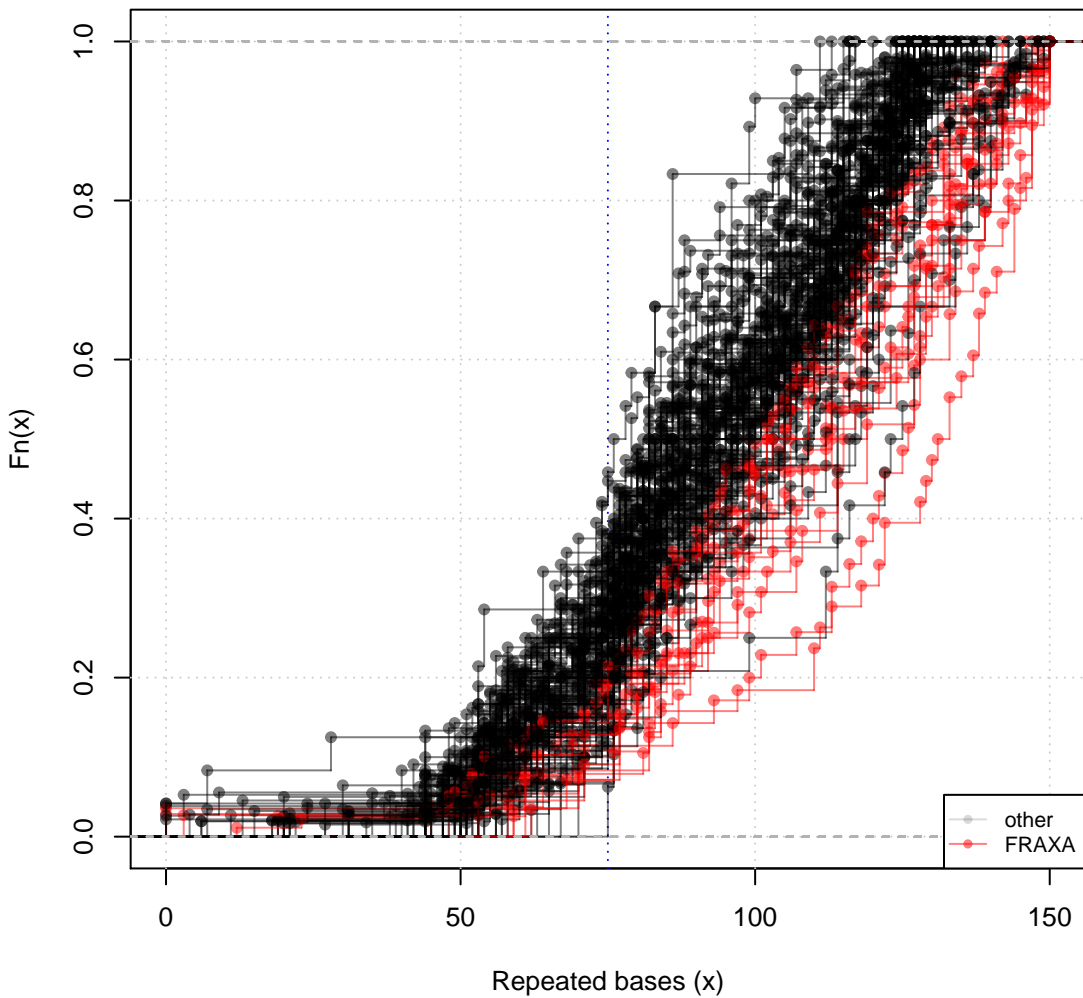
DM1 (3'UTR CTG) norm: 20 (62bp) , exp: 50 (150bp) score ECDF



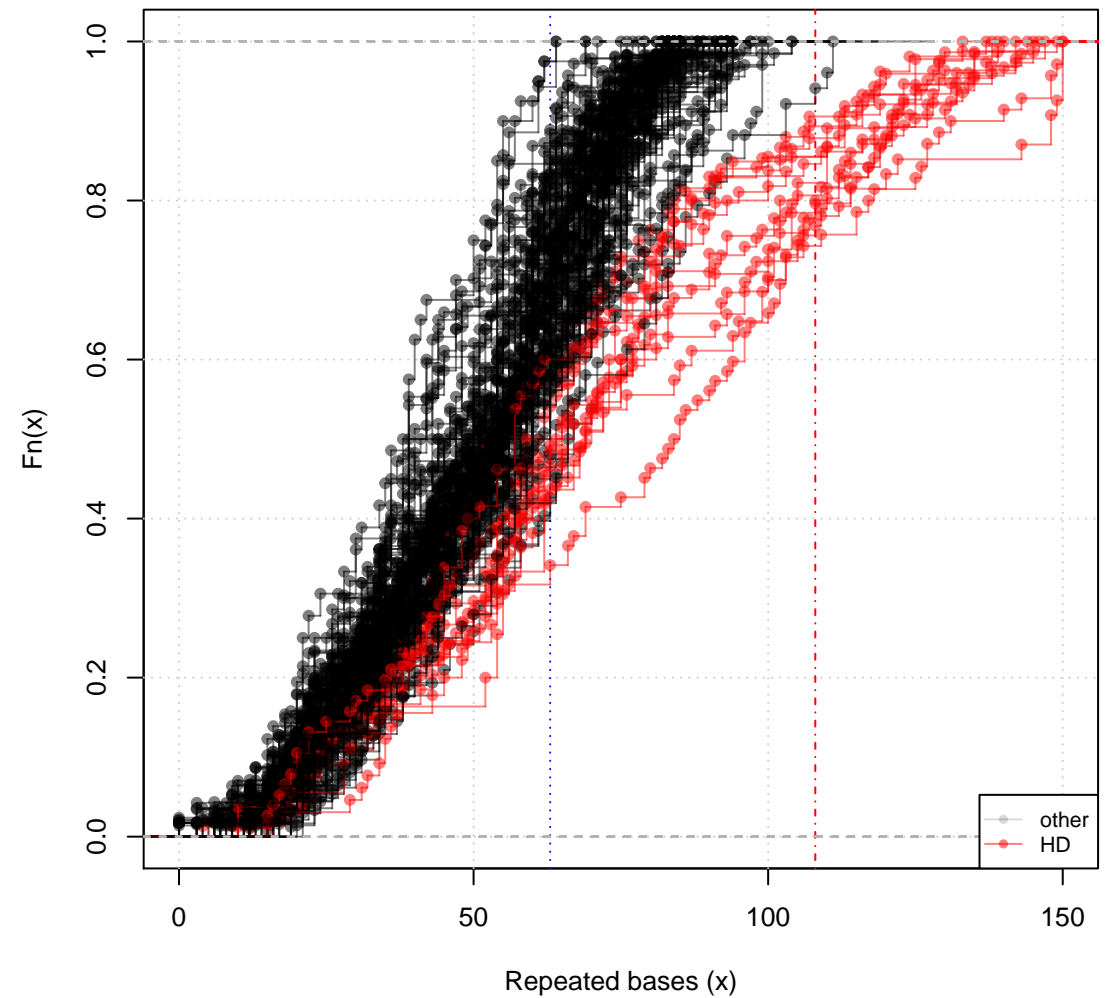
FRDA (intron\_1 GAA) norm: 6 (20bp) , exp: 200 (600bp) score ECDF



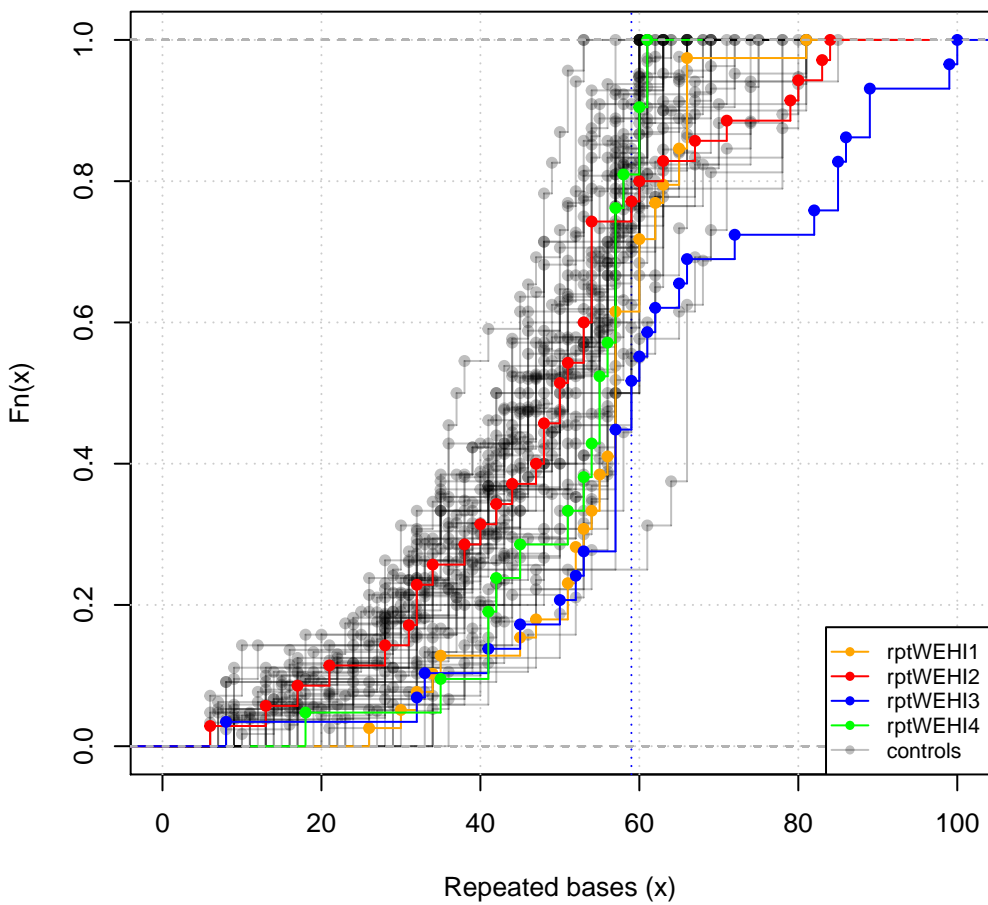
FRAXA (5'UTR CGG) norm: 25 (75bp) , exp: 200 (600bp) score ECDF



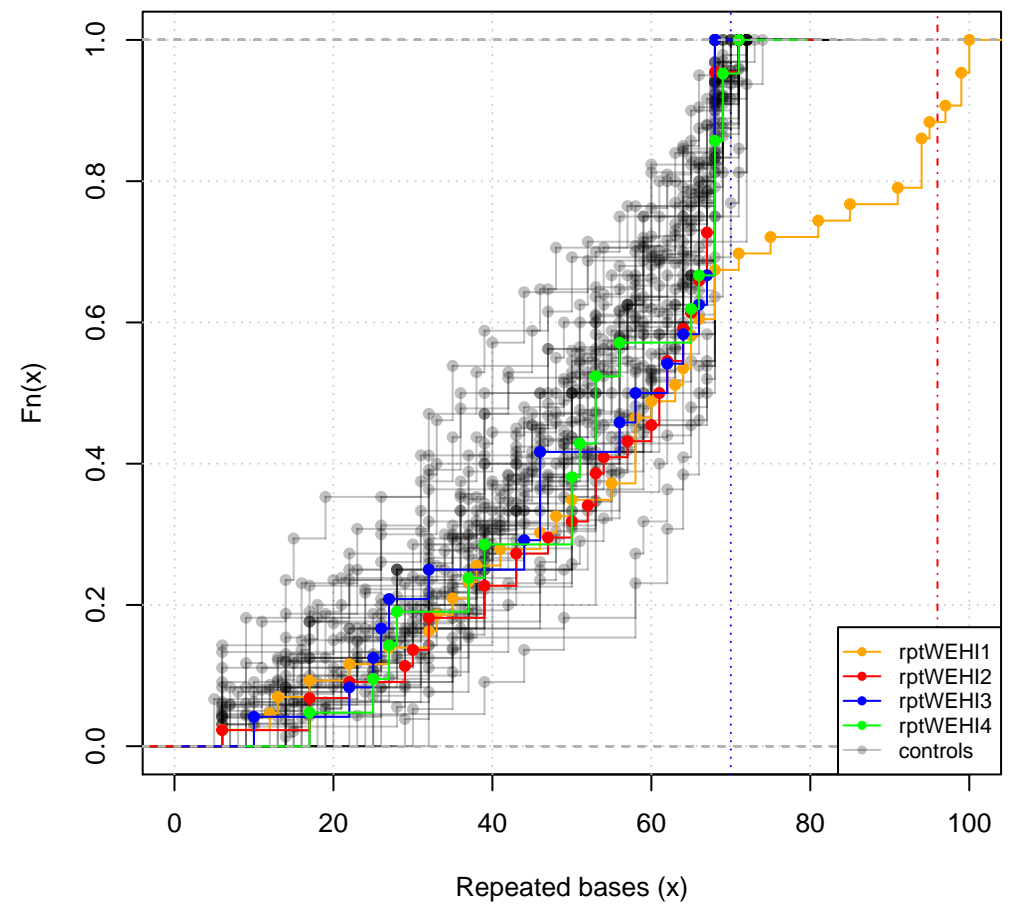
HD (coding CAG) norm: 21 (64bp) , exp: 36 (108bp) score ECDF



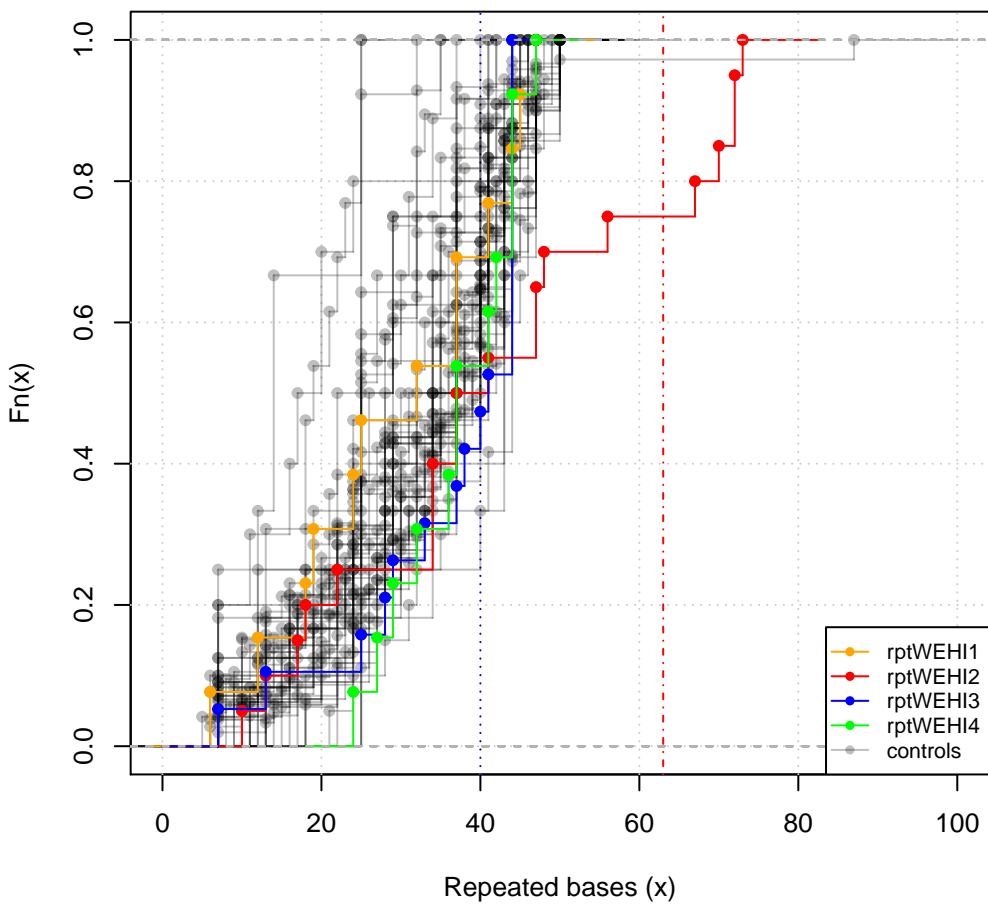
**HD (coding CAG) norm: 19 (59bp) , exp: 36 (108bp) score ECDF**



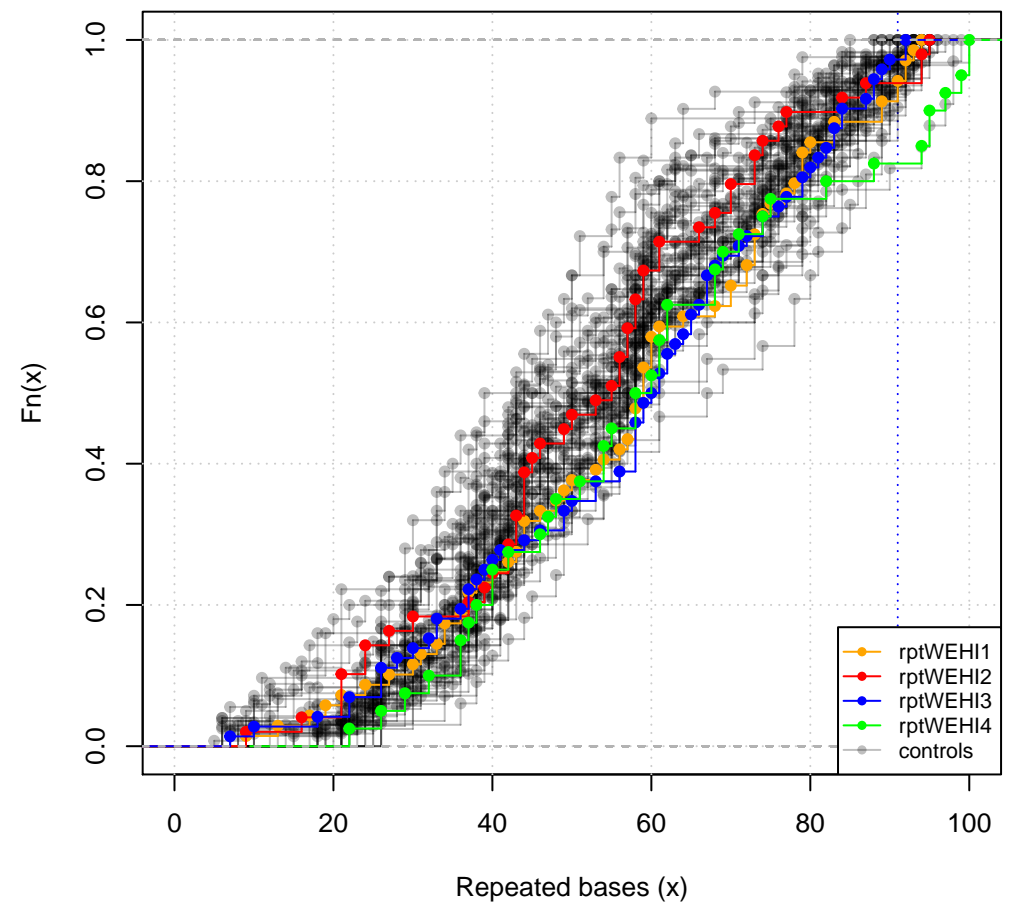
**SCA2 (coding CAG) norm: 23 (70bp) , exp: 32 (96bp) score ECDF**



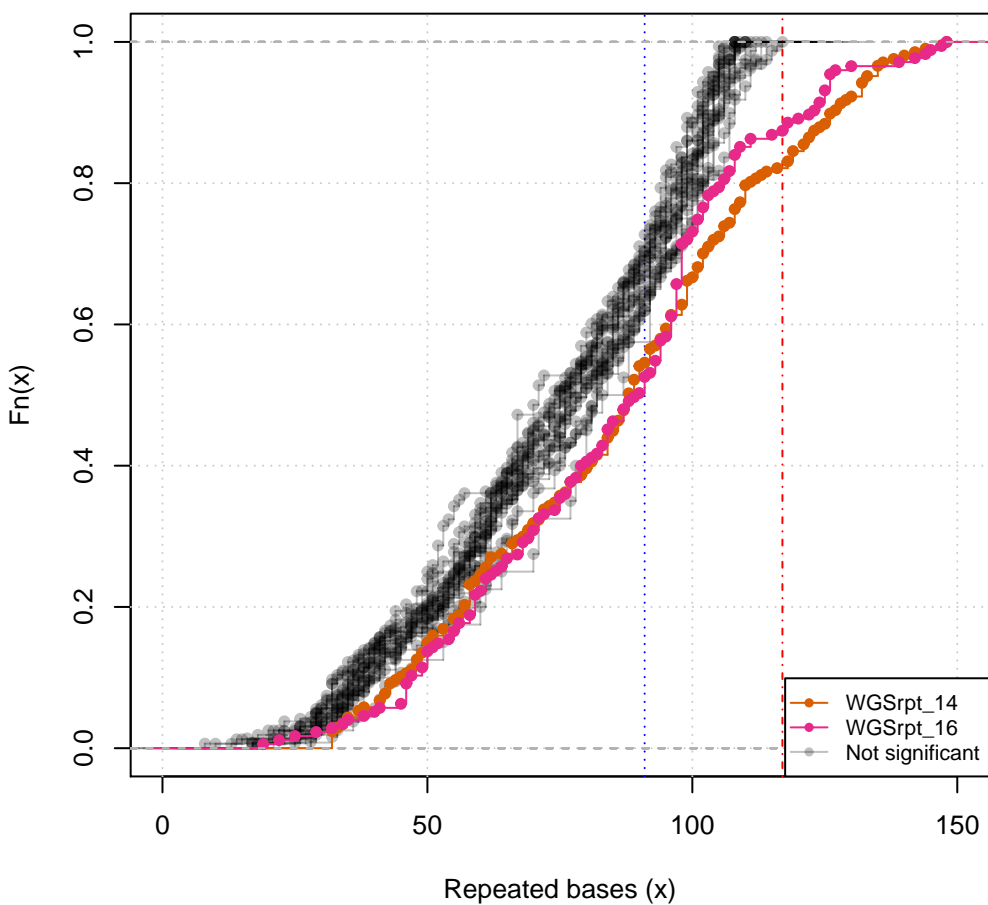
**SCA6 (coding CAG) norm: 13 (40bp) , exp: 21 (63bp) score ECDF**



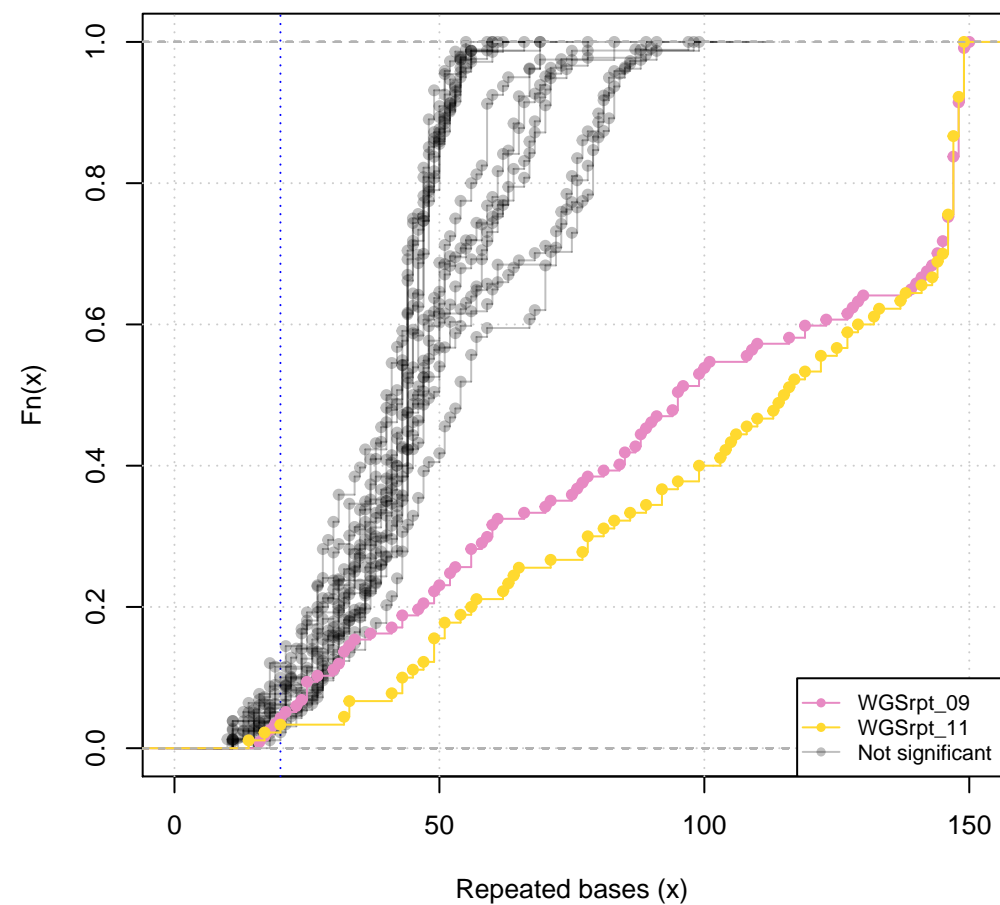
**SCA1 (coding CAG) norm: 30 (91bp) , exp: 39 (117bp) score ECDF**



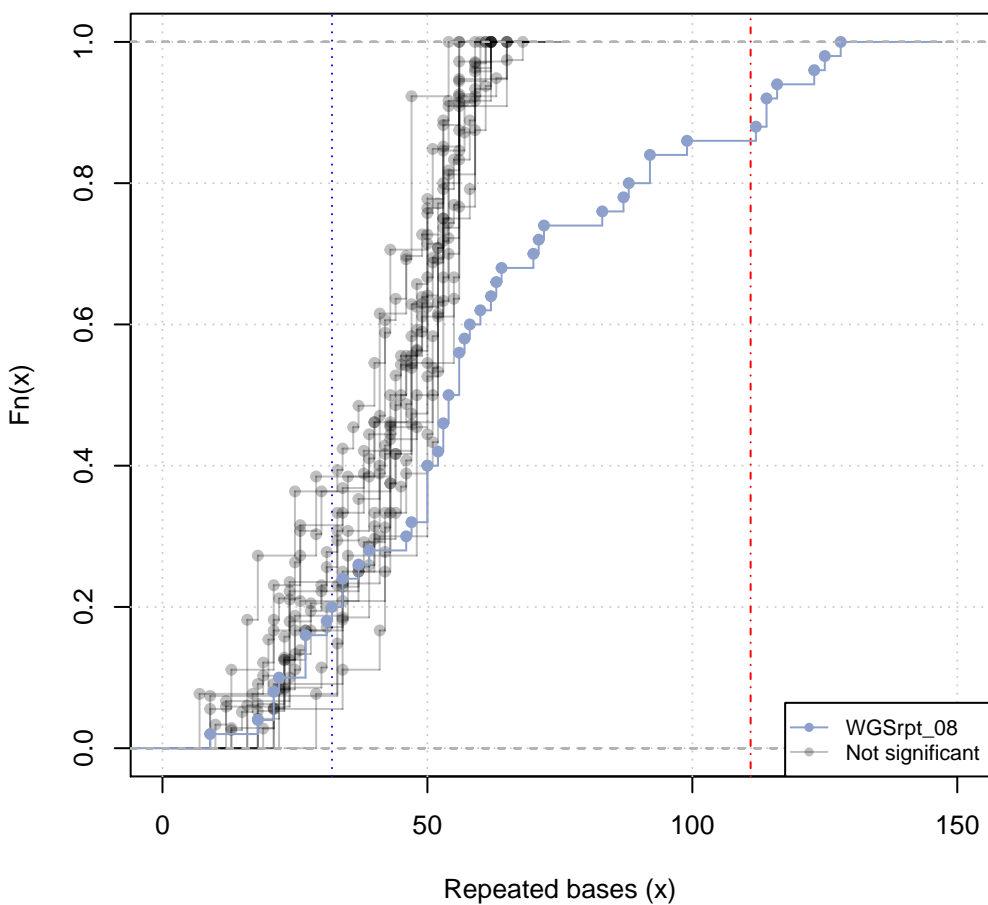
**SCA1 (coding CAG) norm: 30 (91bp) , exp: 39 (117bp) score ECDF**



**FRDA (intron\_1 GAA) norm: 6 (20bp) , exp: 200 (600bp) score ECDF**



**SCA7 (coding CAG) norm: 10 (32bp) , exp: 37 (111bp) score ECDF**



**DM1 (3'UTR CTG) norm: 20 (62bp) , exp: 50 (150bp) score ECDF**

