



## 23 **Abstract**

24           Genetic variation can cause significant differences in gene expression among  
25 individuals. Although quantitative genetic mapping techniques provide ways to identify  
26 genome-wide regulatory loci, they almost entirely focus on single nucleotide variants  
27 (SNVs). Short tandem repeats (STRs) represent a large source of genetic variation with  
28 potential regulatory effects. Here, we leverage the recently generated expression and  
29 STR variation data among wild *Caenorhabditis elegans* strains to conduct a genome-  
30 wide analysis of how STRs affect gene expression variation. We identify thousands of  
31 expression STRs (eSTRs) showing regulatory effects and demonstrate that they explain  
32 missing heritability beyond SNV-based expression quantitative trait loci. We illustrate  
33 specific regulatory mechanisms such as how eSTRs affect splicing sites and alternative  
34 splicing efficiency. We also show that differential expression of antioxidant genes might  
35 affect STR variation systematically. Overall, we reveal the interplay between STRs and  
36 gene expression variation in a tractable model system to ultimately associate STR  
37 variation with differences in complex traits.

## 38 Introduction

39 Genetic variation can cause significant differences in gene expression among  
40 individuals. Mutations in regulatory elements, such as promoters and enhancers, might  
41 only affect the expression of single genes, whereas mutations altering structures and  
42 abundances of diffusible factors, such as transcription factors (TFs) and chromatin  
43 cofactors, might affect the expression of multiple genes across the genome. Quantitative  
44 genetic mapping techniques, including both linkage and genome-wide association  
45 (GWA) mapping studies, enable the identification of genome-wide variants that influence  
46 gene expression and other complex traits. A genomic locus that contains alleles showing  
47 significant association with mRNA expression variation is called an expression  
48 quantitative trait locus (eQTL)<sup>1-5</sup>. Although thousands of eQTL have been detected in  
49 different organisms, associated genetic variants are mostly limited to single nucleotide  
50 variants (SNVs) and short insertions or deletions (indels)<sup>1-11</sup>. Emerging studies  
51 successfully linked gene expression variation to other types of DNA variants, such as  
52 short tandem repeats (STRs) and structural variants<sup>12-19</sup>.

53 STRs are repetitive elements consisting of 1-6 bp DNA sequence motifs<sup>17,20</sup>.  
54 Compared to SNVs and short indels, STR mutations show 1) orders of magnitude higher  
55 mutation rates<sup>20-23</sup>, 2) higher incidence of insertions or deletions, mostly in the number  
56 of repeats<sup>24,25</sup>, 3) more multiallelic sites<sup>26</sup>, and 4) more *de novo* mutations<sup>20,26</sup>. Dozens of  
57 human diseases have been associated with STR mutations<sup>24</sup>. Various effects of STR  
58 variation on regulation of gene expression have also been suggested from both *in vitro*  
59 and *in vivo* studies across a wide range of taxa<sup>27-33</sup>. However, these STRs only  
60 represented a small fraction of STRs in genomes. To our best knowledge, systematic  
61 evaluation of genome-wide associations between STR variation and gene expression  
62 variation have only been applied in humans<sup>12,17,34</sup> and *Arabidopsis thaliana*<sup>16,19</sup>, in part  
63 because of the difficulties in accurately genotyping STRs throughout the genome in large  
64 scales<sup>35</sup>.

65 We have recently studied the natural variation in gene expression<sup>5</sup> and STRs<sup>36</sup>  
66 across wild strains of the nematode *Caenorhabditis elegans*. We collected reliable  
67 expression measurements for 25,849 transcripts of 16,094 genes in 207 *C. elegans*

68 strains using bulk mRNA sequencing and identified 6,545 eQTL underlying expression  
69 variation of 5,291 transcripts of 4,520 genes using GWA mappings<sup>5</sup>. We characterized  
70 9,691 polymorphic STRs (pSTRs) with motif lengths of 1-6 bp across the species,  
71 including the 207 strains above, using high-throughput genome sequencing data<sup>36</sup> and  
72 a bioinformatic tool previously demonstrated to be reliable for large-scale profiling of  
73 STRs<sup>23,35</sup>.

74 In this work, we leveraged the recently generated expression<sup>5</sup> and STR<sup>36</sup> data from  
75 207 wild *C. elegans* strains to conduct a genome-wide scan of how STRs affect gene  
76 expression variation. We identified 3,118 and 1,857 expression STRs (eSTRs) that were  
77 associated with expression of nearby and remote genes, respectively. We found that  
78 eSTRs might help explain missing heritability in SNV-based eQTL studies for both local  
79 and distant eQTL. We also explored specific mechanisms of eSTRs and illustrated how  
80 local eSTRs might have influenced alternative splicing sites to cause differential  
81 transcript usage. We showed that expression of several genes in the same pathway  
82 might be altered because of a distant eSTR in a gene upstream. We also found evidence  
83 that expression variation in an antioxidant gene, *ctl-1*, might underlie STR variation  
84 across wild *C. elegans* strains. We further determined the positive relationship between  
85 endogenous oxidative stress and STR insertions/deletions using three mutation  
86 accumulation line panels. Our results demonstrate the systemic influences of eSTRs on  
87 gene expression and the potential effects of expression variation in antioxidant genes on  
88 STR mutations in *C. elegans*. We reveal the interplay between STRs and gene expression  
89 variation and provide publicly available frameworks to associate STRs with variation in  
90 gene expression and other complex traits in future studies.

91

## 92 **Results**

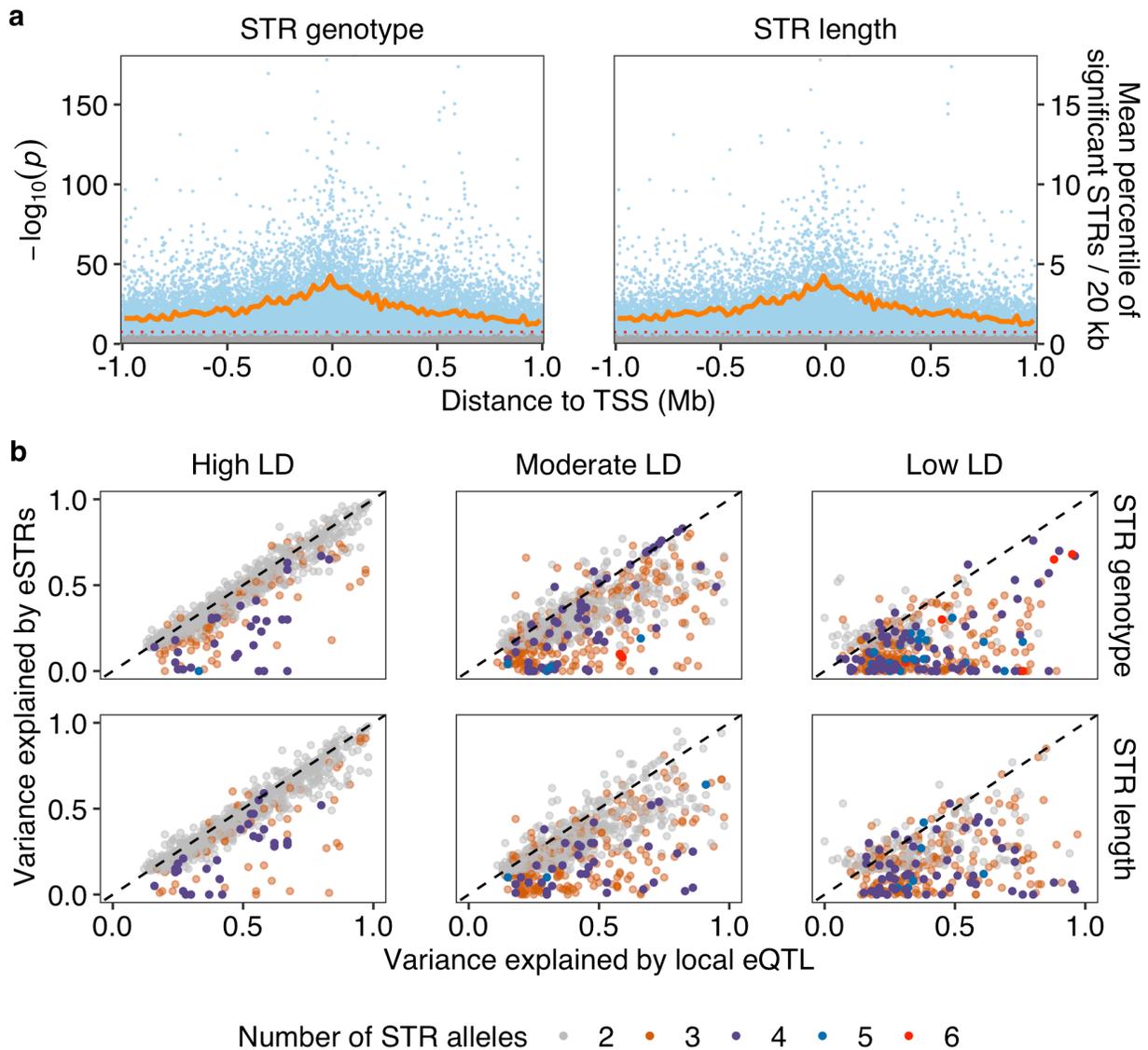
### 93 **Variation in STRs regulates expression in nearby genes**

94 We obtained expression data of 25,849 transcripts<sup>5</sup> of 16,094 genes and 9,691  
95 pSTRs<sup>36</sup> across 207 wild *C. elegans* strains. We investigated the effects of pSTRs on

96 transcript expression of nearby genes using a likelihood-ratio test (LRT) to evaluate the  
97 association between STR variation and transcript expression variation for all pSTRs  
98 within 2 Mb surrounding each transcript and with at least two common alleles (allele  
99 frequency > 0.5). We applied the LRT using both pSTR genotypes and lengths treating  
100 them as factorial variables (See Methods). In total, using STR genotypes, 1,555,828 tests  
101 were performed to test the effect of 3,335 pSTRs on the expression variation of 25,849  
102 transcripts, each of which was tested for a median of 59 STRs (ranging from one to 141)  
103 (Fig. 1a). Using STR lengths, 1,227,485 tests were performed for the effect of 2,607  
104 pSTRs on the expression variation of 25,847 transcripts, each of which was tested for a  
105 median of 47 STRs (ranging from one to 119) (Fig. 1a). For each test, we also performed  
106 another test using permuted STR genotypes or lengths. We identified local eSTRs with  
107 LRT values that passed the Bonferroni threshold ( $3.2E-8$  and  $4.1E-8$  for STR genotypes  
108 and lengths, respectively) and found 3,082 eSTRs for 2,888 transcripts by STR  
109 genotypes and 2,391 eSTRs for 2,791 transcripts by STR lengths, including 2,355 eSTRs  
110 for 2,695 transcripts by both STR genotypes or lengths (Fig. 1a, Supplementary Data 1).  
111 Each transcript had a median of nine eSTRs (ranging from one to 77) and six eSTRs  
112 (ranging from one to 65) by STR genotypes and lengths, respectively. None of the tests  
113 using permuted STRs passed the Bonferroni thresholds (Fig. 1a, Supplementary Data 1).  
114 As expected, we observed that STRs in close proximity to or within a transcript were  
115 more likely to pass the significance threshold than STRs far away from the transcript  
116 (Fig. 1a), indicating a close relationship between STRs and gene expression.

117 In our recent eQTL study<sup>5</sup>, we classified eQTL into local eQTL (located close to  
118 the genes that they influence) and distant eQTL (located farther away from the genes  
119 that they influence). Among the 3,185 transcripts with local eQTL, 2,477 were also found  
120 with eSTRs (Enrichment tested by one-sided Fisher's Exact Test, with  $p = 2.2E-16$ ). To  
121 compare the effects of eQTL and eSTRs in gene regulation, we compared the expression  
122 variance explained by eQTL and the most significant eSTR for each transcript and the  
123 LD between them (Fig. 1b). Most eQTL-eSTR pairs (48%) with high LD ( $r^2 \geq 0.7$ ) explained  
124 similar levels of expression variance (Fig. 1b), suggesting that these eSTRs might be  
125 detected because of the high LD to eQTL or *vice versa*. Among eQTL-eSTR pairs with

126 moderate LD ( $0.3 \leq r^2 < 0.7$ , 35%) or low LD ( $r^2 < 0.3$ , 17%), most eQTL explained more  
127 variance than eSTRs (Fig. 1b), suggesting these eSTRs, in particular multiallelic eSTRs,  
128 might be independent from the eQTL. Although eQTL identified using single-marker  
129 based GWA mappings explained a fraction of the variance in gene expression, eSTRs  
130 might help explain some missing heritability<sup>37</sup>.  
131



132

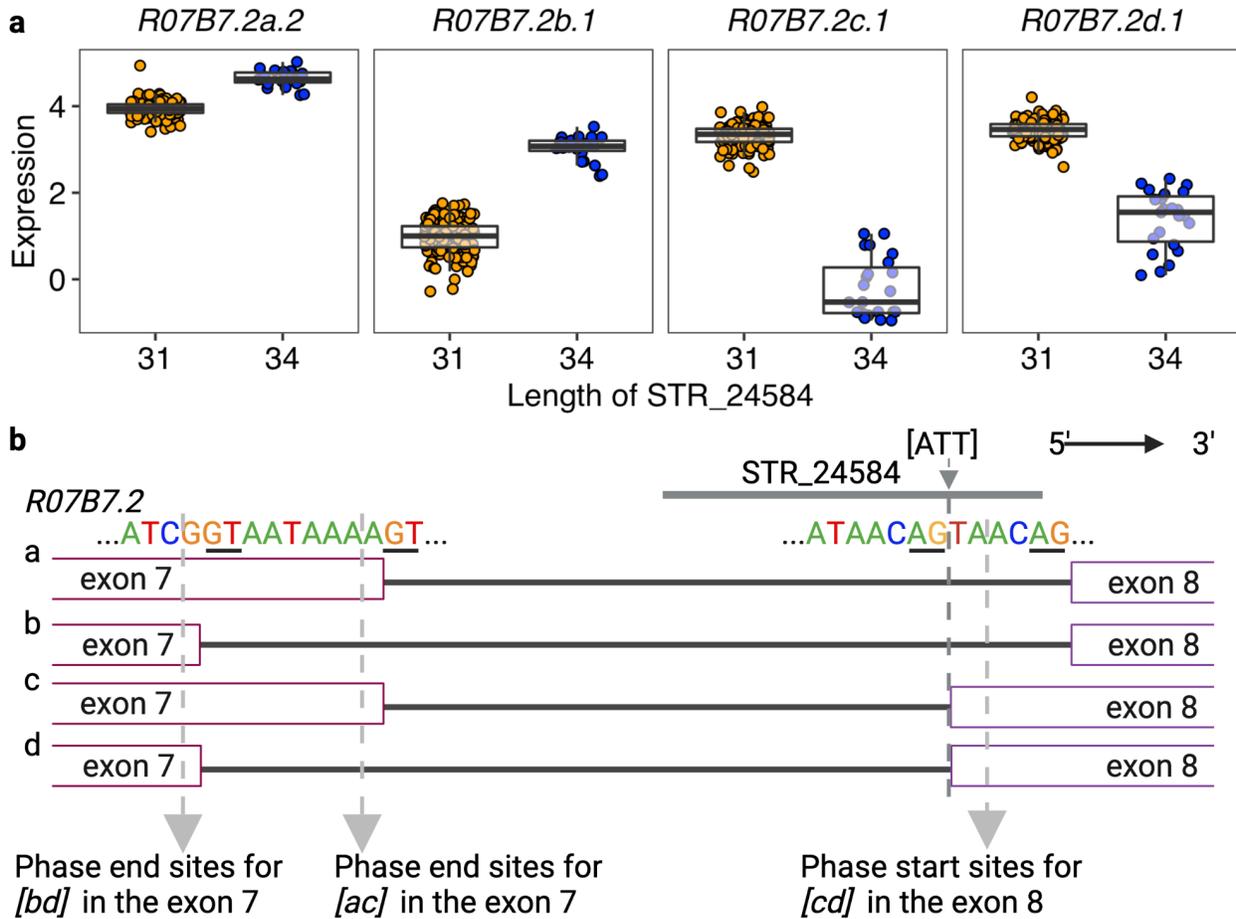
133 **Fig. 1: Expression STRs identified using Likelihood-Ratio Tests.**

134 **a** Identification of expression STRs (eSTRs) using Likelihood-Ratio Tests (LRT) on full  
135 (including STR variation as a variable) and reduced (excluding STR variation as a variable)  
136 models. The effects of STR variation in genotype (left panel) or length (right panel) were  
137 analyzed separately as factorial variables. Each dot represents a test between STR and  
138 transcript expression variation and is plotted with the distance of the STR to the  
139 transcription start site (TSS) of the transcript (x-axis) against its  $-\log_{10}(p)$  value (y-axis  
140 on the left). Blue and gray dots represent tests using real and permuted data of STR  
141 variation, respectively. The red dotted horizontal lines represent Bonferroni thresholds.  
142 The dark orange lines represent the mean percentage of significant tests (real data)  
143 above the Bonferroni thresholds in each 20 kb bin (y-axis on the right). **b** The variance  
144 explained (VE) by local eQTL that were identified using GWA mapping experiments<sup>5</sup> was  
145 plotted against the VE for the most significant eSTRs. Dots are colored by the number  
146 of STR alleles used in eSTR VE calculation. LD ( $r^2$ ) between eQTL and eSTRs were used  
147 to separate panels on the x-axis, with high LD ( $r^2 \geq 0.7$ ), moderate LD ( $0.3 \leq r^2 < 0.7$ ), and  
148 low LD ( $r^2 < 0.3$ ). The dashed lines on the diagonal are shown as visual guides to  
149 represent  $VE_{eQTL} = VE_{eSTRs}$ .  
150

151 **Insertion in a local eSTR affects transcript isoform usage**

152 We next focused on eSTRs that were in genomic features of their target  
153 transcripts and were outside of hyper-divergent regions<sup>38</sup>. We predicted the functional  
154 consequences<sup>39</sup> of these eSTRs and found a total of 13 eSTRs in 16 transcripts of 12  
155 genes that showed high-impact mutations, including missense mutations, in-frame  
156 insertions and deletions, start lost, stop gain, and mutations in splicing regions or  
157 acceptors. Another 17 eSTRs in 21 transcripts of 17 genes were predicted to affect  
158 5'UTRs and 3'UTRs. We identified two enriched motif sequences, ATTTTT and ATGTT,  
159 in these eSTRs by STR genotypes (one-sided Fisher exact test, Bonferroni-corrected  $p$   
160 = 0.04 and 6.8E-5, respectively) or STR lengths (one-sided Fisher exact test, Bonferroni-  
161 corrected  $p$  = 0.03 and 4.6E-5, respectively). Instead of finding multiple eSTRs, the two  
162 motif sequences only came from two eSTRs, STR\_13795 of (ATTTTT)<sub>5</sub> and STR\_24584  
163 of (ATGTT)<sub>6,2</sub>, each of which was associated with multiple transcripts of the same genes.  
164 In particular, STR\_24584 was predicted to have high-impact mutations in the splicing

165 regions of four transcripts of the gene, *R07B7.2*, and was associated with their  
166 expression variation (Fig. 2). Compared to strains with the reference allele, strains with a  
167 3-bp insertion showed significantly higher expression in the isoforms *R07B7.2[ab]* but  
168 significantly lower expression in the isoforms *R07B7.2[cd]* (Fig. 2a). More specifically,  
169 the insertion was located at the 3' splice site in the intron between exon 7 and exon 8 of  
170 *R07B7.2[ab]* and at the junction of the intron and exon 8 for *R07B7.2[cd]* (Fig. 2b). We  
171 speculated that at least two mechanisms might underlie the expression differences  
172 among the four transcripts caused by STR\_24584 variation. First, the insertion [ATT]  
173 changed the 3' splice site of *R07B7.2[ab]* from 5'-GTAACAG-3' to 5'-TTAACAG-3' (Fig.  
174 2b), which became closer to the conserved consensus sequence 5'-UUUUCAG-3' of the  
175 3' splice site in *C. elegans*<sup>40</sup>. Therefore, the insertion might promote splicing efficiency  
176 for *R07B7.2[ab]* in pre-mRNAs of *R07B7.2* and thus increase the expression of the two  
177 transcripts, which consequently would decrease the expression of *R07B7.2[cd]*. Second,  
178 the insertion could cause a frameshift and insertion in the coding regions of *R07B7.2[cd]*,  
179 which caused I474NL (ATA to AATTTA) and V471DL (GTA to GATTTA) in *R07B7.2[c]* and  
180 *R07B7.2[d]* (Fig. 2b), respectively. These mutations might increase mRNA degradation.  
181 Taken together, our results demonstrated the effects of STR variation on gene  
182 expression and provided examples for potential underlying mechanisms.



**Fig. 2: Expression STRs disrupting splicing.**

185 **a** Tukey box plots showing expression variation of four transcripts of the gene *R07B7.2*  
 186 between strains with different lengths of the STR\_24584. Each point corresponds to a  
 187 strain and is colored orange or blue for strains with the N2 reference allele or the  
 188 alternative allele, respectively. Box edges denote the 25th and 75th quantiles of the data;  
 189 and whiskers represent 1.5× the interquartile range. **b** Graphic illustration of sequences  
 190 in the splice site of four transcripts of the gene *R07b7.2[ab]* and the position of STR\_24584.  
 191 The dashed arrow in dark gray indicates the position of a 3-bp insertion in the  
 192 STR\_24584 and the splicing region of *R07b7.2[ab]*. The dashed arrow in light gray  
 193 indicates the phase start and end sites for different exons. Created using  
 194 BioRender.com.

195

## 196 **STR variation underlies distant eQTL hotspots**

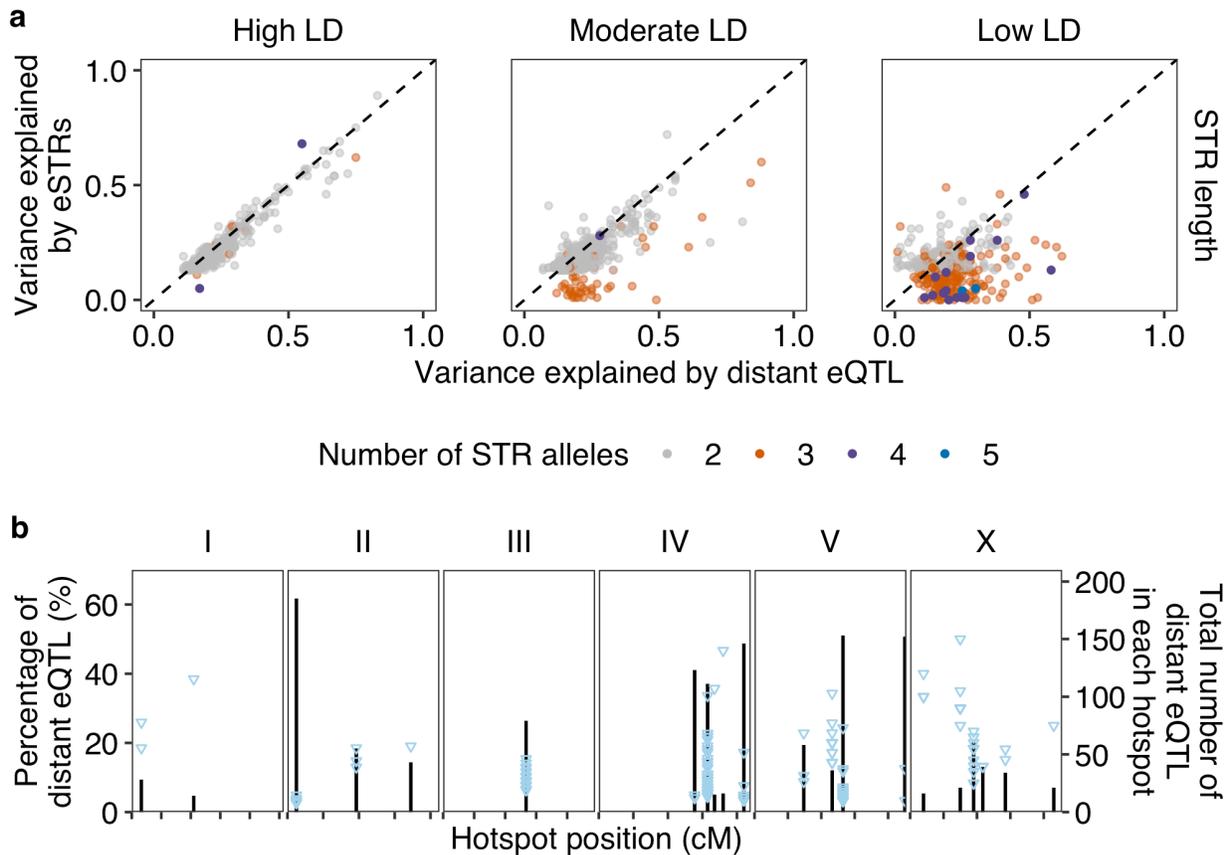
197 In addition to local eQTL, we also identified 3,360 distant eQTL for 2,553  
198 transcripts from 2,382 genes<sup>5</sup>. Genetic variants underlying distant eQTL might affect  
199 genes encoding diffusible factors like TFs to regulate genes across the genome. After  
200 the identification of local eSTRs, we identified distant eSTRs that affect remote genes.  
201 Instead of testing all pSTRs across the genome for each transcript, we selected pSTRs  
202 that are within 2 Mb surrounding the QTL regions of interest for all distant eQTL of each  
203 transcript. We used LRT tests (as above, also see Methods) to associate pSTR length  
204 variation with expression variation. In total, 353,694 tests were performed for the effects  
205 of 2,743 pSTRs on the expression variation of 2,553 transcripts, each of which was  
206 tested for a median of 104 STRs (ranging from one to 1,005). We used the Bonferroni  
207 threshold ( $1.4E-7$ ) to identify 1,857 distant eSTRs for 950 transcripts, with a median of  
208 three distant eSTRs (ranging from one to 127) (Supplementary Data 2). We also  
209 compared the expression variation explained by each distant eQTL and the most  
210 significant distant eSTR, and the LD between them. Different from local eQTL-eSTR pairs  
211 (Fig. 1b), most distant eQTL-eSTR pairs showed moderate (38%) or low (34%) LD,  
212 suggesting a more independent role of distant eSTRs in gene regulation (Fig. 3a). We  
213 have previously identified 46 distant eQTL hotspots that were enriched with distant  
214 eQTL<sup>5</sup> (Fig. 3b). Genetic variants in these hotspots were associated with expression  
215 variation in up to 184 transcripts<sup>5</sup>. Here, we found 229 common distant eSTRs that were  
216 associated with at least five distant eQTL in each hotspot (Fig. 3b). Common eSTRs  
217 might even underlie about half of all the distant eQTL in several hotspots (Fig. 3b).  
218 Altogether, these results suggested the complementary regulatory effects of distant  
219 eSTRs to distant eQTL and hotspots.

220 We next investigated whether any of the common distant eSTRs were in genes  
221 encoding TFs or chromatin cofactors. We found nine TF genes and one chromatin  
222 cofactors genes that harbor common distant eSTRs (Supplementary Data 3). For  
223 example, STR\_12763 was a common eSTR for seven distant eQTL in the hotspot ranging  
224 from 26 to 27.5 cM on chromosome III (Supplementary Data 3). STR\_12763 is in the  
225 3'UTR of the TF gene, *atf-7<sup>41</sup>*, and overlaps with the binding sites of multiple miRNAs

226 (Supplementary Fig. 1). Variation in STR\_12763 could affect the targeting of *atf-7* mRNAs  
227 by miRNAs to alter expression of the six transcripts (genes). However, none of the ten  
228 common distant eSTRs were also identified as local eSTRs for the genes in which they  
229 are located. So, we investigated whether any other common eSTRs, although not in  
230 known regulatory genes, were also identified as local eSTRs.

231 We found ten common distant eSTRs that were also local eSTRs for seven genes  
232 (Supplementary Data 3). We previously mentioned STR\_13795 (ATTTTT)<sub>5</sub> as one of the  
233 two local eSTRs with enriched motif sequences. The variation of STR\_13795 was  
234 associated with two transcripts of the gene, *cls-2*. Strains with STR contraction by about  
235 three repeats (17 bp) in STR\_13795 showed significantly higher expression in both  
236 transcripts of *cls-2* than strains with the reference STR allele (Supplementary Fig. 2a).  
237 Because STR\_13795 was in the 3'UTR of *cls-2*, the 17-bp deletion associated with  
238 expression of *cls-2* might affect targeting by miRNAs<sup>42,43</sup>. STR\_13795 was also identified  
239 as a distant eSTR for another ten transcripts, including the gene *polq-1* (Supplementary  
240 Fig. 2b). STR\_13083 was identified as a local eSTR for *polq-1* and distant eSTRs for  
241 another nine transcripts, of which six had STR\_13795 as an eSTR (Supplementary Fig.  
242 2b, Supplementary Fig. 3). Most strains with length 30 and 13 in the STR\_13795 also  
243 have length 16 and 15, respectively in the STR\_13083 (Supplementary Table 1). Because  
244 STR\_13795 was also associated with *polq-1*, STR\_13795 was more likely to be the  
245 causal candidate than STR\_13083 to alter the expression of the six overlapped target  
246 transcripts. The significant association between STR\_13083 length variation and the  
247 expression variation of the six overlapped transcripts were identified because of the  
248 linkage between STR\_13083 and STR\_13795. The three transcripts that only had  
249 STR\_13083 as their distant eSTRs could also be associated with the length variation of  
250 STR\_13795, which was not tested for the three transcripts because it was too distant  
251 from the genes. Altogether, STR\_13795 might affect the expression of all the 13 remote  
252 transcripts and genes by altering the expression of *cls-2* (Supplementary Fig. 2,  
253 Supplementary Fig. 3b). We performed gene set enrichment analysis for the 13 genes  
254 on WormBase<sup>44</sup> and found significant enrichment in genes related to spindle and  
255 germline defectiveness (Supplementary Table 2). The conserved protein, CLASP/CLS-2,

256 is required for mitotic central spindle stability, oocyte meiotic spindle assembly,  
 257 chromosome segregation, and polar body extrusion in *C. elegans*<sup>45-49</sup>. To summarize,  
 258 variation in STR\_13795 might alter the expression of *cls-2*, which could further affect  
 259 other related genes in the spindle assembly pathways.



260

261 **Fig. 3: Expression STRs underlying distant eQTL hotspots.**

262 **a** The variance explained (VE) by distant eQTL that were identified by GWA mapping  
 263 experiments<sup>5</sup> was plotted against the VE by the most significant eSTRs. Dots are colored  
 264 by the number of STR alleles used in eSTR VE calculation. LD ( $r^2$ ) between eQTL and  
 265 eSTRs were used to separate panels on the x-axis, with high LD ( $r^2 \geq 0.7$ ), moderate LD  
 266 ( $0.3 \leq r^2 < 0.7$ ), and low LD ( $r^2 < 0.3$ ). The dashed lines on the diagonal are shown as visual  
 267 guides to represent  $VE_{eQTL} = VE_{eSTRs}$ . **b** The percentage of distant eQTL (y-axis on the left)  
 268 that were associated with eSTRs in each distant eQTL hotspot<sup>5</sup> across the genome (x-  
 269 axis) is shown. Each blue triangle represents a common eSTR. Black bar indicates the

270 total number of distant eQTL (y-axis on the right) in each hotspot. Tick marks on the x-  
271 axis denote every 10 cM.  
272

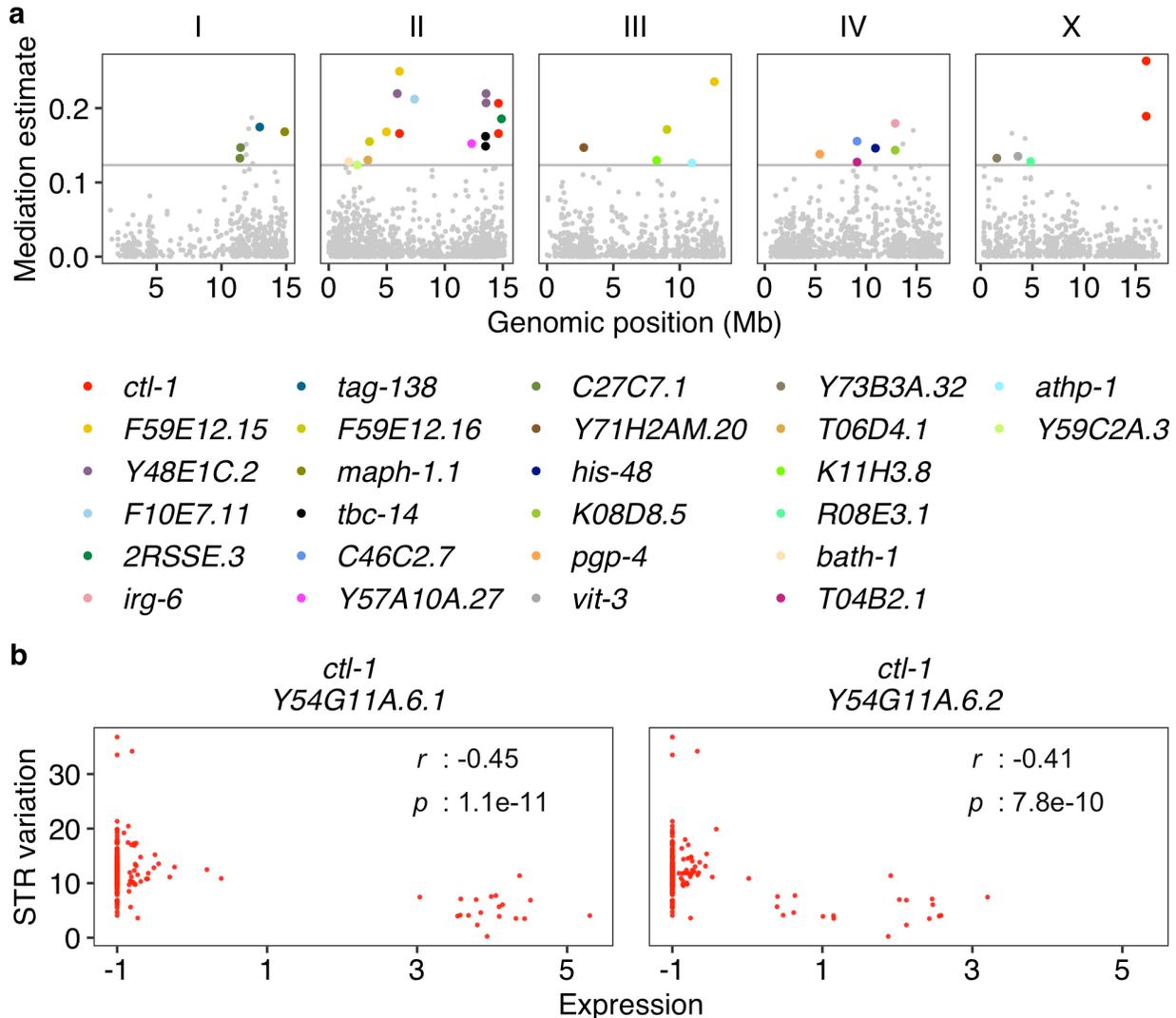
### 273 **Oxidative stress potentially drives STR mutations**

274 To explore the genome-wide influences of STRs on gene expression variation, we  
275 also wondered what factors might affect STR mutations and cause STR variation across  
276 *C. elegans*. DNA strand slippage during replication, DNA repair, and recombination  
277 processes can lead to STR mutations<sup>24</sup>. We reasoned that any genetic or environmental  
278 factors that are able to increase errors during these processes or decrease genome  
279 stability could increase STR mutation rates<sup>50,51</sup>. We hypothesized that, if variation in  
280 genetic factors that affect genomic stability exists, the amount of total STR variation  
281 could be used as a quantitative trait for a GWA mapping study. We recently also  
282 developed a pipeline of mediation analysis to link gene expression variation to  
283 quantitative traits<sup>5</sup>. Thus, we sought to examine potential genetic and mediating factors  
284 underlying STR mutation variation.

285 We first defined an STR variation trait by counting reference and alternative STR  
286 alleles for each of the 207 strains in the 9,691 pSTRs (See Methods) (Supplementary Fig.  
287 4a). Deletions are the predominant mutations in STR mutations across wild *C. elegans*  
288 strains (Supplementary Fig. 4a). We performed GWA mappings using two methods,  
289 LOCO and INBRED<sup>52</sup>, for this trait (see Methods). The INBRED method corrects more  
290 heavily for genetic stratification and many times decreases mapping power more than  
291 the LOCO method<sup>52-54</sup>. We detected six QTL with large QTL regions of interest on five of  
292 the six chromosomes using LOCO but no QTL using INBRED (Supplementary Fig. 4b,  
293 Supplementary Table 3). We next used mediation analysis to link expression differences  
294 with total STR mutation variation. Mediation analysis was performed for any transcripts  
295 with eQTL that overlap with the QTL regions of interest of the six QTL for STR variation.  
296 We identified 31 significant mediator transcripts of 26 genes (Fig. 4a). The top mediator  
297 gene, *ctl-1*, had two transcripts identified as significant mediators by multiple tests using  
298 different pairs of eQTL and QTL (Fig. 4a). We found moderate negative correlations

299 between the expression of the two *ctl-1* transcripts (*Y54G11A.6.1* and *Y54G11A.6.2*) and  
300 STR mutation variation (Fig. 4b), suggesting that the expression level of *ctl-1* might  
301 impact STR mutation variation. We regressed the STR variation trait by the expression  
302 of the transcript *Y54G11A.6.1* and performed GWA mappings. All the QTL mapped using  
303 the raw trait and LOCO disappeared in the mappings using the regressed trait  
304 (Supplementary Fig. 4c, Supplementary Table 3), supporting that the expression  
305 variation of *ctl-1* might affect STR mutation variation. We also identified a new QTL at  
306 the position 14,625,147 on chromosome II in both LOCO and INBRED methods  
307 (Supplementary Fig. 4c, Supplementary Table 3), suggesting that loci other than *ctl-1*  
308 might affect STR mutation variation as well.

309 The gene, *ctl-1*, encodes a cytosolic catalase in the detoxification pathway of  
310 reactive oxygen species (ROS)<sup>55</sup>. Elevated expression of *ctl-1* and other antioxidant  
311 related genes, which likely enhanced resistance to oxidative stresses, was associated  
312 with lifespan elongation in *C. elegans*<sup>56,57</sup>. Oxidative damage can alter DNA secondary  
313 structure, affect genome stability and replication, and cause mutations<sup>58</sup>. Therefore, it is  
314 possible that the group of strains showing high levels of *ctl-1* expression managed to  
315 reduce STR mutations caused by oxidative damage over time and have lower levels of  
316 total STR mutations across the species (Fig. 4b). We have previously detected five (one  
317 local and four distant) and six (one local and five distant) eQTL for expression variation  
318 of the two transcripts of *ctl-1*, *Y54G11A.6.1* and *Y54G11A.6.2*, respectively<sup>5</sup>. Among the  
319 5,291 transcripts with detected eQTL, 4,430 transcripts had a single eQTL detected and  
320 only 30 transcripts were found with equal or more than five eQTL<sup>5</sup>. These results suggest  
321 that the expression of *ctl-1* was highly controlled and might be critical for adaptation to  
322 oxidative stresses.



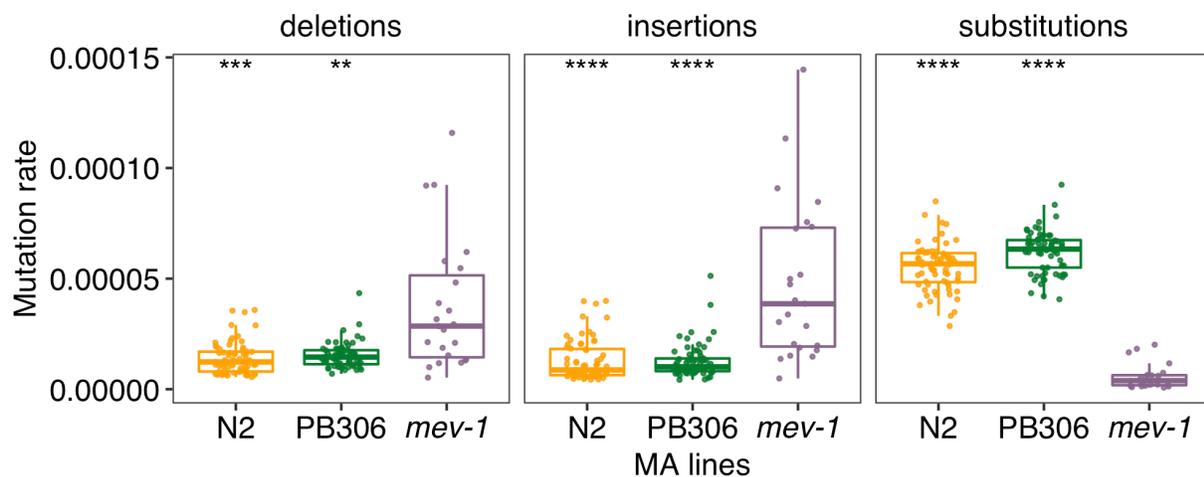
323

324 **Fig. 4: Mediation effects of *ctl-1* expression on STR variation.**

325 **a** Mediation estimates (y-axis) of transcript expression on STR variation are plotted  
 326 against the genomic position (x-axis) of the eQTL. The horizontal gray line represents the  
 327 99<sup>th</sup> percentile of the distribution of mediation estimates. Mediator transcripts with  
 328 adjusted  $p < 0.05$  and interpretable mediation estimate greater than the 99<sup>th</sup> percentile  
 329 estimates threshold are colored by their genes. Other tested mediator transcripts are  
 330 colored gray. **b** The correlation of expression (x-axis) of four mediator transcripts to STR  
 331 variation (y-axis) is shown. Each dot represents a strain and is colored by mediator genes  
 332 as in **a**. The coefficient  $r$  and the  $p$ -value for each correlation using the two-sided  
 333 Pearson's correlation tests are indicated in the top right.

334

335 We further examined potential relationships between oxidative stresses and STR  
336 mutations using three mutation accumulation (MA) line panels<sup>59-62</sup> that have undergone  
337 passage for many generations with minimal selection: 1) 67 MA lines that were derived  
338 from N2 and propagated for ~250 generations; 2) 23 MA lines that were derived from a  
339 mutant strain, *mev-1* (with a missense mutation introgressed into N2, resulting in  
340 elevated oxidative stress), and propagated for ~125 generations; and 3) 67 MA lines that  
341 were derived from PB306 (a wild strain) and propagated for ~250 generations. We  
342 obtained raw sequencing data for these 157 MA lines and their three ancestors and  
343 called STR variation using the same method that we used for wild *C. elegans* strains<sup>36</sup>  
344 (See Methods). We calculated mutation rates for three different mutations (deletions,  
345 insertions, and substitutions) between the ancestor and each derived MA line and  
346 compared mutation rates across the three MA lines (Fig. 5). We found that *mev-1* MA  
347 lines showed significantly higher mutation rates in deletions and insertions but  
348 significantly lower substitution rates than the other two MA lines (Fig. 5, Supplementary  
349 Table 4). The gene *mev-1* encodes a mitochondrial complex II SDHC<sup>63</sup>. The *mev-1*  
350 mutant was found to be highly sensitive to oxidative stress and showed reduced  
351 lifespan<sup>63</sup>. The high deletion and insertion rates in *mev-1* lines might be driven by their  
352 increased endogenous oxidative damage than the other two MA lines. Although the  
353 mutation rate of substitution was low in *mev-1* lines, deletions and insertions likely  
354 contributed most of the variation in STRs (Supplementary Fig. 4A).  
355



356

357 **Fig. 5: STR mutation rates in the MA lines.**

358 Comparison of STR mutation rates in deletions, insertions, and substitutions between  
359 the *mev-1* line (purple) and N2 (orange), and PB306 (green) lines, respectively. Box edges  
360 denote the 25th and 75th quantiles of the data; and whiskers represent 1.5× the  
361 interquartile range. Statistical significance of difference comparisons (Supplementary  
362 Table 4) was calculated using the two-sided Wilcoxon test and *p*-values were adjusted  
363 for multiple comparisons (Bonferroni method). Significance of each comparison is shown  
364 above each comparison pair (\*\*: adjusted  $p \leq 0.01$ ; \*\*\*: adjusted  $p \leq 0.001$ ; \*\*\*\*: adjusted  
365  $p \leq 0.0001$ ).

366  
367 Altogether, these results suggest that oxidative stresses affect variation in STRs.  
368 Although a laboratory mutation in *mev-1* might have increased oxidative stresses and  
369 led to more deletions and insertions in STRs, natural genetic variation that promoted the  
370 expression of *ctl-1* might reduce oxidative stress, which might stabilize STRs to prevent  
371 mutations (Fig. 4b).

372

373 **Discussion**

374 Natural allelic variation in different classes of genomic loci contributes to gene  
375 expression variation<sup>3-5,17-19</sup>. We previously identified thousands of eQTL correlated with  
376 SNVs across wild *C. elegans* strains<sup>5</sup>. Here, we performed genome-wide analysis on how  
377 one of the most polymorphic and abundant repetitive elements, STRs<sup>36</sup>, might affect  
378 expression variation in *C. elegans*. We identified nearly 5,000 associations between STR  
379 variation and expression variation of nearby and remote genes (Fig. 1, Fig. 3). It is  
380 important to note that the number of eSTRs that we detected only represents a  
381 conservative estimate because of the strict significance threshold that we applied.

382 We previously performed genome-wide association analysis on phenotypic  
383 variation in 11 organismal complex traits using pSTR length variation<sup>36</sup> and SNVs<sup>64-73</sup>  
384 respectively. Most of the significant STRs were located within or close to the QTL regions  
385 of interest identified using SNVs and GWA mappings, indicating close relationships

386 between significant STRs and QTL. In the detection of eSTRs, we modeled pSTRs<sup>36</sup>  
387 within 2 Mb surrounding each of the 25,849 transcripts with reliable expression data<sup>5</sup>  
388 (Fig. 1). Close to 84% of transcripts found with local eSTRs were previously detected  
389 with local eQTL<sup>5</sup>, indicating close relationships between eSTRs and eQTL. Therefore, we  
390 further modeled pSTRs within 2 Mb surrounding the QTL regions of interest for  
391 transcripts with detected distant eQTL. Our results revealed important roles of distant  
392 eSTRs underlying distant eQTL and hotspots (Fig. 3). Among transcripts with both eSTRs  
393 and eQTL, 48% of local and 28% of distant eSTR-eQTL pairs showed strong LD with  
394 each other and explained similar amounts of expression variance (Fig. 1b, Fig. 3a). Future  
395 work using simulations and experiments is necessary to partition the contributions of  
396 eSTRs and eQTL to gene regulatory differences. Additionally, we also found 17% of local  
397 and 34% of distant eSTR-eQTL pairs showed low LD with each other (Fig. 1b, Fig. 3a).  
398 Among these low LD eSTR-eQTL pairs, 69% of local and 60% of distant eSTRs had  
399 three to six alleles used in LRT tests (Fig. 1b, Fig. 3a), indicating independent roles of  
400 eSTRs, especially multiallelic STRs, in explaining expression variance. Note that the LD  
401 between eQTL and multiallelic STRs might be overestimated because we transformed  
402 multiallelic STR genotypes to biallelic to calculate LD (See Methods). Therefore,  
403 potentially more multiallelic eSTRs than we reported could have affected expression  
404 independently from eQTL, which could help explain the missing heritability in complex  
405 gene expression traits. One future direction that we did not explore is how epistasis, the  
406 interactions between STRs and SNVs or other mutations, affects gene expression<sup>74,75</sup>.

407 STRs have been proposed to regulate gene expression using various molecular  
408 mechanisms<sup>17,25,76-79</sup>. We found local eSTR variants that caused a variety of mutations in  
409 the target transcripts. We dissected how a 3-bp insertion in an eSTR of the gene  
410 *R07B7.2* altered 3' splice site to change alternative splicing efficiency and cause  
411 differential transcript usage (Fig. 2). The function of the gene *R07B7.2* is not well  
412 understood but the expression of *R07B7.2* was found enriched in neurons, such as AVG  
413 and RIM<sup>44</sup>. Future efforts could investigate the neural consequences of different  
414 transcript usage in the gene *R07B7.2*. Furthermore, we found that distant eSTRs might  
415 affect gene expression by disrupting miRNA binding in the 3'UTRs of genes encoding

416 TFs, such as ATF-7 (Supplementary Fig. 1). Although the variation of STR\_12763 and  
417 expression variation of *atf-7* was not significant in the local eSTR identification, it is  
418 possible that the effects of STR\_12763 variation on the expression of *atf-7* were too  
419 small to be detected using data from 207 strains. But the small changes in the  
420 abundance of the ATF-7 protein might cause strong expression differences in the ATF-  
421 7 targets, which were detectable within the power of our study. In addition to TFs, we  
422 also identified that the eSTR STR\_13795 might affect four genes (*cls-2*, *ddx-23*, *pck-2*,  
423 and *F54E7.9*) in the spindle assembly pathways through both local and distant  
424 regulation. It is possible that *cls-2* is at the upstream of the pathway and its expression  
425 could affect the other three downstream genes. Several mutants of *cls-2* have been  
426 generated<sup>80</sup>. Future work could use these mutants to first examine whether the  
427 expression of *cls-2* affects the other three genes and then validate the role of STR\_13795  
428 mutations in expression regulation.

429 Not only did we observe eSTR that altered gene expression, we also found that  
430 gene expression variation might affect STR mutations. We performed GWA mappings  
431 and mediation analysis on an STR variation trait and identified a candidate gene, *ctl-1*,  
432 that functions in the detoxification pathway of reactive oxygen species (ROS) (Fig. 4,  
433 Supplementary Fig. 4b). We observed low levels of genome-wide STR mutations in  
434 strains with high expression of *ctl-1* (Fig. 4b), which might have increased the antioxidant  
435 capacity in the animal to stabilize the genome and reduce mutations. The effects of ROS  
436 on STR mutations were also revealed by *mev-1* MA lines, which experienced elevated  
437 oxidative stresses and showed higher STR deletion and insertion rates than wild type  
438 MA lines (Fig. 5).

439 Not every strain with low levels of STR mutations had high levels of *ctl-1*  
440 expression (Fig. 4b), suggesting STR mutations are polygenic. For example, other genes  
441 that are responsible for stress response in *C. elegans* might also affect STR mutations.  
442 Fungal infections were found to induce STR expansion in wheat<sup>50</sup>. Various natural  
443 pathogens of *C. elegans* have been discovered<sup>81-84</sup>, future work could compare STR  
444 mutations among *C. elegans* strains isolated from locations with or without known  
445 pathogens. Additionally, genes that are related to DNA replication, repair, or the mitotic

446 process, such as the second top mediator, *F59E12.15* (Fig. 4a), could also cause  
447 genome-wide effects on STR mutations.

448 Altogether, our study provides the first large-scale analysis of associations  
449 between STRs and gene expression variation in wild *C. elegans* strains. We highlighted  
450 the role of eSTRs in explaining expression variation and missing heritability. We also  
451 proposed that oxidative stress might have driven STR mutations globally. STRs have  
452 been proposed to facilitate adaptation and accelerate evolution<sup>12,16–19,25,31,85</sup>. Future work  
453 could use our data and analysis framework to study how STR variation affects complex  
454 traits and facilitates adaptation of *C. elegans* in the wild.

455

## 456 **Methods**

### 457 ***C. elegans* expression and STR data**

458 We obtained summarized expression data of 25,849 transcripts of 16,094 genes and  
459 genotypes of 9,691 polymorphic STRs (pSTRs) in 207 *C. elegans* strains from the original  
460 studies<sup>5,36</sup>. We also obtained 6,545 eQTL positions, their QTL regions of interest, and  
461 eQTL classification from the *C. elegans* eQTL study<sup>5</sup>.

462

### 463 **Expression STRs (eSTRs) identification**

#### 464 *STR genotype transformation*

465 Genotypes of each pSTRs for each strain were transformed as previously described<sup>36</sup>.  
466 Briefly, we used single digits (e.g., “0”, “1”, “2”) to represent STR genotypes in strains  
467 with homozygous alleles (e.g., “0|0”, “1|1”, “2|2”); we chose the smaller digits (e.g., “0”,  
468 “1”, “2”) to represent STR genotypes in strains with heterozygous alleles (e.g., “0|1”,  
469 “1|2”, “3|2”).

470

#### 471 *Selection of STRs for eSTRs identification tests*

472 To identify local eSTRs, we selected pSTRs within 2 Mb surrounding each of the 25,849  
473 transcripts with reliable expression measurements<sup>5</sup>. To identify distant eSTRs, we  
474 selected pSTRs within 2 Mb surrounding the QTL regions of interest for each of the 2,553  
475 transcripts with detected distant eQTL<sup>5</sup>. Among selected pSTRs for each transcript, we  
476 further selected STRs with at least two common variants (frequency > 0.05) among  
477 strains with both STR genotype and expression data, and only retained strains with  
478 common STR variants.

479

#### 480 *Likelihood-ratio test (LRT) to identify eSTRs*

481 We treated STR genotypes as factorial variables and performed LRT on the full model  
482  $lm(\text{expression} \sim \text{STR})$  and the reduced model  $lm(\text{expression} \sim 1)$  using the *lrtest()* function  
483 in the R package *lrtest* (v0.9-39) ([https://cran.r-](https://cran.r-project.org/web/packages/lrtest/index.html)  
484 [project.org/web/packages/lrtest/index.html](https://cran.r-project.org/web/packages/lrtest/index.html)). The Bonferroni threshold was used to  
485 identify significant eSTRs. For each test using real data, we also performed another LRT  
486 using permuted data by shuffling STR genotypes across strains.

487

#### 488 *eSTR identification using STR length variation*

489 Because different alleles of the same STR might have the same length and STR length  
490 variation might have stronger effect on gene expression than substitution, we performed  
491 LRT using the mean allele length of the two copies of each STR for each strain as factorial  
492 variables. We performed STR selection, permutation, LRT, and the Bonferroni threshold  
493 as above to identify eSTRs using STR length variation.

494

#### 495 **LD and variance explained by eQTL and eSTRs**

496 We calculated linkage disequilibrium (LD) between top eSTRs and eQTL for transcripts  
497 with both regulatory sites detected. We used eQTL genotypes and STR genotypes to  
498 calculate LD for eSTRs detected by both STR genotype variation and STR length  
499 variation. Only strains used in eSTR identification were used for LD calculation. We  
500 acquired genotypes of wild strains at each eQTL from the hard-filtered isotype variant

501 call format (VCF) file (CeNDR 20210121 release)<sup>86</sup>. For processed STR genotypes, we  
502 further transformed all multiallelic variants into biallelic variants by converting all non-  
503 reference genotypes (1,2,3, etc.) to 1 and kept reference genotypes as 0. Then, we  
504 calculated LD correlation coefficient  $r^2$  for each STR-SNV and SNV-SNV pairs using the  
505 function `LD()` in the R package *genetics* (v1.3.8.1.2) ([https://cran.r-](https://cran.r-project.org/package=genetics)  
506 [project.org/package=genetics](https://cran.r-project.org/package=genetics)). We also used the generic function `cor()` in R and Pearson  
507 correlation coefficient to calculate the expression variance explained by each QTL and  
508 each top eSTR.  
509

## 510 **Genetic basis of STR variation**

### 511 *STR variation trait*

512 We performed GWA mapping to identify the genetic basis of STR variation in *C. elegans*.  
513 For each of the 207 strains, we counted the total number of STRs with no missing  
514 genotypes among the 9,691 polymorphic STRs and the total number ( $N_{total}$ ) of alternative  
515 alleles ( $N_{alt}$ ) for both copies at each site. The STR variation trait, which is used as the  
516 phenotypic input of GWA mappings, was calculated as  $\log_{10}(N_{alt} / 2N_{total})$ .

517

### 518 *Genome-wide association (GWA) mappings*

519 We performed GWA mappings using the pipeline *Nemascan*  
520 (<https://github.com/AndersenLab/NemaScan>) as previously described<sup>52</sup>. Briefly, we  
521 extracted SNVs of the 207 strains from the hard-filtered isotype VCF (CeNDR 20210121  
522 release)<sup>86</sup> and filtered out variants that had any missing genotype calls and variants that  
523 were below the 5% minor allele frequency using *BCFtools* (v.1.9)<sup>39</sup>. We further pruned  
524 variants with a LD threshold of  $r^2 \geq 0.8$  using `-indep-pairwise 50 10 0.8` in *PLINK* (v1.9)<sup>87,88</sup>  
525 to generate the genotype matrix containing 20,402 markers. We then used two  
526 approaches in the software *GCTA* (v1.93.2)<sup>53,54</sup> to perform GWA mappings: 1) the leave-  
527 one-chromosome-out (LOCO) approach, which uses the `-mlma-loco` function to both  
528 construct a kinship matrix using variants in all chromosomes except the chromosome in  
529 testing and perform the GWA mapping; and 2) the INBRED approach, which uses the -

530 *maker-grm-inbred* function to construct a kinship matrix that is designated for inbred  
531 organisms and the *-fastGWA-Imm-exact* function for the GWA mapping<sup>52-54</sup>. An eigen-  
532 decomposition significance threshold (EIGEN) and a more stringent Bonferroni-  
533 corrected significance threshold (BF) were estimated in *Nemascan* for QTL identification.  
534 For EIGEN, we first estimated the number of independent tests ( $N_{test}$ ) within the genotype  
535 matrix using the R package *RSpectra* (v0.16.0) (<https://github.com/yixuan/RSpectra>) and  
536 *correlateR* (0.1) (<https://github.com/AEBilgrau/correlateR>). EIGEN was calculated as -  
537  $\log_{10}(0.05/N_{test})$ . BF was calculated using all tested markers. Here, QTL were defined by  
538 at least one marker that was above BF. QTL regions of interest were determined by all  
539 markers that were above BF and within 1 kb of one another, and 150 more markers on  
540 each flank.

541

#### 542 *Mediation analysis*

543 We performed mediation analysis that is implemented in *Nemascan* to identify the  
544 mediation effect of gene expression on STR variation as previously described<sup>5</sup>. Briefly,  
545 for each QTL of STR variation, we used the genotype (*Exposure*) at the QTL, transcript  
546 expression traits (*Mediator*) that have eQTL<sup>5</sup> overlapped with the QTL, and STR variation  
547 (*Outcome*) as input to perform mediation analysis using the *medTest()* function in the R  
548 package *MultiMed* (v2.6.0)  
549 (<https://bioconductor.org/packages/release/bioc/html/MultiMed.html>) and the *mediate()*  
550 function in the R package *mediation* (v4.5.0)<sup>89</sup>. Significant mediators were identified as  
551 those with adjusted  $p < 0.05$  and interpretable mediation estimate greater than the 99<sup>th</sup>  
552 percentile of all estimates.

553

#### 554 *GWA mapping for the regressed STR variation trait*

555 We regressed the STR variation trait by the expression of the transcript *Y54G11A.6.1* of  
556 the gene *ctl-1* and performed GWA mappings as described above.

## 557 **STR variants in mutation accumulation (MA) lines**

558 We obtained whole-genome sequence data in the FASTQ format of 160 MA lines,  
559 including N2 MA lines: the N2 ancestor and 67 MA lines; *mev-1* MA lines: the *mev-1*  
560 ancestor and 23 MA lines; and PB306 MA lines: the PB306 ancestor and 67 MA lines  
561 (NCBI Sequence Read Archive projects PRJNA395568, PRJNA429972, and  
562 PRJNA665851)<sup>61,62</sup>. We used the pipelines *trim-fq-nf*  
563 (<https://github.com/AndersenLab/trim-fq-nf>) and *alignment-nf*  
564 (<https://github.com/AndersenLab/alignment-nf>) to trim raw FASTQ files and generate  
565 BAM files for each line, respectively<sup>86</sup>. We called STR variants for the 160 lines using the  
566 pipeline *wi-STRs* (<https://github.com/AndersenLab/wi-STRs>)<sup>36</sup>.

567

## 568 **Mutation rate of polymorphic STRs in MA lines**

569 We calculated the STR mutation rate in MA lines as previously described<sup>36</sup> but using  
570 variant calls before filtering by 10% missing data. Briefly, between each MA line and its  
571 ancestor, we selected STR sites with reliable (“PASS”) calls in both lines. Then, for each  
572 STR, we compared the two alleles in the MA line to the two alleles in the ancestor,  
573 respectively, to identify insertion, deletion, substitution, or no mutation. The mutation  
574 rate (per-allele, per-STR, per-generation)  $\mu$  for each type of mutation was calculated as  
575  $m/2nt$  where  $m$  is the number of the mutation,  $n$  is the total number of reliable STRs, and  
576  $t$  is the number of generations<sup>61,62</sup>.

577

## 578 **Statistical analysis**

579 Statistical significance of difference comparisons were calculated using the Wilcoxon  
580 test and  $p$ -values were adjusted for multiple comparisons (Bonferroni method) using the  
581 *compare\_means()* function in the R package *ggpubr* (v0.2.4)  
582 (<https://github.com/kassambara/ggpubr/>). Enrichment analyses were performed using

583 the one-sided Fisher's exact test and were corrected for multiple comparisons  
584 (Bonferroni method).

585

## 586 **Acknowledgments**

587 We would like to thank Timothy A. Crombie and Ryan McKeown for helpful comments  
588 on the manuscript. G.Z. is supported by the NSF-Simons Center for Quantitative Biology  
589 at Northwestern University (awards Simons Foundation/SFARI 597491-RWC and the  
590 National Science Foundation 1764421). E.C.A. is supported by a National Science  
591 Foundation CAREER Award (IOS-1751035) and a grant from the National Institutes of  
592 Health R01 DK115690. The *C. elegans* Natural Diversity Resource is supported by a  
593 National Science Foundation Living Collections Award to E.C.A. (1930382). We would  
594 also like to thank WormBase because without it these analyses would not have been  
595 possible.

596

## 597 **Author contributions**

598 E.C.A. and G.Z. designed the study. G.Z. analyzed the data. G.Z. and E.C.A. wrote the  
599 manuscript.

600

## 601 **Competing Interests**

602 The authors declare no competing interests.

603

## 604 **Data availability**

605 The datasets for generating all figures can be found at  
606 <https://github.com/AndersenLab/Ce-eSTRs>. Expression and eQTL data of wild  
607 *C. elegans* strains were obtained from <https://github.com/AndersenLab/WI-Ce-eQTL><sup>5</sup>.

608 *C. elegans* STR variation data were obtained from <https://github.com/AndersenLab/WI->  
609 [Ce-STRs](#)<sup>36</sup>. The hard-filtered isotype VCF (20210121 release) was obtained from CeNDR  
610 (<https://www.elegansvariation.org/data/release/20210121>). The raw sequencing data of  
611 MA lines were obtained from the NCBI Sequence Read Archive under accession code  
612 PRJNA395568 [<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA395568>],  
613 PRJNA429972 [<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA429972>], and  
614 PRJNA665851 [<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA665851>]<sup>61,62</sup>.

615

## 616 Code availability

617 The code for generating all figures can be found at <https://github.com/AndersenLab/Ce->  
618 [eSTRs](#).

619

## 620 References

- 621 1. Brem, R. B., Yvert, G., Clinton, R. & Kruglyak, L. Genetic dissection of transcriptional  
622 regulation in budding yeast. *Science* **296**, 752–755 (2002).
- 623 2. West, M. A. L. *et al.* Global eQTL mapping reveals the complex genetic architecture of  
624 transcript-level variation in Arabidopsis. *Genetics* **175**, 1441–1450 (2007).
- 625 3. Albert, F. W., Bloom, J. S., Siegel, J., Day, L. & Kruglyak, L. Genetics of trans-regulatory  
626 variation in gene expression. *Elife* **7**, 1–39 (2018).
- 627 4. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human  
628 tissues. *Science* **369**, 1318–1330 (2020).
- 629 5. Zhang, G., Roberto, N. M., Lee, D., Hahnel, S. R. & Andersen, E. C. The impact of species-  
630 wide gene expression variation on *Caenorhabditis elegans* complex traits. *Nat. Commun.*  
631 **13**, 1–13 (2022).
- 632 6. Zan, Y., Shen, X., Forsberg, S. K. G. & Carlborg, Ö. Genetic Regulation of Transcriptional  
633 Variation in Natural Arabidopsis thaliana Accessions. *G3* **6**, 2319–2328 (2016).
- 634 7. Kita, R., Venkataram, S., Zhou, Y. & Fraser, H. B. High-resolution mapping of cis-  
635 regulatory variation in budding yeast. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E10736–E10744  
636 (2017).
- 637 8. GTEx Consortium *et al.* Genetic effects on gene expression across human tissues. *Nature*  
638 **550**, 204–213 (2017).
- 639 9. Evans, K. S. & Andersen, E. C. The Gene *scb-1* Underlies Variation in *Caenorhabditis*  
640 *elegans* Chemotherapeutic Responses. *G3* **10**, 2353–2364 (2020).
- 641 10. Snoek, B. L. *et al.* The genetics of gene expression in a *Caenorhabditis elegans*  
642 multiparental recombinant inbred line population. *G3* **11**, (2021).
- 643 11. Rockman, M. V., Skrovanek, S. S. & Kruglyak, L. Selection at linked sites shapes heritable

- 644 phenotypic variation in *C. elegans*. *Science* **330**, 372–376 (2010).
- 645 12. Gymrek, M. *et al.* Abundant contribution of short tandem repeats to gene expression  
646 variation in humans. *Nat. Genet.* **48**, 22–29 (2016).
- 647 13. Sekar, A. *et al.* Schizophrenia risk from complex variation of complement component 4.  
648 *Nature* **530**, 177–183 (2016).
- 649 14. Boettger, L. M. *et al.* Recurring exon deletions in the HP (haptoglobin) gene contribute to  
650 lower blood cholesterol levels. *Nat. Genet.* **48**, 359–366 (2016).
- 651 15. Song, J. H. T., Lowe, C. B. & Kingsley, D. M. Characterization of a Human-Specific  
652 Tandem Repeat Associated with Bipolar Disorder and Schizophrenia. *Am. J. Hum. Genet.*  
653 **103**, 421–430 (2018).
- 654 16. Press, M. O., McCoy, R. C., Hall, A. N., Akey, J. M. & Queitsch, C. Massive variation of  
655 short tandem repeats with functional consequences across strains of *Arabidopsis thaliana*.  
656 *Genome Res.* **28**, 1169–1178 (2018).
- 657 17. Fotsing, S. F. *et al.* The impact of short tandem repeat variation on gene expression. *Nat.*  
658 *Genet.* **51**, 1652–1659 (2019).
- 659 18. Jakubosky, D. *et al.* Properties of structural variants and short tandem repeats associated  
660 with gene expression and complex traits. *Nat. Commun.* **11**, 2927 (2020).
- 661 19. Reinar, W. B., Lalun, V. O., Reitan, T., Jakobsen, K. S. & Butenko, M. A. Length variation in  
662 short tandem repeats affects gene expression in natural populations of *Arabidopsis*  
663 *thaliana*. *Plant Cell* **33**, 2221–2234 (2021).
- 664 20. Willems, T. *et al.* Population-Scale Sequencing Data Enable Precise Estimates of Y-STR  
665 Mutation Rates. *Am. J. Hum. Genet.* **98**, 919–933 (2016).
- 666 21. Lynch, M. Rate, molecular spectrum, and consequences of human mutation. *Proc. Natl.*  
667 *Acad. Sci. U. S. A.* **107**, 961–968 (2010).
- 668 22. Sun, J. X. *et al.* A direct characterization of human mutation based on microsatellites. *Nat.*  
669 *Genet.* **44**, 1161–1165 (2012).
- 670 23. Gymrek, M., Willems, T., Reich, D. & Erlich, Y. Interpreting short tandem repeat variations  
671 in humans using mutational constraint. *Nat. Genet.* **49**, 1495–1501 (2017).
- 672 24. Mirkin, S. M. Expandable DNA repeats and human disease. *Nature* **447**, 932–940 (2007).
- 673 25. Gemayel, R., Vincens, M. D., Legendre, M. & Verstrepen, K. J. Variable tandem repeats  
674 accelerate evolution of coding and regulatory sequences. *Annu. Rev. Genet.* **44**, 445–477  
675 (2010).
- 676 26. Gymrek, M. A genomic view of short tandem repeats. *Curr. Opin. Genet. Dev.* **44**, 9–16  
677 (2017).
- 678 27. Weiser, J. N., Love, J. M. & Moxon, E. R. The molecular mechanism of phase variation of  
679 *H. influenzae* lipopolysaccharide. *Cell* **59**, 657–665 (1989).
- 680 28. Rothenburg, S., Koch-Nolte, F., Rich, A. & Haag, F. A polymorphic dinucleotide repeat in  
681 the rat nucleolin gene forms Z-DNA and inhibits promoter activity. *Proc. Natl. Acad. Sci. U.*  
682 *S. A.* **98**, 8985–8990 (2001).
- 683 29. Rockman, M. V. & Wray, G. A. Abundant raw material for cis-regulatory evolution in  
684 humans. *Mol. Biol. Evol.* **19**, 1991–2004 (2002).
- 685 30. Contente, A., Dittmer, A., Koch, M. C., Roth, J. & Dobbelstein, M. A polymorphic  
686 microsatellite that mediates induction of PIG3 by p53. *Nat. Genet.* **30**, 315–320 (2002).
- 687 31. Vincens, M. D., Legendre, M., Caldara, M., Hagihara, M. & Verstrepen, K. J. Unstable  
688 tandem repeats in promoters confer transcriptional evolvability. *Science* **324**, 1213–1216  
689 (2009).
- 690 32. Sureshkumar, S. *et al.* A genetic defect caused by a triplet repeat expansion in  
691 *Arabidopsis thaliana*. *Science* **323**, 1060–1063 (2009).

- 692 33. Yáñez-Cuna, J. O. *et al.* Dissection of thousands of cell type-specific enhancers identifies  
693 dinucleotide repeat motifs as general enhancer features. *Genome Res.* **24**, 1147–1156  
694 (2014).
- 695 34. Quilez, J. *et al.* Polymorphic tandem repeats within gene promoters act as modifiers of  
696 gene expression and DNA methylation in humans. *Nucleic Acids Res.* **44**, 3750–3762  
697 (2016).
- 698 35. Willems, T. *et al.* Genome-wide profiling of heritable and de novo STR variations. *Nat.*  
699 *Methods* **14**, 590–592 (2017).
- 700 36. Zhang, G., Wang, Y. & Andersen, E. Natural variation in *C. elegans* short tandem repeats.  
701 *bioRxiv* 2022.06.25.497600 (2022) doi:10.1101/2022.06.25.497600.
- 702 37. Hannan, A. J. Tandem repeat polymorphisms: modulators of disease susceptibility and  
703 candidates for ‘missing heritability’. *Trends Genet.* **26**, 59–65 (2010).
- 704 38. Lee, D. *et al.* Balancing selection maintains hyper-divergent haplotypes in *Caenorhabditis*  
705 *elegans*. *Nat Ecol Evol* **5**, 794–807 (2021).
- 706 39. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping  
707 and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**,  
708 2987–2993 (2011).
- 709 40. Blumenthal, T. & Steward, K. RNA Processing and Gene Structure. in *C. elegans II* (eds.  
710 Riddle, D. L., Blumenthal, T., Meyer, B. J. & Priess, J. R.) (Cold Spring Harbor Laboratory  
711 Press, 2011).
- 712 41. Kudron, M. M. *et al.* The ModERN Resource: Genome-Wide Binding Profiles for Hundreds  
713 of *Drosophila* and *Caenorhabditis elegans* Transcription Factors. *Genetics* **208**, 937–949  
714 (2018).
- 715 42. Sonenberg, N. & Hinnebusch, A. G. Regulation of translation initiation in eukaryotes:  
716 Mechanisms and biological targets. *Cell* **136**, 731–745 (2009).
- 717 43. Jan, C. H., Friedman, R. C., Ruby, J. G. & Bartel, D. P. Formation, regulation and evolution  
718 of *Caenorhabditis elegans* 3'UTRs. *Nature* **469**, 97–101 (2011).
- 719 44. Harris, T. W. *et al.* WormBase: a modern Model Organism Information Resource. *Nucleic*  
720 *Acids Res.* **48**, D762–D767 (2020).
- 721 45. Dumont, J., Oegema, K. & Desai, A. A kinetochore-independent mechanism drives  
722 anaphase chromosome separation during acentrosomal meiosis. *Nat. Cell Biol.* **12**, 894–  
723 901 (2010).
- 724 46. Espiritu, E. B., Krueger, L. E., Ye, A. & Rose, L. S. CLASPs function redundantly to regulate  
725 astral microtubules in the *C. elegans* embryo. *Dev. Biol.* **368**, 242–254 (2012).
- 726 47. Maton, G. *et al.* Kinetochore components are required for central spindle assembly. *Nat.*  
727 *Cell Biol.* **17**, 953 (2015).
- 728 48. Pelisch, F., Bel Borja, L., Jaffray, E. G. & Hay, R. T. Sumoylation regulates protein  
729 dynamics during meiotic chromosome segregation in *C. elegans* oocytes. *J. Cell Sci.* **132**,  
730 (2019).
- 731 49. Schlientz, A. J. & Bowerman, B. C. *elegans* CLASP/CLS-2 negatively regulates membrane  
732 ingression throughout the oocyte cortex and is required for polar body extrusion. *PLoS*  
733 *Genet.* **16**, e1008751 (2020).
- 734 50. Schmidt, A. L. & Mitter, V. Microsatellite mutation directed by an external stimulus. *Mutat.*  
735 *Res.* **568**, 233–243 (2004).
- 736 51. Cooley, M. B., Carychao, D., Nguyen, K., Whitehand, L. & Mandrell, R. Effects of  
737 environmental stress on stability of tandem repeats in *Escherichia coli* O157:H7. *Appl.*  
738 *Environ. Microbiol.* **76**, 3398–3400 (2010).
- 739 52. Widmayer, S. J., Evans, K. S., Zdraljevic, S. & Andersen, E. C. Evaluating the power and

- 740 limitations of genome-wide association studies in *C. elegans*. *G3* (2022)  
741 doi:10.1093/g3journal/jkac114.
- 742 53. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide  
743 complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 744 54. Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-scale  
745 data. *Nat. Genet.* **51**, 1749–1755 (2019).
- 746 55. Taub, J. *et al.* A cytosolic catalase is needed to extend adult lifespan in *C. elegans* daf-C  
747 and clk-1 mutants. *Nature* **399**, 162–166 (1999).
- 748 56. Lin, C. *et al.* Rosmarinic acid improved antioxidant properties and healthspan via the IIS  
749 and MAPK pathways in *Caenorhabditis elegans*. *Biofactors* **45**, 774–787 (2019).
- 750 57. Song, B., Zheng, B., Li, T. & Liu, R. H. SKN-1 is involved in combination of apple peels and  
751 blueberry extracts synergistically protecting against oxidative stress in *Caenorhabditis*  
752 *elegans*. *Food Funct.* **11**, 5409–5419 (2020).
- 753 58. Poetsch, A. R. The genomics of oxidative DNA damage, repair, and resulting mutagenesis.  
754 *Comput. Struct. Biotechnol. J.* **18**, 207–219 (2020).
- 755 59. Joyner-Matos, J., Bean, L. C., Richardson, H. L., Sammeli, T. & Baer, C. F. No evidence of  
756 elevated germline mutation accumulation under oxidative stress in *Caenorhabditis*  
757 *elegans*. *Genetics* **189**, 1439–1447 (2011).
- 758 60. Matsuba, C. *et al.* Invariance (?) of mutational parameters for relative fitness over 400  
759 generations of mutation accumulation in *Caenorhabditis elegans*. *G3* **2**, 1497–1503 (2012).
- 760 61. Saxena, A. S., Salomon, M. P., Matsuba, C., Yeh, S.-D. & Baer, C. F. Evolution of the  
761 Mutational Process under Relaxed Selection in *Caenorhabditis elegans*. *Mol. Biol. Evol.* **36**,  
762 239–251 (2019).
- 763 62. Rajaei, M. *et al.* Mutability of mononucleotide repeats, not oxidative stress, explains the  
764 discrepancy between laboratory-accumulated mutations and the natural allele-frequency  
765 spectrum in *C. elegans*. *Genome Res.* **31**, 1602–1613 (2021).
- 766 63. Ishii, T. *et al.* Model animals for the study of oxidative stress from complex II. *Biochim.*  
767 *Biophys. Acta* **1827**, 588–597 (2013).
- 768 64. Cook, D. E. *et al.* The Genetic Basis of Natural Variation in *Caenorhabditis elegans*  
769 Telomere Length. *Genetics* **204**, 371–383 (2016).
- 770 65. Zdraljevic, S. *et al.* Natural variation in a single amino acid substitution underlies  
771 physiological responses to topoisomerase II poisons. *PLoS Genet.* **13**, e1006891 (2017).
- 772 66. Hahnel, S. R. *et al.* Extreme allelic heterogeneity at a *Caenorhabditis elegans* beta-tubulin  
773 locus explains natural resistance to benzimidazoles. *PLoS Pathog.* **14**, e1007226 (2018).
- 774 67. Lee, D. *et al.* Selection and gene flow shape niche-associated variation in pheromone  
775 response. *Nat Ecol Evol* **3**, 1455–1463 (2019).
- 776 68. Brady, S. C. *et al.* A Novel Gene Underlies Bleomycin-Response Variation in  
777 *Caenorhabditis elegans*. *Genetics* **212**, 1453–1468 (2019).
- 778 69. Zdraljevic, S. *et al.* Natural variation in *C. elegans* arsenic toxicity is explained by  
779 differences in branched chain amino acid metabolism. *Elife* **8**, e40260 (2019).
- 780 70. Na, H., Zdraljevic, S., Tanny, R. E., Walhout, A. J. M. & Andersen, E. C. Natural variation in  
781 a glucuronosyltransferase modulates propionate sensitivity in a *C. elegans* propionic  
782 acidemia model. *PLoS Genet.* **16**, e1008984 (2020).
- 783 71. Evans, K. S. *et al.* Natural variation in the sequestosome-related gene, sqst-5, underlies  
784 zinc homeostasis in *Caenorhabditis elegans*. *PLoS Genet.* **16**, e1008986 (2020).
- 785 72. Evans, K. S. *et al.* Two novel loci underlie natural differences in *Caenorhabditis elegans*  
786 abamectin responses. *PLoS Pathog.* **17**, e1009297 (2021).
- 787 73. Zhang, G., Mostad, J. D. & Andersen, E. C. Natural variation in fecundity is correlated with

- 788 species-wide levels of divergence in *Caenorhabditis elegans*. *G3* (2021)  
789 doi:10.1093/g3journal/jkab168.
- 790 74. Undurraga, S. F. *et al.* Background-dependent effects of polyglutamine variation in the  
791 *Arabidopsis thaliana* gene ELF3. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 19363–19367 (2012).
- 792 75. Press, Maximilian O., Carlson, K. D. & Queitsch, C. The overdue promise of short tandem  
793 repeat variation for heritability. *Trends Genet.* **30**, 504–512 (2014).
- 794 76. Raveh-Sadka, T. *et al.* Manipulating nucleosome disfavoring sequences allows fine-tune  
795 regulation of gene expression in yeast. *Nat. Genet.* **44**, 743–750 (2012).
- 796 77. Afek, A., Schipper, J. L., Horton, J., Gordân, R. & Lukatsky, D. B. Protein-DNA binding in  
797 the absence of specific base-pair recognition. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 17140–  
798 17145 (2014).
- 799 78. Conlon, E. G. *et al.* The C9ORF72 GGGGCC expansion forms RNA G-quadruplex  
800 inclusions and sequesters hnRNP H to disrupt splicing in ALS brains. *Elife* **5**, (2016).
- 801 79. Liu, X. S. *et al.* Rescue of Fragile X Syndrome Neurons by DNA Methylation Editing of the  
802 FMR1 Gene. *Cell* **172**, 979–992.e6 (2018).
- 803 80. Munoz, N. R., Black, C. J., Young, E. T. & Chu, D. S. New alleles of *C. elegans* gene *cls-2*  
804 (*R107.6*), called *xc3*, *xc4*, and *xc5*. *MicroPubl Biol* **2017**, (2017).
- 805 81. Troemel, E. R., Félix, M.-A., Whiteman, N. K., Barrière, A. & Ausubel, F. M. Microsporidia  
806 are natural intracellular parasites of the nematode *Caenorhabditis elegans*. *PLoS Biol.* **6**,  
807 2736–2752 (2008).
- 808 82. Félix, M.-A. *et al.* Natural and experimental infection of *Caenorhabditis* nematodes by  
809 novel viruses related to nodaviruses. *PLoS Biol.* **9**, e1000586 (2011).
- 810 83. Zhang, G. *et al.* A Large Collection of Novel Nematode-Infecting Microsporidia and Their  
811 Diverse Interactions with *Caenorhabditis elegans* and Other Related Nematodes. *PLoS*  
812 *Pathog.* **12**, e1006093 (2016).
- 813 84. Luallen, R. J. *et al.* Discovery of a Natural Microsporidian Pathogen with a Broad Tissue  
814 Tropism in *Caenorhabditis elegans*. *PLoS Pathog.* **12**, e1005724 (2016).
- 815 85. King, D. G. Indirect selection of implicit mutation protocols. *Ann. N. Y. Acad. Sci.* **1267**,  
816 45–52 (2012).
- 817 86. Cook, D. E., Zdraljevic, S., Roberts, J. P. & Andersen, E. C. CeNDR, the *Caenorhabditis*  
818 *elegans* natural diversity resource. *Nucleic Acids Res.* **45**, D650–D657 (2017).
- 819 87. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based  
820 linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 821 88. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer  
822 datasets. *Gigascience* **4**, 7 (2015).
- 823 89. Tingley, D., Yamamoto, T., Hirose, K., Keele, L. & Imai, K. mediation: R Package for Causal  
824 Mediation Analysis. *Journal of Statistical Software, Articles* **59**, 1–38 (2014).

## Supplementary Information

### Genome-wide regulatory effects of STRs stabilized by elevated expression of antioxidant genes in *C. elegans*

Gaotian Zhang<sup>1</sup> and Erik C. Andersen<sup>1,\*</sup>

1. Department of Molecular Biosciences, Northwestern University, Evanston, IL 60208, USA

\*Corresponding author. E-mail: erik.andersen@gmail.com(E.C.A.)

## Description of Additional Supplementary Files

File Name: Supplementary Data 1

Description: List of eSTRs for nearby genes.

File Name: Supplementary Data 2

Description: List of eSTRs for remote genes.

File Name: Supplementary Data 3

Description: List of common distant eSTRs that are in genes encoding TFs or chromatin cofactors, or the distant eSTRs are also local eSTRs for the genes in which they are located.

## Supplementary Tables

### Supplementary Table 1

Number of strains with REF or ALT allele lengths in STR\_13795 and STR\_13083. Only 186 strains with expression data and genotypes at both STR sites are included.

| STR       | Allele | Length | STR_13795   |            |
|-----------|--------|--------|-------------|------------|
|           |        |        | REF         | ALT        |
|           |        | Length | 30          | 13         |
| STR_13083 | REF    | 16     | 133 strains | 15 strains |
|           | ALT    | 15     | 6 strains   | 32 strains |

## Supplementary Table 2

GSEA results of ten genes that were associated with STR\_13795 in the gene *cls-2*.

| Enrichment term                             | Expected count | Observed count | Enrichment Fold Change | <i>p</i> value | Adjusted <i>p</i> value | Enriched gene                                       |
|---|----------------|----------------|------------------------|----------------|-------------------------|---|
| Oocytes disorganized                        | 0.13           | 2              | 16                     | 0.0003         | 0.072                   | <i>ddx-23</i> ,<br><i>F37C12.1</i>                  |
| Tumorous germline                           | 0.16           | 2              | 13                     | 0.00054        | 0.072                   | <i>ddx-23</i> ,<br><i>F37C12.1</i>                  |
| spindle orientation variant                 | 0.18           | 2              | 11                     | 0.0008         | 0.072                   | <i>ddx-23</i> ,<br><i>F54E7.9</i>                   |
| spindle defective early embryo              | 0.42           | 3              | 7.1                    | 0.0009         | 0.072                   | <i>ddx-23</i> ,<br><i>F54E7.9</i> ,<br><i>pck-2</i> |
| microtubule organization biogenesis variant | 0.5            | 3              | 6                      | 0.0016         | 0.079                   | <i>ddx-23</i> ,<br><i>F54E7.9</i> ,<br><i>pck-2</i> |

### Supplementary Table 3

GWA QTL and regions of interest for the raw and regressed STR variation traits.

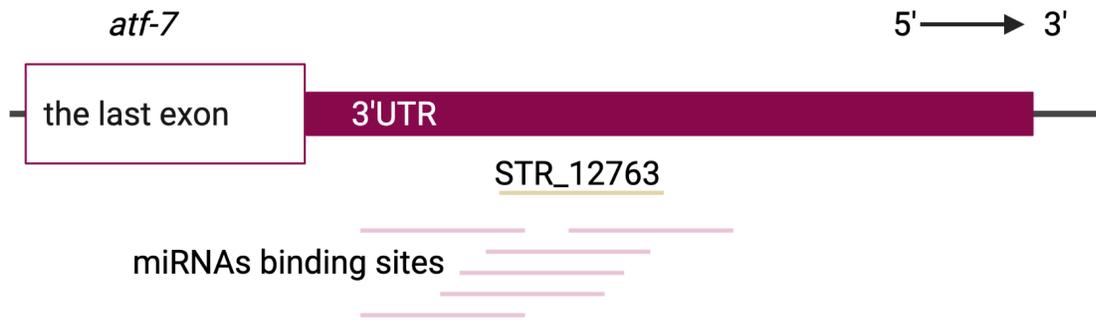
| Trait                   | GWA method | QTL chromosome | QTL peak position | Start position of QTL region of interest | End position of QTL region of interest |
|-------------------------|------------|----------------|-------------------|--|--|
| STR variation           | LOCO       | I              | 12153609          | 4754801                                  | 13671576                               |
| STR variation           | LOCO       | II             | 2706647           | 1477894                                  | 15272855                               |
| STR variation           | LOCO       | III            | 4892173           | 2807201                                  | 13718059                               |
| STR variation           | LOCO       | IV             | 13760657          | 2278611                                  | 15378958                               |
| STR variation           | LOCO       | X              | 2603542           | 826476                                   | 8910667                                |
| STR variation           | LOCO       | X              | 14551237          | 14217837                                 | 17696299                               |
| Regressed STR variation | LOCO       | II             | 11566198          | 10505390                                 | 11842443                               |
| Regressed STR variation | LOCO       | II             | 14625147          | 13968708                                 | 15272855                               |
| Regressed STR variation | INBRED     | II             | 14625147          | 13968708                                 | 15272855                               |

## Supplementary Table 4

Comparison of mutation rates among MA lines using two-sided Wilcoxon tests and Bonferroni method for multiple testing correction.

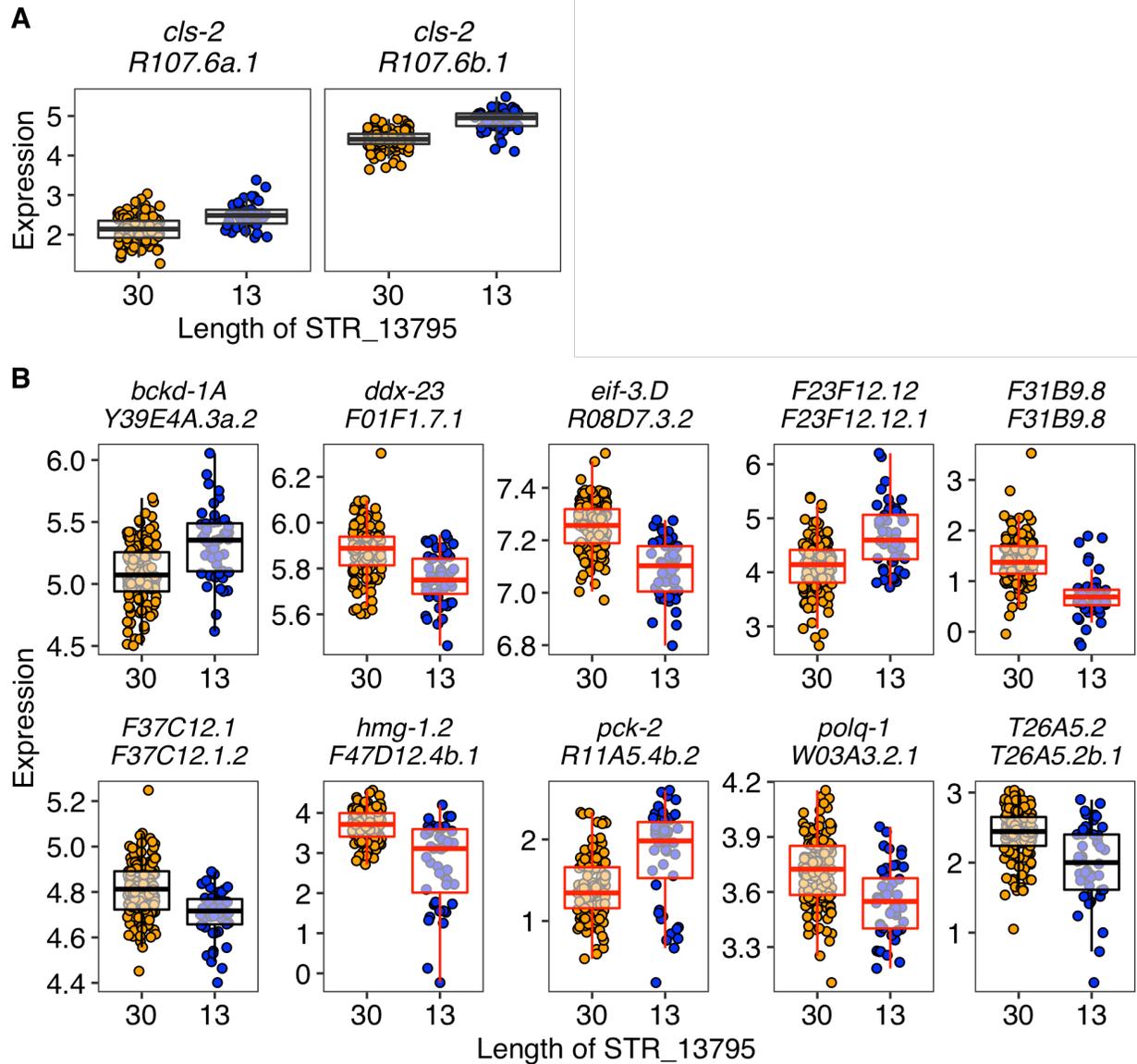
| Mutation      | group1       | group2 | <i>p</i> value | Adjusted <i>p</i> value |
|---------------|--------------|--------|----------------|-------------------------|
| deletions     | <i>mev-1</i> | N2     | 2.46E-05       | 0.00015                 |
| deletions     | <i>mev-1</i> | PB306  | 1.79E-04       | 0.0011                  |
| insertions    | <i>mev-1</i> | N2     | 1.56E-07       | 9.4E-07                 |
| insertions    | <i>mev-1</i> | PB306  | 1.67E-08       | 1E-07                   |
| substitutions | <i>mev-1</i> | N2     | 1.06E-12       | 6.3E-12                 |
| substitutions | <i>mev-1</i> | PB306  | 1.06E-12       | 6.3E-12                 |

## Supplementary Figures



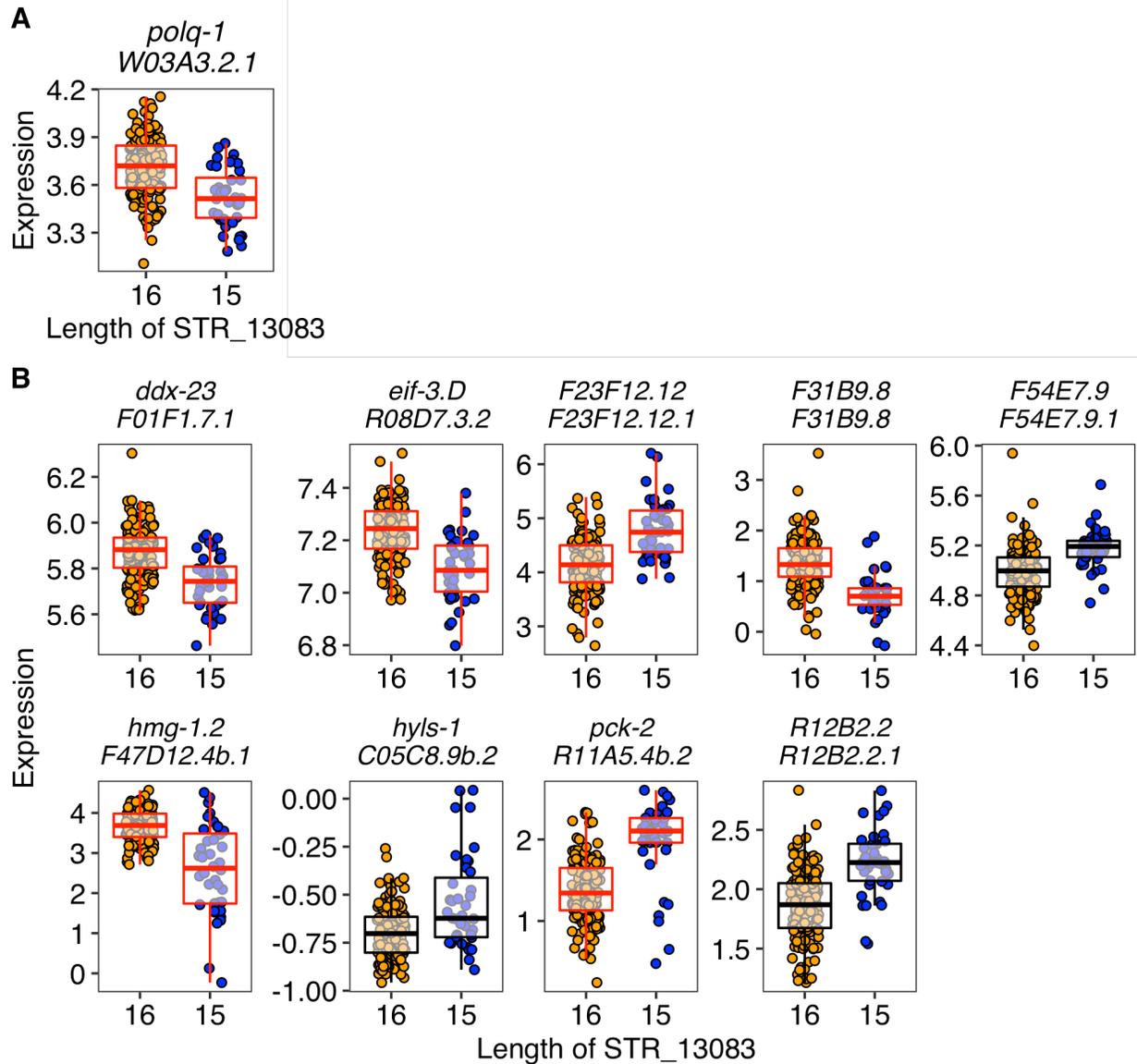
### Supplementary Fig. 1

**STR\_12763 in 3'UTR of the TF gene, *atf-7*, might affect miRNA binding sites.** Graphic illustration of the 3'UTR of *atf-7*, the STR\_12763 (the light brown line), and predicted binding sites of miRNAs (pink lines) based on WormBase<sup>1</sup>. Created using BioRender.com.



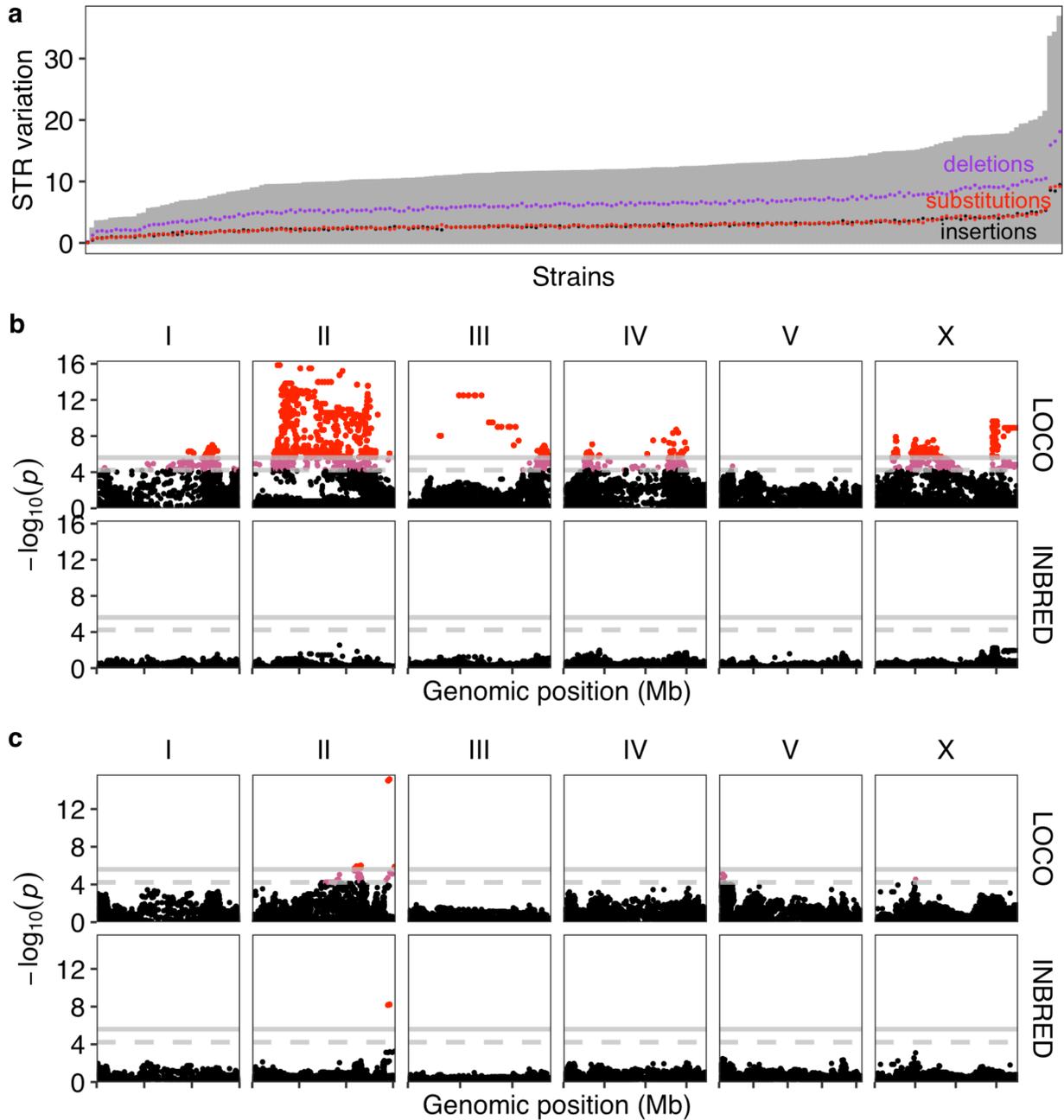
### Supplementary Fig. 2

**The local and distant eSTR, STR\_13795.** STR\_13795 was identified as local eSTRs for two transcripts of the gene *cls-2* (**A**) and distant eSTRs for ten other transcripts (**B**). Tukey box plots showing expression variation of the 12 transcripts between strains with different lengths of the STR\_13795 are shown and colored red for those transcripts with STR\_13083 as an eSTR (Supplementary Fig. 3). Each point corresponds to a strain and is colored orange and blue for strains with the N2 reference allele and the alternative allele, respectively. Box edges denote the 25th and 75th quantiles of the data; and whiskers represent 1.5× the interquartile range.



### Supplementary Fig. 3

**The local and distant eSTR, STR\_13083.** STR\_13083 was identified as local eSTRs for the transcript of the gene *polq-1* (**A**) and distant eSTRs for nine other transcripts (**B**). Tukey box plots showing expression variation of the ten transcripts between strains with different lengths of the STR\_13083 are shown and colored red for those transcripts with STR\_13795 as an eSTR (Supplementary Fig. 2). Each point corresponds to a strain and is colored orange and blue for strains with the N2 reference allele and the alternative allele, respectively. Box edges denote the 25th and 75th quantiles of the data; and whiskers represent 1.5× the interquartile range.



#### Supplementary Fig. 4

**Genetic basis underlying STR variation.** **a** The distribution of an STR variation trait across 207 strains is shown. The STR variation traits calculated by deletions, insertions, and substitutions for each strain are shown as dots and colored purple, black, and red, respectively. **b** Manhattan plots indicating the GWA mapping results for STR variation across 207 strains using LOCO and INBRED approaches are shown, respectively. **c** Manhattan plots indicating the GWA mapping results for STR variation regressed by the expression of *Y54G11A.6.1* of the gene *ctl-1* across 206 strains using LOCO and

INBRED approaches are shown, respectively. In **b** and **c**, each point represents an SNV that is plotted with its genomic position (x-axis) against its  $-\log_{10}(p)$  value (y-axis) in mapping. SNVs that pass the genome-wide EIGEN threshold (the dashed gray horizontal line) and the genome-wide Bonferroni threshold (the solid gray horizontal line) are colored pink and red, respectively. QTL were identified using the Bonferroni threshold.

## REFERENCES

1. Harris, T. W. *et al.* WormBase: a modern Model Organism Information Resource. *Nucleic Acids Res.* **48**, D762–D767 (2020).