

1 **Transposon-mediated genic rearrangements underlie variation in small RNA pathways**

2

3

Gaotian Zhang<sup>1,\*</sup>, Marie-Anne Félix<sup>1,\*</sup>, and Erik C. Andersen<sup>2,\*</sup>

4

5

1. Institut de Biologie de l'École Normale Supérieure, Paris, Île-de-France, France

6

2. Biology Department, Johns Hopkins University, Baltimore, MD, USA

7

8

ORCID IDs: 0000-0001-6468-1341 (G.Z.); 0000-0003-0229-9651 (E.C.A.)

9

10 \*CORRESPONDENCE: [gzhang@bio.ens.psl.eu](mailto:gzhang@bio.ens.psl.eu), [felix@bio.ens.psl.eu](mailto:felix@bio.ens.psl.eu), [erik.andersen@gmail.com](mailto:erik.andersen@gmail.com)

11 **Abstract**

12 Transposable elements (TEs) are parasitic DNA sequences that insert into the host genome and  
13 can cause alterations in host gene structure and expression. Host organisms cope with the often  
14 detrimental consequences caused by recent transposition and develop mechanisms that repress  
15 TE activities. In the nematode *Caenorhabditis elegans*, a small interfering RNA (siRNA) pathway  
16 dependent on the helicase ERI-6/7 primarily silences long terminal repeat retrotransposons and  
17 recent genes of likely viral origin. By studying gene expression variation among wild *C. elegans*  
18 strains, we discovered that structural variants and transposon remnants at the *eri-6/7* locus alter  
19 its expression in *cis* and underlie a *trans*-acting expression quantitative trait locus affecting non-  
20 conserved genes and pseudogenes. Multiple insertions of the *Polinton* DNA transposon (also  
21 known as *Mavericks*) reshuffled the *eri-6/7* locus in different configurations, separating the *eri-6*  
22 and *eri-7* exons and causing the inversion of *eri-6* as seen in the reference N2 genome. In the  
23 inverted configuration, gene function was previously shown to be repaired by unusual *trans*-  
24 splicing mediated by direct repeats flanking the inversion. We show that these direct repeats  
25 originated from terminal inverted repeats specific to *C. elegans Polintons*. This *trans*-splicing  
26 event occurs infrequently compared to *cis*-splicing to novel downstream exons, thus affecting  
27 the production of ERI-6/7. Diverse *Polinton*-induced structural variations display regulatory  
28 effects within the locus and on targets of ERI-6/7-dependent siRNA pathways. Our findings  
29 highlight the role of host-transposon interactions in driving rapid host genome diversification  
30 among natural populations and shed light on evolutionary novelty in genes and splicing  
31 mechanisms.

## 32 Main

33 Transposable elements (TEs) are ubiquitous mobile DNA sequences. With their parasite-like  
34 nature and the invasive mechanisms of transposition, these selfish genetic elements propagate  
35 in host genomes and cause diverse mutations, ranging from point mutations to genome  
36 rearrangements and expansions<sup>1-3</sup>. They can even transfer horizontally across individuals and  
37 species, leading to movement of genetic material between widely diverged taxa<sup>4,5</sup>. To the hosts,  
38 recent TE insertions are mostly deleterious. Various pathways have evolved in hosts to repress  
39 expression and transposition of TEs<sup>6-9</sup>. By contrast, hosts can also benefit from TEs, because  
40 TE sequences can serve as building blocks for the emergence of protein-coding genes, non-  
41 coding RNAs, centromeres, and *cis*-regulatory elements<sup>10-12</sup>.

42 Small RNAs are widely used to repress expression of TEs and other genes<sup>6,7,9</sup>. In the  
43 nematode *Caenorhabditis elegans*, the helicase ERI-6/7-dependent small interfering RNAs  
44 (siRNAs) primarily target long terminal repeat (LTR) retrotransposons and pairs or groups of non-  
45 conserved genes and pseudogenes that show extensive homology and have likely viral  
46 origins<sup>9,13</sup>. The closest known species of *C. elegans*, *Caenorhabditis inopinata*, lost the *eri-6/7*  
47 related small RNA pathway, which was suggested to have caused the expansion of transposons  
48 in its genome compared to *C. elegans* and another related species, *Caenorhabditis briggsae*<sup>9,14</sup>.  
49 In *C. elegans*, ERI-6/7 is required for the biogenesis of the Argonaute ERGO-1-associated  
50 endogenous siRNAs (Fig. 1a)<sup>13</sup>. Likely because endogenous and exogenous siRNA pathways  
51 share and compete for downstream resources<sup>15</sup>, mutants of *eri-6/7* display enhanced RNA  
52 interference (RNAi) responses to exogenous dsRNAs<sup>16</sup>. Competition also exists among different  
53 endogenous siRNA pathways. Within the *eri-6/7* locus, three other local open reading frames  
54 (*eri-6[e]*, *eri-6[f]*, and *sosi-1*) act independently of one another in a feedback loop to modulate  
55 the expression of ERI-6/7 and maintain a balance between different endogenous siRNAs (Fig.  
56 1a, b)<sup>17</sup>.

57 In addition to the vital role of ERI-6/7 in RNAi pathways, its discovery<sup>16</sup> revealed a highly  
58 unusual expression mechanism. Fischer and Ruvkun showed that *eri-6* and *eri-7*, two adjacent  
59 genes oriented in opposing genomic directions in the *C. elegans* reference strain N2, employ a  
60 *trans*-splicing mechanism to generate fused *eri-6/7* mRNAs encoding the helicase ERI-6/7 (Fig.  
61 1a). They further demonstrated that a direct repeat flanking *eri-6* facilitated the *trans*-splicing  
62 process (Fig. 1a). Remarkably, they also noticed variation of the locus within and between  
63 species: a single contiguous gene structure at the *eri-6/7* locus was found in some wild

64 *C. elegans* strains and the *C. briggsae* reference strain AF16. However, the evolutionary history  
65 and consequence of the polymorphic variation remained unknown.

66 Expression quantitative trait loci (eQTL) are genomic loci that explain variation in gene  
67 expression across a species<sup>18</sup>. We recently conducted a genome-wide eQTL analysis among  
68 207 wild *C. elegans* strains, using single nucleotide variants (SNVs) as markers<sup>19</sup> (Extended Data  
69 Fig. 1a). Here, we show that the *cis*-acting eQTL of the *eri-6/7* locus is associated with a genomic  
70 hotspot enriched for *trans*-acting eQTL of non-conserved genes and pseudogenes, including  
71 known ERI-6/7-dependent siRNA targets. We identify structural variation underlying the *eri-6/7*  
72 eQTL, including a distinct gene structure and multiple TE remnants. Our results further  
73 demonstrate that the insertion of multiple copies of the virus-like DNA transposon, *Polinton*<sup>20,21</sup>,  
74 might have caused gene inversion and fission of a single ancestral *eri-6-7* gene. Although some  
75 wild strains still possess the single *eri-6-7* gene, other strains such as N2 evolved the *eri-6/7*  
76 *trans*-splicing mechanism to compensate for the *eri-6* inversion. The direct repeats used for  
77 *trans*-splicing originated from the terminal inverted repeats (TIRs) of *Polintons*. The neighboring  
78 putative genes *eri-6[e]*, *eri-6[f]*, and *sosi-1* are affected by other *Polinton*-induced structural  
79 variants and could have acquired their regulatory functions because of the inversions. Taken  
80 together, the *eri-6/7* gene structure polymorphisms and further structural variants at the locus  
81 impart sophisticated regulatory effects on the biogenesis of the ERI-6/7 helicase, downstream  
82 siRNAs, and the expression of their novel gene targets.

## 83 **Results**

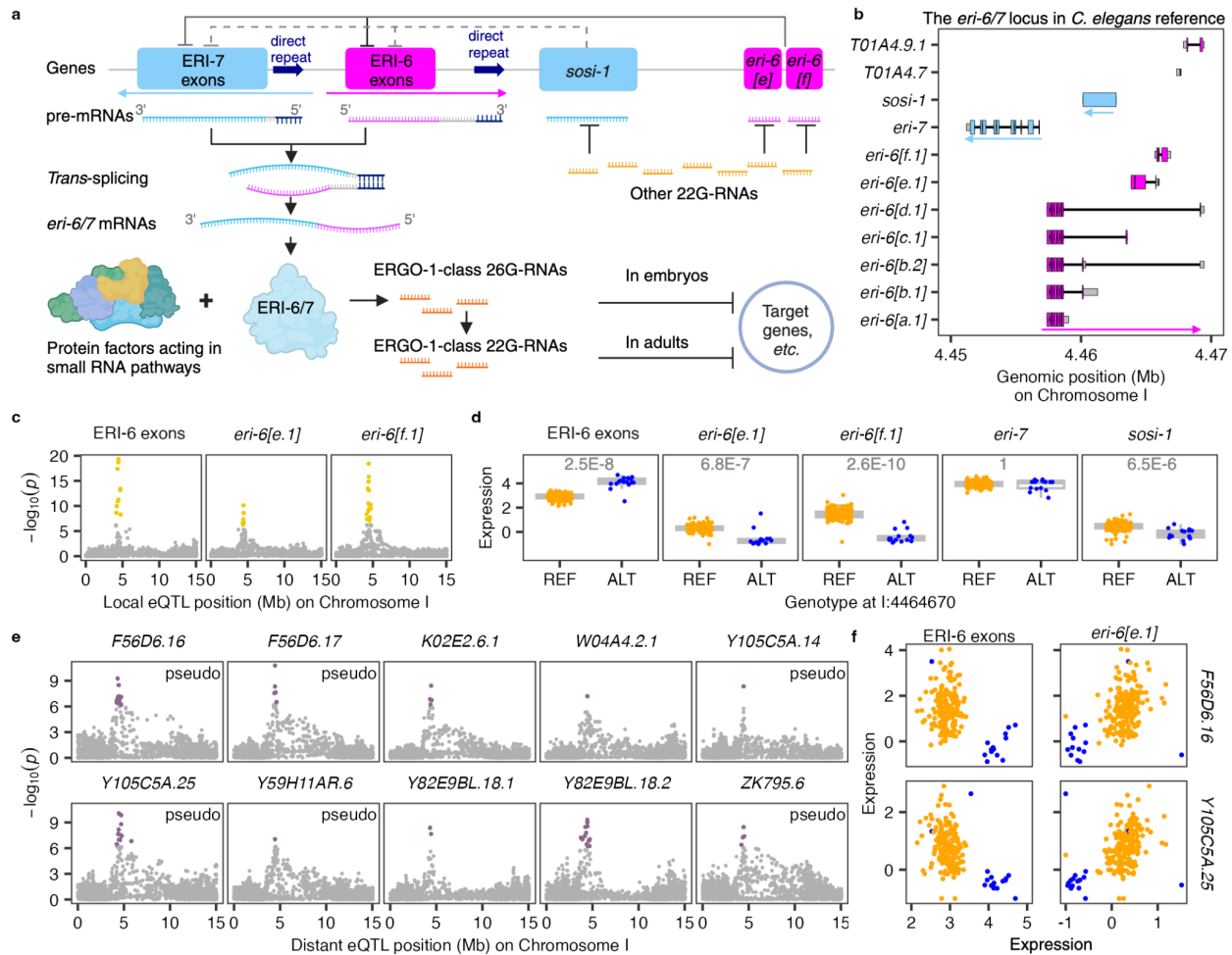
### 84 **Natural variation in *eri-6* underlies differential expression of non-conserved genes and** 85 **pseudogenes**

86 The genes *eri-6* and *eri-7* are next to each other in an opposite head-to-head orientation at 4.45-  
87 4.47 Mb on chromosome I in the N2 reference genome (WS283)<sup>22</sup> (Fig. 1b). The *eri-6* gene has  
88 had a changing transcript annotation in Wormbase<sup>22</sup> because of a variety of rare splicing events.  
89 Presently, it includes six isoforms [*a-f*] that do not all share exons: *eri-6[a-d]* share their first seven  
90 exons (hereafter “ERI-6 exons”, which encode the ERI-6 portion of ERI-6/7) and short  
91 downstream exons, some of them quite distant; *eri-6[e]* and *eri-6[f]* do not share ERI-6 exons  
92 but are transcribed from distinct downstream exons (Fig. 1b). Because the small downstream  
93 exons of *eri-6[a-d]* do not contribute many RNA-seq reads, we used the combined expression

94 of *eri-6[a-d]* as a proxy for the total expression of ERI-6 exons (Extended Data Fig. 1b). We  
95 investigated the genetic basis of expression variation (eQTL) for ERI-6 exons, *eri-6[e]*, *eri-6[ff]*,  
96 and other protein-coding genes in *C. elegans* (See Methods and our previous study<sup>19</sup>,  
97 Supplementary Tables 1, 2). Here, we focused on eQTL related to the *eri-6/7* locus.

98 We classified eQTL into local and distant eQTL based on the location of the QTL in the  
99 genome relative to its expression targets<sup>19</sup> (Extended Data Fig. 1a, Supplementary Table 2). At  
100 the threshold used (see Methods), we detected local eQTL for expression variation in ERI-6  
101 exons, *eri-6[e]* and *eri-6[ff]* (Fig. 1c, Supplementary Table 2). Fine mappings of these local eQTL  
102 identified the top candidate variant (l: 4,464,670), a missense mutation (259D>259Y) in the  
103 coding region of *eri-6[e]*. Strains with the alternative allele at this site showed significantly lower  
104 *eri-6[e]* and *eri-6[ff]* expression than strains with the reference allele but higher expression in ERI-  
105 6 exons (Fig. 1d). Because *eri-6[e]* was found to repress the expression of ERI-6 exons (Fig. 1a)<sup>17</sup>,  
106 it is possible that the alternative non-synonymous allele at the *eri-6[e]* variant could repress the  
107 expression of *eri-6[e]*, which then would enhance expression of ERI-6 exons.

108 Expression variation in ERI-6 exons could further affect the production of the ERI-6/7  
109 helicase, the biogenesis of siRNAs in the ERGO-1 pathway, and finally the expression of target  
110 genes (Fig. 1a). We found that 13 transcripts of 12 genes across the genome, including four  
111 known targets of ERI-6/7-dependent siRNAs<sup>13</sup>, have their distant eQTL (l: 4.3-4.7 Mb) located  
112 nearby the *eri-6/7* locus (Fig. 1e, Extended Data Table 1, Supplementary Table 2). Fine  
113 mappings of these distant eQTL also identified the l: 4,464,670 *eri-6[e]* variant as the top  
114 candidate (Extended Data Table 1). These transcripts showed significantly lower expression in  
115 strains with the alternative allele than strains with the reference allele (Extended Data Fig. 1c).  
116 Their expression also exhibited negative correlations with ERI-6 exons but positive correlations  
117 with *eri-6[e]* expression (Fig. 1f). As mentioned above, pseudogenes and non-conserved genes  
118 are among the primary targets of the ERI-6/7-dependent siRNAs<sup>9,13</sup>. Nine of 12 genes are  
119 pseudogenes and seven of them lack known orthologs in other species<sup>22</sup> (Extended Data Table  
120 1). Taken together, all these 12 genes are potential targets of ERI-6/7-dependent siRNAs.  
121 Genetic variation at the *eri-6/7* locus functions as a *trans*-acting hotspot to regulate expression  
122 of target genes across the genome using the siRNA pathways.



123  
 124 **Fig. 1: Expression variation in *eri-6* potentially mediates a *trans*-acting eQTL hotspot.** **a**,  
 125 Graphic illustration of the ERI-6/7-dependent siRNA pathways and the feedback loop. Dark blue  
 126 arrows indicate direct repeats. Pink and blue rectangles indicate exons on the plus and minus  
 127 strand, respectively (The same color scheme is used in the following figures). Created using  
 128 BioRender. **b**, Structures of genes and isoforms at the *eri-6/7* locus in the reference genome  
 129 (WS283)<sup>22</sup>. **c**, **e**, Manhattan plots indicating the GWAS mapping results of transcript expression  
 130 traits on chromosome I for ERI-6 exons, *eri-6[e]*, and *eri-6[f]* (**c**) and ten transcripts across the  
 131 genome (**e**). Each point represents a SNV that is plotted with its genomic position (x-axis) against  
 132 its  $-\log_{10}(p)$  value (y-axis) in mappings. SNVs that pass the 5% FDR threshold are colored gold  
 133 and purple for local and distant eQTL, respectively. Transcripts of pseudogenes are indicated.  
 134 **d**, Tukey box plots showing expression ( $-\log_2(\text{normalized TPM}+0.5)$ ) variation of five transcripts  
 135 at the *eri-6/7* locus between strains with different alleles at the top candidate SNV (I: 4,464,670).  
 136 Statistical significance of each comparison is shown above and was calculated using the two-  
 137 sided Wilcoxon test and was corrected for multiple comparisons using the Bonferroni method.  
 138 **f**, Correlations of expression variation of two transcripts to expression variation of ERI-6 exons  
 139 and *eri-6[e]*. Each point (**d**, **f**) represents a strain and is colored orange and blue for strains with  
 140 the reference (REF) or the alternative (ALT) allele at the SNV, respectively.  
 141

142 We hypothesized above that the *eri-6[e]* candidate variant could affect the expression of  
 143 *eri-6[e]*, ERI-6 exons, and potential siRNA targets. However, it was unclear why the variant was

144 also associated with *eri-6[ff]* and *sosi-1* expression variation (Fig. 1d). We used CRISPR-Cas9  
145 genome editing to individually introduce the two alleles of the candidate *eri-6[e]* variant into  
146 different genetic backgrounds and showed that this variant did not underlie the local eQTL of  
147 *eri-6* (Extended Data Fig. 2) nor the distant eQTL of potential targets.

148 Two of the strains (CB4856 and MY18) in our expression dataset with an alternative allele  
149 at the *eri-6[e]* variant were previously found to have *eri-6* and *eri-7* on the same (Crick) strand,  
150 similar to the *eri-7* ortholog in the reference genomes of the species *C. briggsae* and *C. brenneri*  
151 (Fig. 2)<sup>16,22</sup>. We thus focused on structural variants, which were not included in the eQTL mapping  
152 because of the difficulty in characterizing them. We first studied them at the genomic level to  
153 uncover the diversity of structural variants, then discovered their transposon origin and finally  
154 demonstrated the association of these structural polymorphisms with a diversity of gene  
155 expression phenotypes.

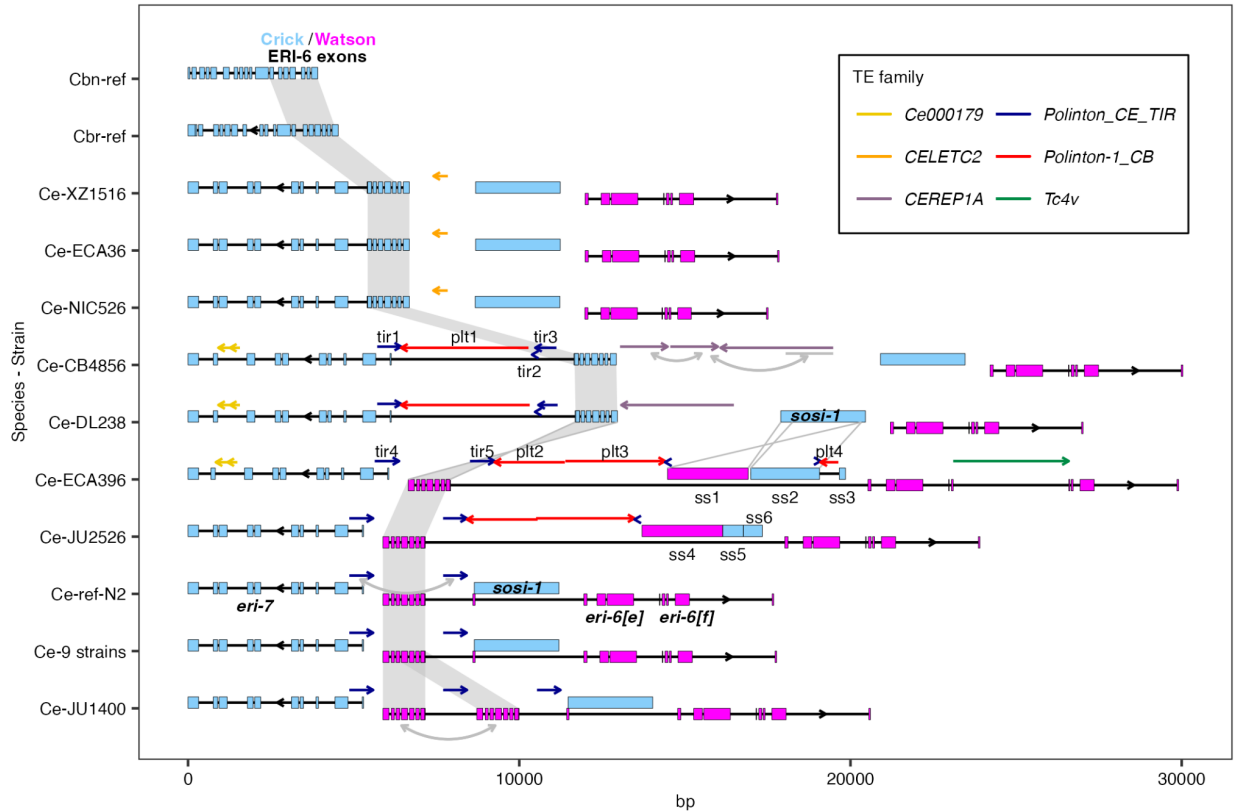
## 156 **High diversity of structural variants and TE insertions throughout the *eri-6/7* locus**

157 Long-read genome assemblies of 17 wild *C. elegans* strains are presently available<sup>23–26</sup>, in  
158 addition to the reference strain N2. We first performed a multiple pairwise alignment of the *eri-*  
159 *6/7* region among these strains (Fig. 2, Extended Data Fig. 3a)<sup>23–27</sup>. Nine of the 17 strains are  
160 approximately identical to the reference strain N2 in this region, with *eri-6* on the Watson strand  
161 (pink in figures) and *eri-7* on the Crick strand (blue in figures). Hereafter, the first seven exons of  
162 *eri-6* in the N2 reference orientation are called “Watson ERI-6 exons”. The strain JU1400 has a  
163 2.8 kb duplication that includes the Watson ERI-6 exons and one copy of the direct repeats that  
164 flank ERI-6 exons (Fig. 2).

165 The other seven strains harbor a large diversity of deletions, insertions, and inversions  
166 compared to the reference genome. The two strains ECA396 and JU2526 have a largely inverted  
167 *sosi-1* gene compared to the N2 strain, two different *sosi-1* fragments, and several other  
168 insertions (Fig. 2, Extended Data Figs. 3a, 4a). The remaining five strains show inversion of ERI-  
169 6 exons compared to the N2 strain (hereafter “Crick ERI-6 exons” when in the same orientation  
170 as *eri-7*): the strains XZ1516, ECA36, and NIC526 also lack the direct repeats that flank ERI-6  
171 exons and include a ~1.7 kb insertion between their Crick ERI-6 exons and *sosi-1*; the strains  
172 CB4856 and DL238 have retained most of the direct repeat sequences and show multiple large  
173 insertions with sizes up to ~8 kb within *eri-7* and surrounding the Crick ERI-6 exons (Fig. 2,  
174 Extended Data Fig. 3a). The Crick orientation of the ERI-6 exons in these five strains likely



175 represents the ancestral genetic structure at the *eri-6/7* locus, based on the following: 1) *eri-6-7*  
 176 orthologs in *C. briggsae* and *C. brenneri* show a simple continuous structure on a single strand  
 177 (Fig. 2); 2) the XZ1516, ECA36, CB4856, and DL238 strains were found to have patterns of  
 178 ancestral genetic diversity in the *C. elegans* species<sup>28-30</sup> (Extended Data Fig. 5).  
 179



180  
 181 **Fig. 2: Hyper-variable structural variants and TEs at the *eri-6/7* locus.** Graphic illustration of  
 182 DNA sequence alignment at the *eri-6/7* locus in the 18 *C. elegans* (Ce) strains with genome  
 183 assemblies. The gene structures of the *C. briggsae* reference (Cbn-ref) *eri-7* and its best match  
 184 homolog in *C. brenneri* reference (Cbr-ref) are shown on top. The exon structures of the  
 185 *C. elegans* strains are shown based on the reference N2 genome. Regions with a potential  
 186 transposon origin are indicated as colored single-headed arrows, with the color indicating the  
 187 type of transposon and the arrow direction representing their potential coding orientation when  
 188 inserted. Double-headed arrows indicate duplications. ERI-6 exons are shaded gray. Detailed  
 189 alignment to the reference of regions with labels “tir1-5” (for terminal inverted repeats), “plt1-4”  
 190 (for *Polintons*), and “ss1-6” (for *sosi-1*) are shown in Extended Data Fig. 4.  
 191

192 This structural diversity corresponds to an astonishing diversity of polymorphic TEs within  
 193 the 18 kb locus (Fig. 2, Extended Data Fig. 3a). First, a 435-bp fragment of *CELETC2* (a  
 194 nonautonomous *Tc2*-related DNA transposon)<sup>22</sup> resides in the ~1.7 kb insertion on the right of  
 195 Crick ERI-6 exons in the strains XZ1516, ECA36, and NIC526. Second, two different fragments



196 (354-bp and 299-bp) of the unclassified transposon *Ce000179*<sup>22</sup> constitute most of the 838-bp  
197 insertion within *eri-7* in the strains CB4856, DL238, and ECA396. Third, a full-length *CEREP1A*  
198 (a putative nonautonomous 3.4-kb DNA transposon likely using *HAT*-related transposase for  
199 propagation)<sup>22</sup> was found in both the CB4856 and DL238 strains, and the CB4856 strain has two  
200 other *CEREP1A* fragments immediately upstream in the opposite orientation. Fourth, the strain  
201 ECA396 has a full-length *Tc4v* (a variant class of the DNA transposon *Tc4*)<sup>22,31</sup> within the first  
202 exon of *eri-6[ff]*. Fifth, we found multiple TE insertions from a family of autonomous double-  
203 stranded DNA transposons derived from viruses, called *Polintons*<sup>20,21</sup>. Four different sizes of  
204 *Polinton* remnants were identified at this locus in the strains CB4856, DL238, ECA396, and  
205 JU2526 (Fig. 2).

### 206 **The direct repeats allowing *eri-6/7* trans-splicing originate from *Polintons***

207 *Polintons* (a.k.a. *Mavericks*) were identified across unicellular and multicellular eukaryotes and  
208 proposed to transpose through protein-primed self-synthesis<sup>5,20,32</sup>. They code numerous  
209 proteins, including two core components, a protein-primed DNA polymerase B (pPolB1) and a  
210 retroviral-like integrase (INT), and different capsid proteins<sup>20,21</sup>. The different *Polinton* remnants  
211 we found at the *eri-6/7* locus in wild strains are all likely from the pPolB1 end of the *Polinton-*  
212 *1\_CB* (named after the *Polintons* in *C. briggsae*, Extended Data Figs. 4b)<sup>22</sup>. In the reference  
213 genome of *C. elegans*, three partial copies of *Polinton-1\_CB* have been identified at 10.30-10.32  
214 Mb (WBTransposon00000738) and 13.08-13.10 Mb (WBTransposon00000637) on chromosome  
215 I and at 17.25-17.27 Mb (WBTransposon00000739) on chromosome X, with lengths ranging from  
216 13.4 to 15.4 kb<sup>22</sup>. We found 744-bp inverted repeats perfectly flanking WBTransposon00000738  
217 (Extended Data Figs. 4b, Supplementary Table 3) and partially flanking the other two *Polintons*  
218 in the genome of the reference strain N2. We hypothesized that these inverted repeats were  
219 specific terminal inverted repeats (TIRs) of *Polintons* in *C. elegans*. They were previously not  
220 regarded as *Polintons* because *C. briggsae* *Polinton* consensus sequences were used to identify  
221 *Polintons* in *C. elegans*. To examine the validity and species specificity of the TIRs, we first  
222 identified potential *Polintons* by searching colocalization (within 20 kb) of pPolB1 and INT in the  
223 genomes of 18 *C. elegans* and three *C. briggsae* strains (Extended Data Figs. 6). We identified  
224 three to nine potential *Polintons* in each *C. elegans* strain and 13 to 15 in each *C. briggsae* strain.  
225 Complete or partial sequences of the 744-bp TIRs were flanking 63 of the total 107 *Polintons* in  
226 the 18 *C. elegans* strains but none in the three *C. briggsae* strains (Extended Data Figs. 6). We

227 also found colocalization of pPolB1 and the TIR but not INT at 10 loci, including but not limited  
228 to the *eri-6/7* locus in *C. elegans* genomes of both N2-like strains and the divergent strains  
229 (Extended Data Figs. 6a). Furthermore, all significant NCBI BLAST<sup>33</sup> results in the query of the  
230 TIR sequence are from *C. elegans*. Taken together, the 744-bp TIRs are components of *Polintons*  
231 specifically in *C. elegans*, termed *Polinton\_CE\_TIR*. We distinguish them from the annotated  
232 *Caenorhabditis Polinton-1\_CB*.

233 The *Polinton\_CE\_TIR* sequences are present as direct repeats instead of inverted repeats  
234 exclusively at the *eri-6/7* locus in the reference N2, the nine N2-like strains, JU1400, JU2526,  
235 and ECA396 (Fig. 2, Extended Data Figs. 6a). In fact, ~700 bp of the ~930-bp direct repeats that  
236 facilitate *trans*-splicing are exactly *Polinton\_CE\_TIR* (Extended Data Figs. 3b, Supplementary  
237 Table 3). The repeat sequences also include new putative TF binding sites for transcriptional  
238 regulation (Extended Data Fig. 3c). Therefore, strains such as the reference N2 use components  
239 of *Polintons* to compensate for the disruptive gene inversion that was likely caused by the  
240 *Polintons* themselves.

#### 241 **Multiple *Polinton* copies likely mediated inversions and other structural rearrangements**

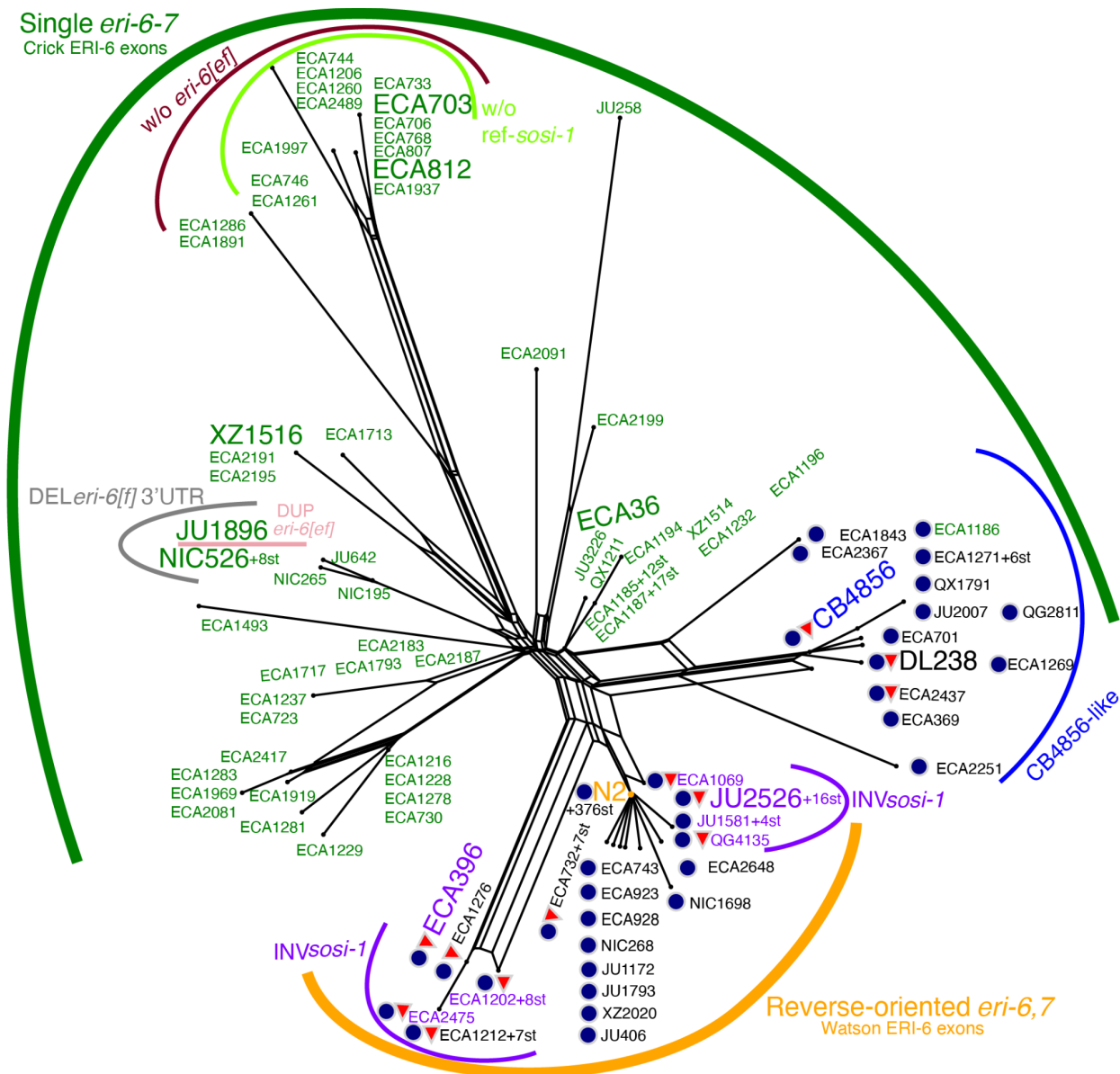
242 To evaluate the diversity of this locus using a larger set of strains, we obtained short-read whole-  
243 genome sequencing (WGS) data of 550 isotype strains, aligned to the reference N2, representing  
244 1384 wild strains from the *Caenorhabditis* Natural Diversity Resource (CaeNDR, 20220216  
245 release)<sup>34</sup>. We aimed to detect inversions and other structural variants in the species using  
246 information of split reads and mapping coverages (See Methods) and relate them to the SNV  
247 haplotypes in the region.

248 We identified diverse structural variants within the *eri-6/7* locus among the 550 wild  
249 strains (Extended Data Figs. 3d, 7, Supplementary Table 4): (1) inversions, 93 strains have Crick  
250 ERI-6 exons and 34 strains have partial inversions of *sosi-1* (*INVsosi-1*) (Extended Data Fig. 7a);  
251 (2) *Polinton* insertions, 48 strains likely have partial remnants of the pPolB1 end of the *Polinton-*  
252 *1\_CB* (Extended Data Fig. 7a); (3) lack of reference genes (which might result from deletion or  
253 maybe an ancestral lack of insertion), 14 strains lack the reference *sosi-1*, *eri-6[e]*, and *eri-6[ff]*,  
254 whereas two strains only lack *eri-6[e]* and *eri-6[ff]* (Extended Data Fig. 7b, d); (4) deletions, 13  
255 strains showed a ~250-bp deletion mostly spanning the 3'UTR of *eri-6[ff]*; (5) duplications, the  
256 strain JU1896, might have duplications of *eri-6[e]* and *eri-6[ff]*; (6) high heterozygosity in *sosi-1*,  
257 80 strains with the reference *sosi-1* might have a second copy of *sosi-1* beyond the locus, which

258 was also possessed by three of the 14 strains lacking the reference *sosi-1* (Supplementary Table  
259 4).

260 The short-read data are limited in their ability to detect the full extent of structural variants.  
261 However, we observed *Polintons* (*Polinton\_CE\_TIR* and *Polinton-1\_CB*) at multiple sites  
262 throughout the *eri-6/7* locus (Extended Data Fig. 7a), especially at flanking regions of ERI-6 exons  
263 and *sosi-1*. TEs have been associated with chromosomal rearrangements since their first  
264 discoveries<sup>1</sup>. Ectopic recombination between TE copies or alternative transposition mechanisms  
265 could cause structural variants such as inversions, duplications, or deletions<sup>2</sup>. We reasoned that  
266 the inversions of ERI-6 exons and *sosi-1* were possibly induced by homologous recombination  
267 between the flanking *Polintons* or simply the TIRs.

268 To understand the evolutionary relationships of the 550 strains at *eri-6/7* and group them,  
269 we performed a haplotype network analysis using the 95 SNVs within the locus (Fig. 3). We  
270 observed and defined two major groups, “Single *eri-6-7*” and “Reverse-oriented *eri-6,7*”, with  
271 112 and 438 strains, respectively (Fig. 3). As expected, a Crick orientation of ERI-6 exons was  
272 detected for all strains in the “Single *eri-6-7*” group, except 17 strains that were clustered with  
273 CB4856 and DL238. We hypothesized that all these 17 strains also have the original Crick  
274 orientation of ERI-6 exons, but with large *Polinton* remnants in between them and *eri-7*: we  
275 defined them as “CB4856-like” strains together with the strains DL238 and ECA1186 (Fig. 3).  
276 The “Reverse-oriented *eri-6,7*” group of strains includes the reference strain N2 and likely all  
277 have Watson ERI-6 exons and the direct repeats for *trans*-splicing (Extended Data Figs. 3d, 7b,  
278 c, Supplementary Table 4). Most strains in this group are clustered with N2, whereas the strain  
279 ECA396 and 19 other strains formed a second cluster based on SNVs and likely all have *INV**sosi-*  
280 *1* (Fig. 3). Remnants of *Polinton-1\_CB* were found in both groups, but mostly in CB4856-like  
281 strains and strains with *INV**sosi-1* (Fig. 3). Strains with deletion polymorphisms in *eri-6[e]*, *eri-6[f]*,  
282 and *sosi-1* formed two clusters exclusively in the “Single *eri-6-7*” group (Fig. 3). It is challenging  
283 to associate these structural variants with *Polintons* or other TEs. Nevertheless, these deletions  
284 and duplications might also affect expression of *eri-6/7* and siRNA pathways.



285  
 286 **Fig. 3: Haplotype network with clustered strains sharing structural variation.** Neighbor-  
 287 joining net depicting 550 strains based on 95 SNVs within the *eri-6/7* locus. Two major groups,  
 288 “Single *eri-6-7*” and “Reverse-oriented *eri-6,7*”, were defined based on orientation of ERI-6  
 289 exons and denoted with dark green and orange curves. Subgroups with other structural  
 290 variations were indicated using thin curves and labels (“w/o” for deletions or no insertions, “DEL”  
 291 for deletions, and “DUP” for duplications). Strain names are colored in green and purple for  
 292 detection of Crick ERI-6 exons and inversion of *sosi-1* (*INVsosi-1*), respectively, using short-read  
 293 WGS data in Extended Data Fig. 7a. Dark blue circles and red triangles next to strain names  
 294 represent strains with *Polinton\_CE\_TIR* (TIRs only) and *Polinton-1\_CB* (TIRs excluded) insertions,  
 295 respectively, based on Extended Data Figs. 3d, 7 and manual inspection of genome alignments.  
 296 Some strains (st) share all alleles of the 95 SNVs and all detected structural variations are  
 297 collapsed to only show a representative strain followed by the number of strains with this *eri-6/7*  
 298 haplotype (e.g. “N2 +376st”). Trapezoidal junctions indicate that some recombination occurred  
 299 within the locus.

## 300 **Cis- and trans-effects of *Polinton*-induced structural variants on gene expression**

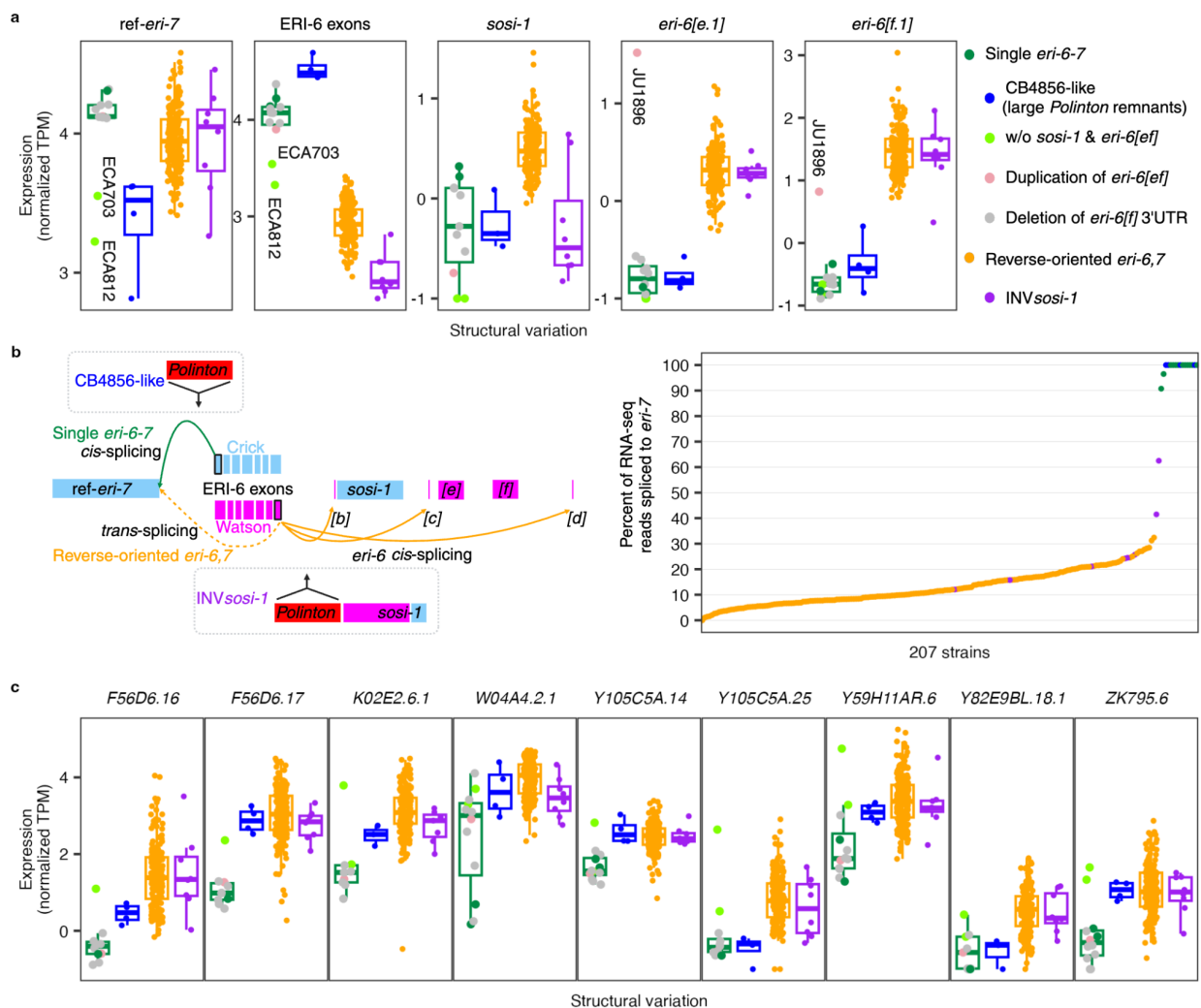
301 Among the 550 wild *C. elegans* strains, ~20% likely have a single "classical" *eri-6-7* gene  
302 to encode the ERI-6/7 protein, as in *C. briggsae* and *C. brenneri*. The remaining ~80% strains  
303 make a fused *eri-6/7* mRNA by some amount of *trans*-splicing between the pre-mRNAs of the  
304 Watson ERI-6 exons and *eri-7* as in the reference strain N2. Though *trans*-splicing compensates  
305 inversion of ERI-6 exons to continue ERI-6/7 production, Fischer and Ruvkun could not consider  
306 whether the reverse-oriented *eri-6/7* gene structure might represent a hypomorphic form of the  
307 locus compared to the ancestral, compact gene. We thus turned our focus back to gene  
308 expression consequences of structural variants, which could affect expression at two levels: the  
309 expression abundances of different exons and their splicing.

310 We first examined local regulatory effects at the *eri-6/7* and *sosi-1* locus, starting with  
311 diversity among strains having the Crick ERI-6 exons and *eri-7*. The strains with a potential  
312 compact *eri-6-7* gene (green box color in Fig. 4a) expressed both parts of the gene at similar  
313 levels, as expected, and expressed low levels of *eri-6[e]*, *eri-6[f]*, and *sosi-1*. The exception in  
314 this group is the two strains ECA703 and ECA812, which do not have *eri-6[e]*, *eri-6[f]*, and *sosi-*  
315 *1* and showed low expression in ERI-6 exons and ref-*eri-7* (mRNA sequences of *eri-6[a-d]* and  
316 *eri-7* in the N2 reference, respectively) (Figs. 3, 4a, Supplementary Table 4). Because *eri-6[e]*, *eri-*  
317 *6[f]*, and *sosi-1* were found to repress *eri-6/7* expression in the reference strain N2<sup>17</sup>, their putative  
318 deletions could cause elevated expression of ERI-6 exons and *eri-7*. These observations suggest  
319 that other linked genetic variation at the locus reduces expression of ERI-6 exons and ref-*eri-7*  
320 or that *eri-6[e]*, *eri-6[f]*, and *sosi-1* function differently in strains of the "Single *eri-6-7*" group. The  
321 strain JU1896, which likely has a duplication in *eri-6[e]* and *eri-6[f]* showed higher expression in  
322 both (Figs. 3, 4a, Extended Data Fig. 7d). The subgroup of CB4856-like strains (blue color), with  
323 large *Polinton* remnants between ERI-6 exons and the downstream ERI-7 exons (Fig. 2),  
324 exhibited significantly elevated expression in ERI-6 exons and significantly decreased  
325 expression in ref-*eri-7*: the large intronic insertion likely affects transcription of the downstream  
326 exons, *i.e.*, *eri-7*.

327 The second large group of strains, those in the "Reverse-oriented *eri-6,7*" group (orange  
328 and purple colors), showed significantly lower expression in ERI-6 exons and significantly higher  
329 expression in *eri-6[e]*, *eri-6[f]*, and *sosi-1* than strains in the "Single *eri-6-7*" group (Fig. 4a,  
330 Supplementary Table 5). The lower expression ERI-6 exons might be the result of either  
331 enhancer/promoter rearrangement or deficiencies in splicing or poly-A tail formation making the



332 mRNA less stable. By contrast, these strains exhibited a similar level of expression of the ref-*eri-7*  
 333 7 to the “Single *eri-6-7*” group. The subgroup of strains with INV*sosi-1* (purple color) showed  
 334 significantly lower expression in both *sosi-1* and ERI-6 exons than other strains in the “Reverse-  
 335 oriented *eri-6,7*” group. Those strains with genome assemblies show large *Polinton* remnants  
 336 upstream of *sosi-1* (Fig. 2), which could explain the lower expression of *sosi-1* and perhaps  
 337 render mRNAs of ERI-6 exons unstable. In summary, the diverse structural variations correlate  
 338 with their expected effect on the *eri-6/7* locus between and within the two large structural variant  
 339 groups.  
 340



341  
 342 **Fig. 4: Structural variations at the *eri-6/7* locus regulate genes in *cis* and *trans*.** **a,c**, Tukey  
 343 box plots showing expression ( $-\log_2(\text{normalized TPM}+0.5)$ ) variation of five transcripts at the *eri-6/7*  
 344 locus (**a**) and nine transcripts (**c**) across the genome that include known targets of siRNAs  
 345 requiring the ERI-6/7 helicase, among strains with major and minor structural variations (SVs)  
 346 within the locus. Each data point represents a strain, color-coded by SVs. Each box is colored

347 by major SVs. Box edges denote the 25<sup>th</sup> and 75<sup>th</sup> quantiles of the data; and whiskers represent  
348 1.5× the interquartile range. Statistical pairwise comparison results using two-sided Wilcoxon  
349 tests and Bonferroni corrections were presented in Supplementary Table 5. **b**, Percent of  
350 spanning RNA-seq reads at the end of the last (seventh) ERI-6 exon that were spliced to *eri-7*  
351 when mapped to the reference genome, for 207 strains. Each point represents one strain and is  
352 colored by SVs. Graphic illustration of structural variation within the *eri-6/7* locus was created  
353 using BioRender.  
354

355 Different splicing mechanisms between the two groups further alter the efficiency of the  
356 ERI-6/7-dependent siRNA pathways. In strains with a single *eri-6-7* gene, the ERI-6/7 protein is  
357 produced through standard transcription and translation. In contrast, strains with reverse-  
358 oriented *eri-6,7* perform separate transcription of pre-mRNAs in opposite orientation and *trans*-  
359 splicing<sup>16</sup>, which could reduce the efficiency of ERI-6/7 production. We analyzed spanning reads  
360 in the RNA-seq data of 207 strains to compare their splicing efficiency between the seventh exon  
361 of *eri-6* and the first of *eri-7* (Fig. 4b). In strains with a single *eri-6-7* gene, most split RNA-seq  
362 reads at the end of the Crick ERI-6 exons should have their chimeric alignment to ERI-7 exons  
363 through *cis*-splicing. In strains with the Watson ERI-6 exons, split RNA-seq reads at the end of  
364 ERI-6 exons could splice to downstream exons for *eri-6[b-d]* or partially map to ERI-7 exons  
365 because of *trans*-spliced *eri-6/7* mRNAs<sup>16</sup> (Fig. 4b). Indeed, among the 207 strains in our RNA-  
366 seq dataset, all 16 strains with a single *eri-6-7* gene showed higher than 90% and mostly 100%  
367 splicing between ERI-6 and ERI-7 exons. Instead, the 183 strains with reverse-oriented *eri-6,7*  
368 but not *INV**sosi-1* showed a median of 10% and a maximum of 32% *trans*-splicing (Fig. 4b). In  
369 conclusion, the evolutionary inversion of *eri-6* does affect the synthesis of full-length *eri-6/7*  
370 mRNA.

371 Together, the expression level of ERI-6 and ERI-7 exons and their splicing rate alter the  
372 biogenesis of the helicase ERI-6/7. Strains with a single *eri-6-7* gene but no extra insertions or  
373 deletions might generate the most abundant ERI-6/7 because of their high expression in ERI-6/7  
374 exons and mostly 100% *cis*-splicing (Fig. 4a, b). The reverse-oriented *eri-6/7* gene structure  
375 represents a hypomorphic form of the locus, because strains in this group showed decreased  
376 expression of ERI-6 exons and low splicing rate between ERI-6/7 exons (Fig. 4a, b), which likely  
377 causes reduced ERI-6/7 protein.

378 The structural variants showed various local effects on gene expression but their  
379 influences likely extend beyond the locus because of the pivotal role of ERI-6/7 in *C. elegans*  
380 endogenous siRNA pathways (Fig. 1a). Differences in ERI-6/7 abundances will affect the  
381 generation of ERGO-1 dependent siRNAs and repression of their target genes. Among the



382 putative targets of ERI-6/7-dependent siRNAs from our eQTL analysis, we observed significantly  
383 lower expression in strains in the “Single *eri-6-7*” group than strains in the “Reverse-oriented *eri-*  
384 *6,7*” group (Fig. 4c, Supplementary Table 5). We also found potential effects of structural variants  
385 in the CB4856-like strains on target expression variation within the “Single *eri-6-7*” group.  
386 Altogether, these results demonstrate that diverse structural variants at the *eri-6/7* locus altered  
387 *C. elegans* endogenous siRNA pathways from the production of the ERI-6/7 helicase to the  
388 expression of target genes.

## 389 Discussion

### 390 Evolutionary genomic history of the *eri-6/7* locus driven by *Polintons*

391 Most strains with a single *eri-6-7* gene were isolated from the Hawaiian Islands or the Pacific  
392 region, where the highest known genetic diversity in the *C. elegans* species is found, (Fig. 3,  
393 Extended Data Figs. 5, 8), which likely reflects the retention of ancestral diversity<sup>28-30</sup>. Strains  
394 with an inversion of ERI-6 exons, however, are more widely distributed over the world and  
395 predominant in Europe. This set of strains show reduced genetic diversity at the locus, in  
396 agreement with an evolutionary-derived inversion of ERI-6 exons from the Crick to the Watson  
397 strand within the species (Extended Data Figs. 5, 8)<sup>30</sup>.

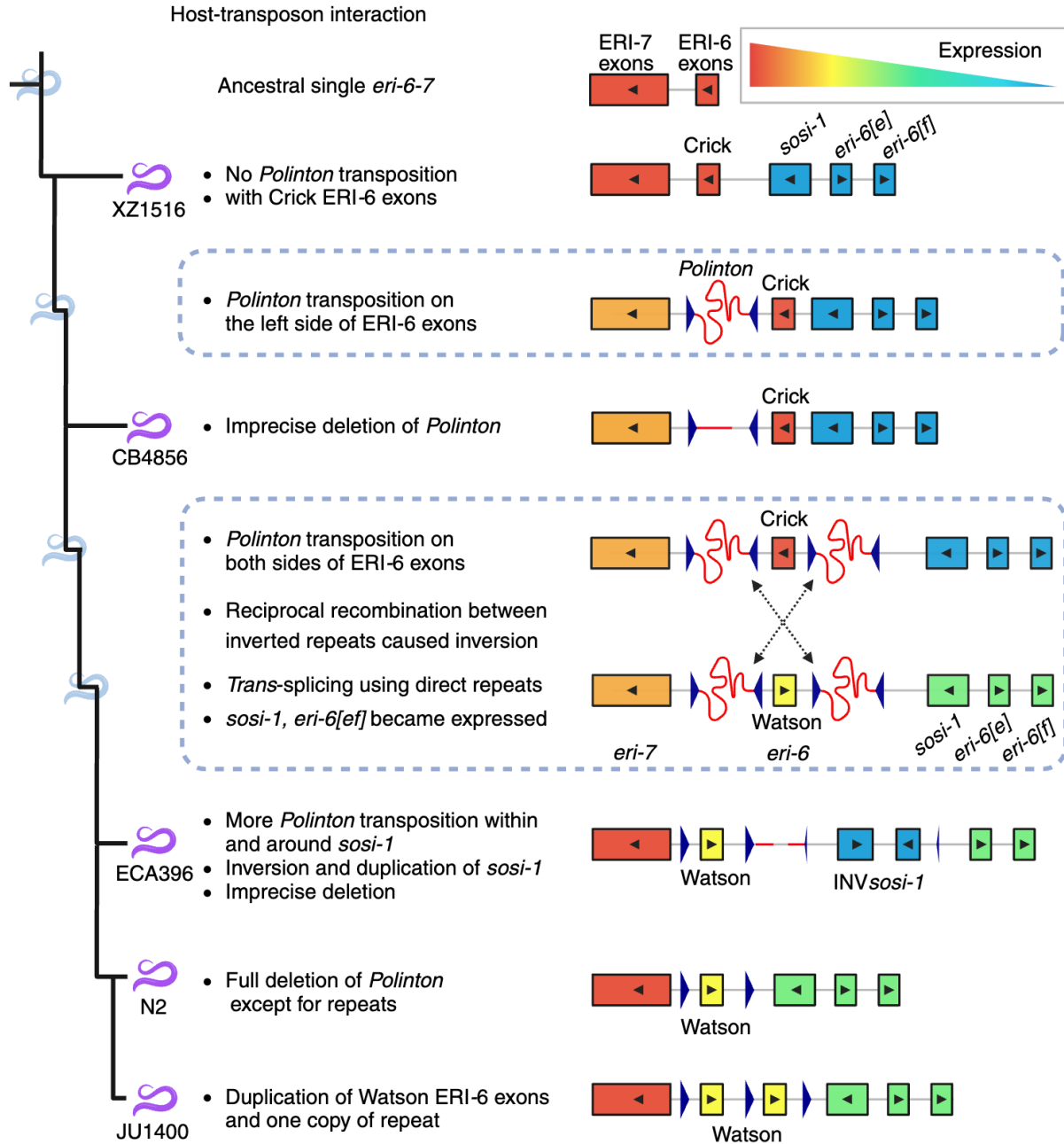
398 We thus favor the following scenario of evolution at the *eri-6/7* locus (Fig. 5). The *eri-6/7*  
399 gene was ancestrally coded as a single gene as in *C. briggsae* and *C. brenneri*, without *Polinton*  
400 insertions. The lack of *eri-6/7* homolog in *C. inopinata*<sup>14</sup> prevents us from using it as a closer  
401 outgroup. The ancestor of all *C. elegans* strains likely conserved the compact single *eri-6-7* gene  
402 structure as in *C. briggsae* and *C. brenneri*. Some strains, such as XZ1516, likely kept this  
403 ancestral single *eri-6-7* gene with no trace of *Polintons* (Figs. 2, 5). Alternatively, in these strains,  
404 the *Polintons* were fully eliminated from the *eri-6/7* locus, yet the parsimonious explanation is  
405 that *Polintons* invaded the locus after the speciation of *C. elegans*.

406 We found *Polinton* remnants in the genome of every *C. elegans* strain with available WGS  
407 data at CaenDR (Extended Data Fig. 6). At some time during the evolutionary history of *C.*  
408 *elegans*, a *Polinton* copy transposed, likely from another location in the genome or through  
409 horizontal transfer, and interrupted the *eri-6/7* gene with a large insertion on the left side of ERI-  
410 6 exons. No strain in our dataset retains a full *Polinton* at the locus, thus this *Polinton* was either  
411 a partial copy when it transposed or subsequently became largely deleted. In strains such as

412 CB4856, the still large *Polinton* remnants (~5 kb in CB4856) appear to impair *eri-7* transcription  
413 (Figs. 4a, 5).

414 Further *Polintons* insertions occurred in the vicinity, including perhaps to the right side of  
415 ERI-6 exons (Figs. 2, 5). The occurrence of several *Polinton* copies at the same locus may have  
416 favored ectopic recombination between inverted sequences and the ERI-6 exon inversion (Fig.  
417 5). Surviving descendants of this inversion, such as the ECA396 and N2 strains, use repeats from  
418 *Polintons* for *trans*-splicing and thus maintain a hypomorphic *eri-6/7* function (Figs. 4c, 5).  
419 Meanwhile, the inversion activated *eri-6[e]* and *eri-6[ff]*, which were barely expressed in most  
420 strains with a single *eri-6-7* gene, at least in the tested conditions (Figs. 4a, 5). Ancestors of the  
421 reference strain N2 eliminated other *Polinton* fragments from the locus, except for the direct  
422 repeats that are necessary for *trans*-splicing. Strains such as JU1400 evolved a duplication of  
423 the Watson ERI-6 exons and one copy of the direct repeat, which may increase the number of  
424 correctly spliced *eri-6/7* transcripts (Figs. 2, 5). *Polintons* might have caused more structural  
425 variations such as *INV**sosi-1* (Figs. 2, 5, Extended Data Fig. 7).

426 The actual evolutionary process within this locus must be more complex than the model  
427 proposed above. The *Polinton* insertions could have occurred through sudden bursts of  
428 transposition instead of gradually. Sudden environmental stress might have caused the high  
429 transposition rate of *Polintons* and the other four TEs (Fig. 2). Overall, the large number of  
430 transposon insertions at this locus regulating small RNA pools and thereby transposons support  
431 the hypothesis of a presumed battle between TE insertions and genomic rearrangement to  
432 preserve ERI-6/7 function to combat further TE activity. Only through further investigations of  
433 gene expression and TE positions in *de novo* assemblies will we learn more about the broad  
434 evolutionary significance of this type of battle.



435

436 **Fig. 5: Possible scenario for evolution at the *eri-6/7* locus with *Polintons*.** Purple and light  
 437 blue worms on the tree represent nodes with or without actual strains, respectively, to the best  
 438 of our knowledge. Rectangles for different segments of *eri-6/7* were filled with gradient colors to  
 439 indicate expression level across segments and branches on the tree. Black triangles inside  
 440 rectangles represent orientation of gene segments. Dark blue triangles represent repeats. Red  
 441 curved lines indicate *Polintons* other than the repeats. Created using BioRender.

## 442 **Phenotypic effect of the structural variation at *eri-6/7* on siRNA pathways and their targets**

443 With the ERGO-1 Argonaute, the ERI-6/7 helicase is required for production of endogenous  
444 primary 26G siRNAs by non-canonical Dicer processing of target mRNAs<sup>13</sup>. Secondary siRNAs  
445 are produced by an amplification machinery, for which different pools of primary siRNAs  
446 compete<sup>15,35</sup>, including endo-siRNAs dependent on Argonautes ERGO-1 and ALG-3/4, the  
447 genomically encoded piRNAs, and the siRNAs derived from exogenous double-stranded  
448 RNAs<sup>13,16,36–38</sup>. Depending on the genomic and environmental contexts, genetic variation favoring  
449 one or the other primary siRNA pathway could have been selected<sup>39–42</sup>. Research in mammals  
450 has shown the importance of dosage of the orthologous MOV10 helicase on retrovirus  
451 silencing<sup>43</sup>. We showed here that natural structural variants at the *eri-6/7* locus were a major  
452 driver of variation in ERGO-1 pathway activity and mRNA levels of its downstream regulated  
453 targets. Two events, likely driven by *Polintons*, lowered ERI-6/7 pathway activity, and increased  
454 piRNA-dependent and exogenous RNAi pathways: (1) the initial insertion of a *Polinton* within the  
455 *eri-6/7* gene and (2) the inversion of ERI-6 exons. Other events might have acted in the reverse  
456 direction: the deletion of most of the intervening *Polintons*, the retention of direct repeats used  
457 in *trans*-splicing and, in the strain JU1400, the duplication of the inverted ERI-6 exons. Because  
458 ERI-6/7-dependent siRNAs primarily target retrotransposons and unconserved, duplicated  
459 genes, with few introns, potentially of viral origins<sup>9,13</sup>, the insertion of the *Polintons* and the  
460 resulting inversion could have at least transiently increased expression of novel genes and  
461 retrotransposons, while repressing exogenous dsRNAs.

462 However, it is unclear what the effect might have been on *Polintons* themselves. Since  
463 their recent discovery in *C. elegans*, their possible regulation by small RNAs remains to be  
464 studied. The DNA polymerase of *Polintons* might be an ancient target of ERI-6/7-dependent  
465 siRNAs, because the gene *E01G4.5*, a known target of ERI-6/7-dependent siRNAs in *C. elegans*,  
466 encodes a protein that has homology to viral DNA polymerases<sup>9,13</sup>. *Polintons* might also bring  
467 novel genes within them<sup>5</sup>, which are potential targets of the ERGO-1 or piRNA pathways. The  
468 genes *sosi-1*, *eri-6[e]*, and *eri-6[f]* are absent at the *eri-6/7* locus in a subset of Hawaiian strains  
469 showing the most divergent *eri-6/7* region based on SNVs (Fig. 3). It is tempting to suggest that  
470 they appeared at this locus during the evolution of the species. The *eri-6[f]* exons are highly  
471 similar to another locus in the genome<sup>17</sup>. The gene *sosi-1* keeps additional copies in some wild  
472 strains and is a distant paralog of *eri-7* and other helicases in its C-terminal part. Further research

473 can test whether *sosi-1*, *eri-6[e]*, and *eri-6[ff]* have been carried by a *Polinton* transposon.  
474 Similarly, the mode of duplication of the ERI-6/7 targets remains to be investigated.

475 Detailed genetic studies in the N2 reference strain have uncovered intricate regulatory  
476 interactions at the *eri-6/7/sosi-1* locus and between this locus and the splicing machinery. First,  
477 in the N2 strain, in part through matching piRNAs, *eri-6[e]*, *eri-6[ff]*, and *sosi-1* are strong ERI-6/7-  
478 independent siRNA targets<sup>17</sup>. Their downregulation by MUT-16-dependent siRNAs enables *eri-*  
479 *6/7* expression, perhaps by spreading chromatin marks<sup>17</sup>. This regulation has been proposed to  
480 act as a negative feedback loop balancing ERGO-1 dependent secondary siRNAs and other  
481 secondary siRNA classes. Second, the use of the *Polinton* repeats as *trans*-splicing signal  
482 partially rescues the production of ERI-6/7. This peculiar mechanism of *eri-6/7 trans*-splicing was  
483 proposed to act as a compensatory sensor of the splicing machinery, enabling more exogenous  
484 RNAi when an overwhelmed splicing machinery increases endo-siRNA production on poorly  
485 spliced genes<sup>44</sup>. It remains unclear whether these seemingly intricate effects on siRNA pools in  
486 the N2 reference strain are an evolutionary leftover of transposon-driven structural variation at  
487 the locus. We hypothesize that across the evolutionary history of *C. elegans*, different siRNA  
488 pools may have been successively favored by natural selection. Alternatively, successive  
489 structural variants could have endowed the *eri-6/7* locus with physiological regulatory loops used  
490 in balancing the different siRNA classes downstream of environmental and organismal inputs.

491 To conclude, our work dissected a distant eQTL hotspot and identified diverse TEs and  
492 structural variations within the *eri-6/7* locus underlying variation in *C. elegans* endogenous siRNA  
493 pathways. This locus appears to have been the target of a large number of TE insertions including  
494 multiple copies of the otherwise rare *Polinton* transposon, which may have caused high genetic  
495 diversity at the locus through genome rearrangements. Some *C. elegans* strains evolved an odd  
496 *trans*-splicing mechanism to maintain hypomorphic function of the locus, using *Polinton* TIRs  
497 that came to form direct repeats. The remarkable interactions between hosts and TEs play a  
498 major role in genome rearrangements and the regulation of gene expression.

## 499 **Methods**

### 500 **Genomic and transcriptomic data**

501 We obtained the reference genomes of *C. elegans* (N2) and *C. briggsae* (AF16), the GTF files of  
502 *C. elegans*, *C. briggsae*, and *C. brenneri*, from WormBase (WS283)<sup>22</sup>; the *de novo* assemblies of  
503 17 wild *C. elegans* strains (CB4856, DL226, DL238, ECA36, ECA396, EG4725, JU310, JU1395,

504 JU1400, JU2526, JU2600, MY2147, MY2693, NIC2, NIC526, QX1794, XZ1516) and two wild  
505 *C. briggsae* strains (QX1410, VX34) from the NCBI Sequence Read Archive (SRA projects  
506 PRJNA523481, PRJNA622250, PRJNA692613, PRJNA784955, and PRJNA819174)<sup>23-27</sup>, the  
507 alignment of whole-genome sequence data in the BAM format of 550 wild *C. elegans* strains, the  
508 soft- and hard-filtered isotype VCF from the *Caenorhabditis* Natural Diversity Resource  
509 (CaeNDR, 20220216 release)<sup>34</sup>; the Illumina RNA-seq FASTQ files of 608 samples of 207 wild *C.*  
510 *elegans* strains from the NCBI SRA (projects PRJNA669810)<sup>19</sup>.

511

## 512 RNA-seq mapping and eQTL analysis

513 To put transcriptomic data on the same page with the genomic data, we re-mapped RNA-seq  
514 reads using the *C. elegans* reference genome (WS283), the GTF file (WS283), and the pipeline  
515 *PEmRNA-seq-nf* (v1.0) (<https://github.com/AndersenLab/PEmRNA-seq-nf>). Then, we selected  
516 reliably expressed transcripts, filtered outlier samples, and normalized expression abundance  
517 across samples using the R scripts *counts5strains10.R*, *nonDivergent\_clustered.R*, and  
518 *norm\_transcript\_gwas.R* (<https://github.com/AndersenLab/WI-Ce-eQTL/tree/main/scripts>),  
519 respectively, as previously described<sup>19</sup>. In summary, we collected reliable expression abundance  
520 for 23,349 transcripts of 16,172 genes (15,449 protein-coding genes and 723 pseudogenes) from  
521 560 samples of 207 strains. We also used *STAR* (v2.7.5)<sup>45</sup> to identify chimeric RNA-seq reads in  
522 the 560 samples.

523 We further used our recently developed GWAS mapping pipeline, *Nemascan*<sup>46</sup>, to identify  
524 eQTL for the 23,349 transcript expression traits (Supplementary Table 1), following the steps  
525 outlined previously<sup>19</sup>. Briefly, we randomly selected 200 traits and permuted each of them 200  
526 times. For each of the 40,000 permuted traits, we used the leave-one-chromosome-out (LOCO)  
527 approach and the INBRED approach in the *GCTA* software (v1.93.2)<sup>47,48</sup>, and calculated the  
528 eigen-decomposition significance (EIGEN) threshold as  $-\log_{10}(0.05/N_{test})$  to identify QTL.

529 We determined the 5% false discovery rate (FDR) significance threshold for LOCO and  
530 INBRED, respectively, by calculating the 95<sup>th</sup> percentile of the significance of all detected QTL  
531 above using each approach. We then performed GWAS mapping on all 23,349 traits using LOCO  
532 and INBRED approaches and identified eQTL that passed their respective 5% FDR thresholds.  
533 Overall, we detected 10,291 eQTL for 5668 transcript expression traits, with 4899 eQTL for 4254  
534 traits in LOCO, 5392 eQTL for 4700 traits in INBRED (Supplementary Table 2). Fine-mappings  
535 were further performed on each eQTL using *Nemascan*.



536 We classified eQTL as local (within 2 Mb surrounding the transcript) or distant (non-local).  
537 For distant eQTL located outside of the common hyper-divergent regions among the 207  
538 strains<sup>19,25</sup>, we identified hotspot regions enriched with distant eQTL for LOCO and INBRED  
539 results, respectively<sup>19</sup>.

540 The genomic region harboring the *eri-6/7* locus at 21 cM on chromosome I was identified  
541 as a distant eQTL hotspot in both LOCO and INBRED in this study and in our previous study<sup>19</sup>.

542

### 543 **DNA alignment**

544 We aligned each of the 17 *de novo* PacBio assemblies of wild *C. elegans* strains to the N2  
545 reference genome using *MUMmer* (v3.1)<sup>49</sup> and extracted sequences that were aligned to the N2  
546 *eri-6/7* locus using *BEDTools* (v2.29.2)<sup>50</sup>. Then, we performed pairwise alignments among these  
547 sequences and to the *eri-6/7* N2 reference sequence using *Unipro UGENE* (v.47.0)<sup>51</sup>. Large  
548 insertions (>50 bp) in the wild strains to the reference were blasted in WormBase<sup>22</sup> to identify  
549 potential transposon origins.

550

### 551 **Scan for *Polinton* and TIRs in genome assemblies**

552 We obtained the amino acid sequences of pPolB1 and INT in *C. briggsae Polinton-1*  
553 (WBTransposon00000832)<sup>22</sup> using *ORFfinder* (<https://www.ncbi.nlm.nih.gov/orffinder/>) and the  
554 744 bp DNA sequence for the TIRs from 10,302,516 to 10,303,259 bp on chromosome I in the  
555 *C. elegans* (N2) reference genome. We searched for the *Polinton* and TIRs sequences in the 21  
556 genome assemblies using *tblastn* and *blastn* in BLAST (v2.14.0)<sup>52</sup>, respectively. We filtered the  
557 results by a maximum e-value of 0.001 and a minimum bitscore of 50<sup>32</sup>. We merged pPolB1, INT,  
558 and TIR hits within 4 kb, 2 kb, and 2 kb, respectively, with consideration of strandedness.  
559 *Polinton* insertions were identified by the presence of both pPolB1 and INT within 20 kb.

560 We also searched for *sosi-1* outside of the *eri-6/7* locus in the genome assemblies using  
561 DNA sequence of *sosi-1* in the reference and found an additional copy in the strains JU2526,  
562 ECA396, XZ1516, and JU1400, and two additional copies in the strains ECA36 and QX1794 in  
563 their PacBio genome assemblies. Genomic locations surrounding these additional copies in the  
564 six strains correspond to ~0.31 Mb on the chromosome III in the reference N2 genome. The  
565 additional copies of *sosi-1* outside the *eri-6/7* locus in the six strains share most alleles compared  
566 to the *sosi-1* within the *eri-6/7* locus.

567

### 568 **Identification of SVs using short-read WGS data**



569 We extracted information of split reads mapped to the reference *eri-6/7* locus (I: 4,451,194 -  
570 4,469,460 bp) and with a minimum quality score equal of 20 from the BAM files of the 550 wild  
571 *C. elegans* strains. 1): To identify potential inversions in the *eri-6/7* locus, we first selected split  
572 reads with both the primary and chimeric alignments mapped to this region but to different  
573 strands. We assigned the primary and chimeric alignment positions of each split read into 200-  
574 bp bins and required at least four reads that had the primary and chimeric alignments in the  
575 same pair of bins for a relatively reliable inversion event in each strain. We focused on inversions  
576 spanning at least three bins and found in more than 10 strains. 2): To identify potential sites of  
577 *Polinton* remnants, we selected the split reads outside of the direct repeats at the *eri-6/7* locus  
578 and with the chimeric alignment mapped to *Polinton* (*Polinton-1\_CB*, WBTransposon00000738)  
579 and its surrounding *Polinton\_CE\_TIR* on chromosome I from 10,302,516 to 10,319,657 bp. At  
580 least two reads were required. The primary alignment of these reads indicated the potential sites  
581 of *Polinton* remnants in the *eri-6/7* locus in wild strains.

582 Furthermore, we counted the coverage per bp in the *eri-6/7* locus for each short-read  
583 WGS BAM file using *BEDTools* (v2.29.2)<sup>50</sup>. We calculated the percentage of the coverage at each  
584 bp to the mean coverage within the *eri-6/7* locus in each strain. Then, we performed a sliding  
585 window analysis with a 200-bp window size and a 100-bp step size for each strain. A 173-bp  
586 tandem repeat region from 4,465,414 to 4,465,586 bp on chromosome I was masked in the  
587 results.

588 To identify additional copies and haplotypes of *sosi-1* among the 550 wild strains, we  
589 focused on 93 variants of the 101 SNVs tagged “high heterozygosity” within the *sosi-1* region in  
590 the soft-filtered isotype VCF (CaeNDR, 20220216 release)<sup>34</sup>. We used the following threshold to  
591 define *sosi-1* haplotype and copy numbers among the 550 strains: 449 strains show  
592 homozygous reference alleles at all 93 SNVs (except one strain at 92 SNVs), indicating they only  
593 have the reference haplotype *sosi-1*; 80 strains show heterozygous alleles at more than 60 SNVs,  
594 indicating two copies of *sosi-1* with divergent haplotypes; three strains have homozygous  
595 alternative alleles at more than 90 SNVs, indicating missing of the reference *sosi-1* in the *eri-6/7*  
596 locus and the existence of the alternative *sosi-1* copy; 11 strains show undetected genotype at  
597 60 to 93 SNVs and extreme low coverages in *sosi-1* (Extended Data Fig. 7d), indicating they may  
598 lack *sosi-1* in the genomes; the *sosi-1* haplotype and copy number of the remaining seven strains  
599 are unclear as they have numbers of homozygous and homozygous alleles in between the above  
600 threshold (Supplementary Table 4).

601

602 **Genetic relatedness**

603 Genetic variation data across the genome among the 550 *C. elegans* strains were extracted from  
604 the hard-filtered VCF above using *BCFtools* (v.1.9)<sup>53</sup>. These variants were pruned to the  
605 1,199,944 biallelic SNVs without missing genotypes. We converted this pruned VCF file to a  
606 PHYLIP file using the *vcf2phylip.py* script<sup>54</sup>. The unrooted neighbor-joining tree was made using  
607 the R packages *phangorn* (v2.5.5)<sup>55</sup> and *ggtree* (v1.14.6)<sup>56</sup>.

608 A second PHYLIP file was built by the same method above but only with 95 SNVs within  
609 the *eri-6/7* locus. A haplotype network was generated using this PHYLIP file and *SplitsTree CE*  
610 (v6.1.16)<sup>57</sup>.

611 **Acknowledgements**

612 This research was supported by the National Science Foundation (1764421) and Simons  
613 Foundation (597491) Research Center for Mathematics of Complex Biological Systems and a  
614 Human Frontiers Research Program grant (RGP0001/2019) to E.C.A. G.Z. is supported by a  
615 grant from the Fondation pour la Recherche Médicale (ARF202209015859). M.-A.F. is supported  
616 by the Centre National de la Recherche Scientifique and a grant from the Agence Nationale de  
617 la Recherche (ANR-19-CE12-0025). We like to thank Emily Koury for the efforts in CRISPR-Cas9  
618 genome editing. We like to thank Nicolas D. Moya for suggestions on useful bioinformatic tools.  
619 We would like to thank members of the Andersen Lab for helpful comments on the manuscript.  
620 We also like to thank CaenDR (supported by NSF Capacity Grant 2224885 to E.C.A.) and  
621 WormBase for which these analyses would not have been possible.

622 **Author contributions**

623 G.Z., M.-A.F., and E.C.A. conceived of the study. G.Z. analyzed the data. G.Z., M.-A.F., and  
624 E.C.A. wrote the manuscript.

625 **Competing interests**

626 The authors declare no competing interests.

627 **Data and code availability**

628 The datasets and code for generating all figures can be found at  
629 <https://github.com/AndersenLab/Ce-eri-67>

630 **Supplementary information**

631 Description of Additional Supplementary Files

632 Supplementary Table 1

633 23,349 GWAS gene expression traits

634 Supplementary Table 2

635 eQTL summary

636 Supplementary Table 3

637 Sequences alignment of *Polinton\_CE\_TIR* and the direct repeat

638 Supplementary Table 4

639 Structural variants at the *eri-6/7* locus of 550 strains

640 Supplementary Table 5

641 Statistical pairwise comparison results in Fig. 4a, c

642 **References**

643 1. McClintock, B. The origin and behavior of mutable loci in maize. *Proceedings of the*  
644 *National Academy of Sciences* **36**, 344–355 (1950).

645 2. Gray, Y. H. It takes two transposons to tango: transposable-element-mediated  
646 chromosomal rearrangements. *Trends Genet.* **16**, 461–468 (2000).

647 3. Bourque, G. *et al.* Ten things you should know about transposable elements. *Genome*  
648 *Biol.* **19**, 199 (2018).

649 4. Gilbert, C. & Feschotte, C. Horizontal acquisition of transposable elements and viral  
650 sequences: patterns and consequences. *Curr. Opin. Genet. Dev.* **49**, 15–24 (2018).

651 5. Widen, S. A. *et al.* Virus-like transposons cross the species barrier and drive the evolution  
652 of genetic incompatibilities. *Science* **380**, eade0705 (2023).

- 653 6. Brennecke, J. *et al.* Discrete small RNA-generating loci as master regulators of transposon  
654 activity in *Drosophila*. *Cell* **128**, 1089–1103 (2007).
- 655 7. Ito, H. Small RNAs and transposon silencing in plants. *Dev. Growth Differ.* **54**, 100–107  
656 (2012).
- 657 8. Deniz, Ö., Frost, J. M. & Branco, M. R. Regulation of transposable elements by DNA  
658 modifications. *Nat. Rev. Genet.* **20**, 417–431 (2019).
- 659 9. Fischer, S. E. J. & Ruvkun, G. *Caenorhabditis elegans* ADAR editing and the ERI-  
660 6/7/MOV10 RNAi pathway silence endogenous viral elements and LTR retrotransposons.  
661 *Proc. Natl. Acad. Sci. U. S. A.* **117**, 5987–5996 (2020).
- 662 10. Rebollo, R., Romanish, M. T. & Mager, D. L. Transposable elements: an abundant and  
663 natural source of regulatory sequences for host genes. *Annu. Rev. Genet.* **46**, 21–42  
664 (2012).
- 665 11. Gao, D., Jiang, N., Wing, R. A., Jiang, J. & Jackson, S. A. Transposons play an important  
666 role in the evolution and diversification of centromeres among closely related species.  
667 *Front. Plant Sci.* **6**, 216 (2015).
- 668 12. Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory activities of transposable elements:  
669 from conflicts to benefits. *Nat. Rev. Genet.* **18**, 71–86 (2017).
- 670 13. Fischer, S. E. J. *et al.* The ERI-6/7 helicase acts at the first stage of an siRNA amplification  
671 pathway that targets recent gene duplications. *PLoS Genet.* **7**, e1002369 (2011).
- 672 14. Kanzaki, N. *et al.* Biology and genome of a newly discovered sibling species of  
673 *Caenorhabditis elegans*. *Nat. Commun.* **9**, 3216 (2018).
- 674 15. Lee, R. C., Hammell, C. M. & Ambros, V. Interacting endogenous and exogenous RNAi  
675 pathways in *Caenorhabditis elegans*. *RNA* **12**, 589–597 (2006).
- 676 16. Fischer, S. E. J., Butler, M. D., Pan, Q. & Ruvkun, G. Trans-splicing in *C. elegans*

- 677 generates the negative RNAi regulator ERI-6/7. *Nature* **455**, 491–496 (2008).
- 678 17. Rogers, A. K. & Phillips, C. M. A Small-RNA-Mediated Feedback Loop Maintains Proper  
679 Levels of 22G-RNAs in *C. elegans*. *Cell Rep.* **33**, 108279 (2020).
- 680 18. Albert, F. W. & Kruglyak, L. The role of regulatory variation in complex traits and disease.  
681 *Nat. Rev. Genet.* **16**, 197–212 (2015).
- 682 19. Zhang, G., Roberto, N. M., Lee, D., Hahnel, S. R. & Andersen, E. C. The impact of species-  
683 wide gene expression variation on *Caenorhabditis elegans* complex traits. *Nat. Commun.*  
684 **13**, 1–13 (2022).
- 685 20. Kapitonov, V. V. & Jurka, J. Self-synthesizing DNA transposons in eukaryotes. *Proc. Natl.*  
686 *Acad. Sci. U. S. A.* **103**, 4540–4545 (2006).
- 687 21. Pritham, E. J., Putliwala, T. & Feschotte, C. Mavericks, a novel class of giant transposable  
688 elements widespread in eukaryotes and related to DNA viruses. *Gene* **390**, 3–17 (2007).
- 689 22. Harris, T. W. *et al.* WormBase: a modern Model Organism Information Resource. *Nucleic*  
690 *Acids Res.* **48**, D762–D767 (2020).
- 691 23. Kim, C. *et al.* Long-read sequencing reveals intra-species tolerance of substantial  
692 structural variations and new subtelomere formation in *C. elegans*. *Genome Res.* **29**,  
693 1023–1035 (2019).
- 694 24. Bubrig, L. T., Sutton, J. M. & Fierst, J. L. *Caenorhabditis elegans* dauers vary recovery in  
695 response to bacteria from natural habitat. *Ecol. Evol.* **10**, 9886–9895 (2020).
- 696 25. Lee, D. *et al.* Balancing selection maintains hyper-divergent haplotypes in *Caenorhabditis*  
697 *elegans*. *Nat Ecol Evol* 1–14 (2021).
- 698 26. Lee, B. Y., Kim, J. & Lee, J. Long-read sequencing infers a mechanism for copy number  
699 variation of template for alternative lengthening of telomeres in a wild *C. elegans* strain.  
700 *MicroPubl Biol* **2022**, (2022).

- 701 27. Stevens, L. *et al.* Chromosome-Level Reference Genomes for Two Strains of  
702 *Caenorhabditis briggsae*: An Improved Platform for Comparative Genomics. *Genome Biol.*  
703 *Evol.* **14**, (2022).
- 704 28. Crombie, T. A. *et al.* Deep sampling of Hawaiian *Caenorhabditis elegans* reveals high  
705 genetic diversity and admixture with global populations. *Elife* **8**, e50465 (2019).
- 706 29. Crombie, T. A. *et al.* Local adaptation and spatiotemporal patterns of genetic diversity  
707 revealed by repeated sampling of *Caenorhabditis elegans* across the Hawaiian Islands.  
708 *Mol. Ecol.* **31**, 2327–2347 (2022).
- 709 30. Andersen, E. C. *et al.* Chromosome-scale selective sweeps shape *Caenorhabditis elegans*  
710 genomic diversity. *Nat. Genet.* **44**, 285–290 (2012).
- 711 31. Li, W. & Shaw, J. E. A variant Tc4 transposable element in the nematode *C. elegans* could  
712 encode a novel protein. *Nucleic Acids Res.* **21**, 59–67 (1993).
- 713 32. Jeong, D.-E. *et al.* DNA polymerase diversity reveals multiple incursions of Polintons  
714 during nematode evolution. *bioRxiv* 2023.08.22.554363 (2023)  
715 doi:10.1101/2023.08.22.554363.
- 716 33. Sayers, E. W. *et al.* Database resources of the national center for biotechnology  
717 information. *Nucleic Acids Res.* **50**, D20–D26 (2022).
- 718 34. Cook, D. E., Zdraljevic, S., Roberts, J. P. & Andersen, E. C. CeNDR, the *Caenorhabditis*  
719 *elegans* natural diversity resource. *Nucleic Acids Res.* **45**, D650–D657 (2017).
- 720 35. Duchaine, T. F. *et al.* Functional proteomics reveals the biochemical niche of *C. elegans*  
721 DCR-1 in multiple small-RNA-mediated pathways. *Cell* **124**, 343–354 (2006).
- 722 36. Claycomb, J. M. *et al.* The Argonaute CSR-1 and its 22G-RNA cofactors are required for  
723 holocentric chromosome segregation. *Cell* **139**, 123–134 (2009).
- 724 37. Gu, W. *et al.* Distinct argonaute-mediated 22G-RNA pathways direct genome surveillance

- 725 in the *C. elegans* germline. *Mol. Cell* **36**, 231–244 (2009).
- 726 38. Conine, C. C. *et al.* Argonautes ALG-3 and ALG-4 are required for spermatogenesis-  
727 specific 26G-RNAs and thermotolerant sperm in *Caenorhabditis elegans*. *Proc. Natl. Acad.*  
728 *Sci. U. S. A.* **107**, 3588–3593 (2010).
- 729 39. Tijsterman, M., Okihara, K. L., Thijssen, K. & Plasterk, R. H. A. PPW-1, a PAZ/PIWI protein  
730 required for efficient germline RNAi, is defective in a natural isolate of *C. elegans*. *Curr.*  
731 *Biol.* **12**, 1535–1540 (2002).
- 732 40. Pollard, D. A. & Rockman, M. V. Resistance to germline RNA interference in a  
733 *Caenorhabditis elegans* wild isolate exhibits complexity and nonadditivity. *G3* **3**, 941–947  
734 (2013).
- 735 41. Ashe, A. *et al.* A deletion polymorphism in the *Caenorhabditis elegans* RIG-I homolog  
736 disables viral RNA dicing and antiviral immunity. *Elife* **2**, e00994 (2013).
- 737 42. Chou, H. T. *et al.* Diversification of small RNA pathways underlies germline RNAi  
738 incompetence in wild *C. elegans* strains. *Genetics* (2023) doi:10.1093/genetics/iyad191.
- 739 43. Guan, Y. *et al.* The MOV10 RNA helicase is a dosage-dependent host restriction factor for  
740 LINE1 retrotransposition in mice. *PLoS Genet.* **19**, e1010566 (2023).
- 741 44. Newman, M. A. *et al.* The surveillance of pre-mRNA splicing is an early step in *C. elegans*  
742 RNAi of endogenous genes. *Genes Dev.* **32**, 670–681 (2018).
- 743 45. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 744 46. Widmayer, S. J., Evans, K. S., Zdraljevic, S. & Andersen, E. C. Evaluating the power and  
745 limitations of genome-wide association studies in *Caenorhabditis elegans*. *G3* **12**, (2022).
- 746 47. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide  
747 complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 748 48. Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-scale



- 749 data. *Nat. Genet.* **51**, 1749–1755 (2019).
- 750 49. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol.* **5**,  
751 R12 (2004).
- 752 50. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic  
753 features. *Bioinformatics* **26**, 841–842 (2010).
- 754 51. Okonechnikov, K., Golosova, O., Fursov, M. & UGENE team. Unipro UGENE: a unified  
755 bioinformatics toolkit. *Bioinformatics* **28**, 1166–1167 (2012).
- 756 52. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421  
757 (2009).
- 758 53. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping  
759 and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**,  
760 2987–2993 (2011).
- 761 54. Ortiz, E. M. *vcf2phylip v2.0: convert a VCF matrix into several matrix formats for*  
762 *phylogenetic analysis.* (2019). doi:10.5281/zenodo.2540861.
- 763 55. Schliep, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592–593 (2011).
- 764 56. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T. Ggtree : An r package for visualization  
765 and annotation of phylogenetic trees with their covariates and other associated data.  
766 *Methods Ecol. Evol.* **8**, 28–36 (2017).
- 767 57. Huson, D. H. & Bryant, D. Application of phylogenetic networks in evolutionary studies.  
768 *Mol. Biol. Evol.* **23**, 254–267 (2006).

769

## Extended data figures and tables

770

771

### **Transposon-mediated genic rearrangements underlie variation in small RNA pathways**

772

773

Gaotian Zhang<sup>1</sup>, Marie-Anne Félix<sup>1</sup>, and Erik C. Andersen<sup>2</sup>

774

775

1. Institut de Biologie de l'École Normale Supérieure, Paris, Île-de-France, France

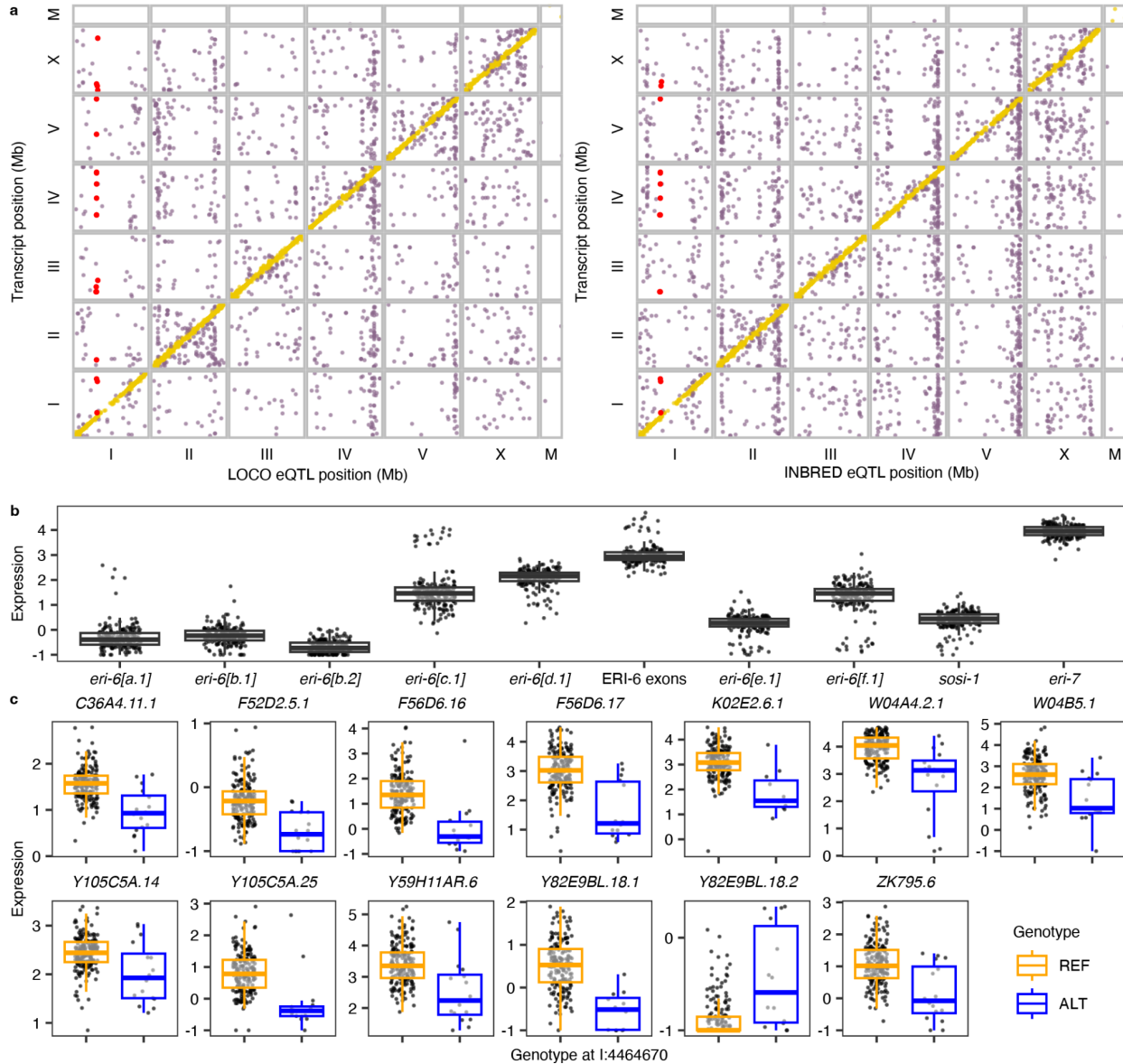
776

2. Biology Department, Johns Hopkins University, Baltimore, MD, USA

Transcript	Gene	Genomic position	Biotype	eQTL method	Distant eQTL peak position	cM bin on Chr I	Candidate variants in <i>eri-6</i>	Known targets of ERI-6/7 dependent siRNAs	With known orthologs in other nematodes
<i>C36A4.11.1</i>	<i>C36A4.11</i>	III: 3840579-3841248	protein-coding	LOCO	I:4696683	21	4464670		No
<i>F52D2.5.1</i>	<i>F52D2.5</i>	X: 1973511-1974799	protein-coding	LOCO	I:4443624	21	4464670		No
<i>F56D6.16</i>	<i>F56D6.16</i>	IV: 3898478-3899806	pseudogene	INBRED, LOCO	I:4316016, I:4434616	20.5, 21	4464670		No
<i>F56D6.17</i>	<i>F56D6.17</i>	IV: 3895963-3897997	pseudogene	INBRED, LOCO	I:4443624, I:4443624	21, 21	4464670		No
<i>K02E2.6.1</i>	<i>K02E2.6</i>	V: 20380561-20382036	protein-coding	INBRED, LOCO	I:4434616, I:4434616	21, 21	4464670	Yes	Yes
<i>W04A4.2.1</i>	<i>W04A4.2</i>	I:13682137-13684547	protein-coding	INBRED, LOCO	I:4443624, I:4443624	21, 21	4464670		Yes
<i>W04B5.1</i>	<i>W04B5.1</i>	III: 2428204-2429430	pseudogene	LOCO	I:4370151	20.5	4464670		No
<i>Y105C5A.14</i>	<i>Y105C5A.14</i>	IV:15646970-15648298	pseudogene	INBRED, LOCO	I:4443624, I:4443624	21, 21	4464670	Yes	No
<i>Y105C5A.25</i>	<i>Y105C5A.25</i>	IV:15860221-15861442	pseudogene	INBRED, LOCO	I:4446767, I:4446767	21, 21	4464670	Yes	No
<i>Y59H11AR.6</i>	<i>Y59H11AR.6</i>	IV: 8597762-8598727	pseudogene	INBRED, LOCO	I:4443624, I:4443624	21, 21	4464670	Yes	No
<i>Y82E9BL.18.1</i>	<i>Y82E9BL.18</i>	III: 1328590-1331066	protein-coding	INBRED, LOCO	I:4316016, I:4316016	20.5, 20.5	4464670		Yes
<i>Y82E9BL.18.2</i>	<i>Y82E9BL.18</i>	III: 1326791-1331066	protein-coding	INBRED, LOCO	I:4434616, I:4434616	21,21	4464670		Yes
<i>ZK795.6</i>	<i>ZK795.6</i>	IV: 12555311-12556482	pseudogene	INBRED, LOCO	I:4443624, I:4443624	21,21	4464670		No

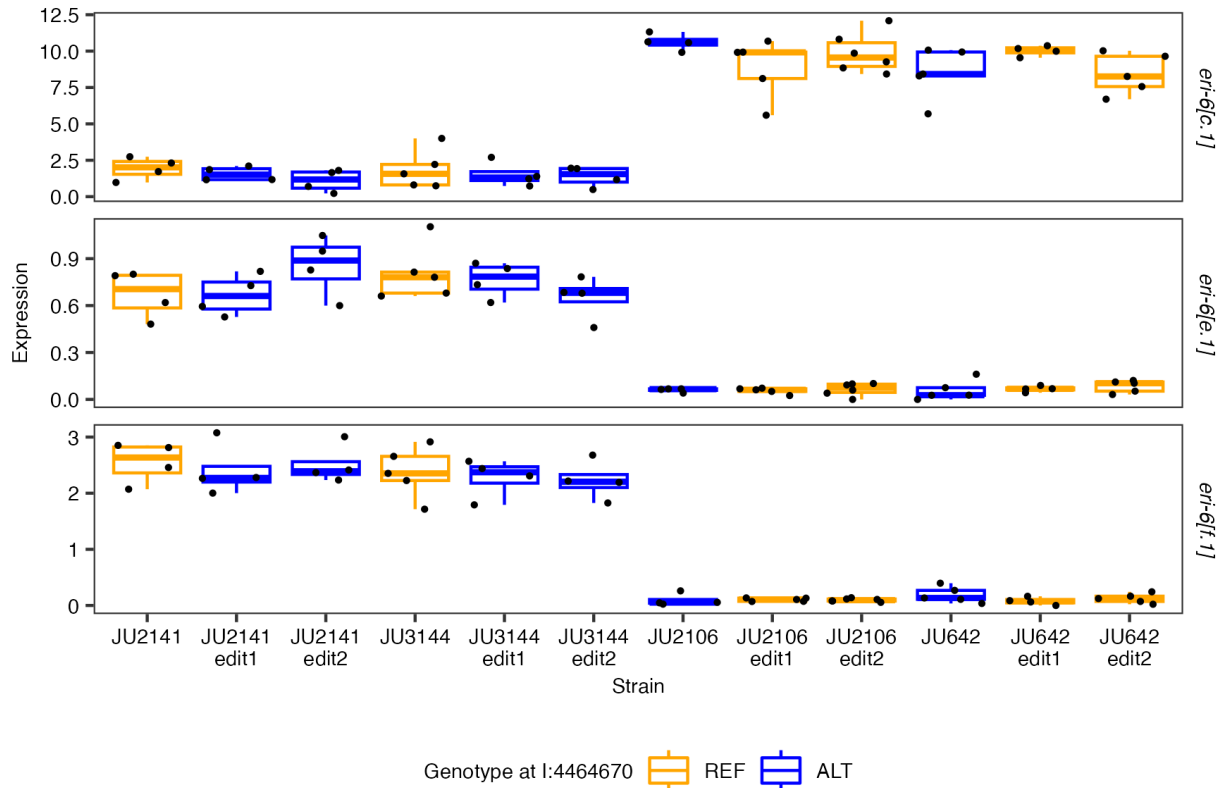
777 **Extended Data Table 1 | Transcript expression traits associated with *eri-6*.**

778 In Fig. 1e, we showed Manhattan plots for the ten traits identified with distant eQTL using both  
779 INBRED and LOCO methods. In Extended Data Fig. 1c, we showed the phenotype by genotype  
780 plots for all of ten gene expression traits at the top candidate variant. Because of the negative  
781 correlation in expression between the two transcripts *Y82E9BL.18.1* and *Y82E9BL.18.2* of the  
782 gene *Y82E9BL.18*, only one of these transcripts was depicted in Fig. 4c for clarity purposes.

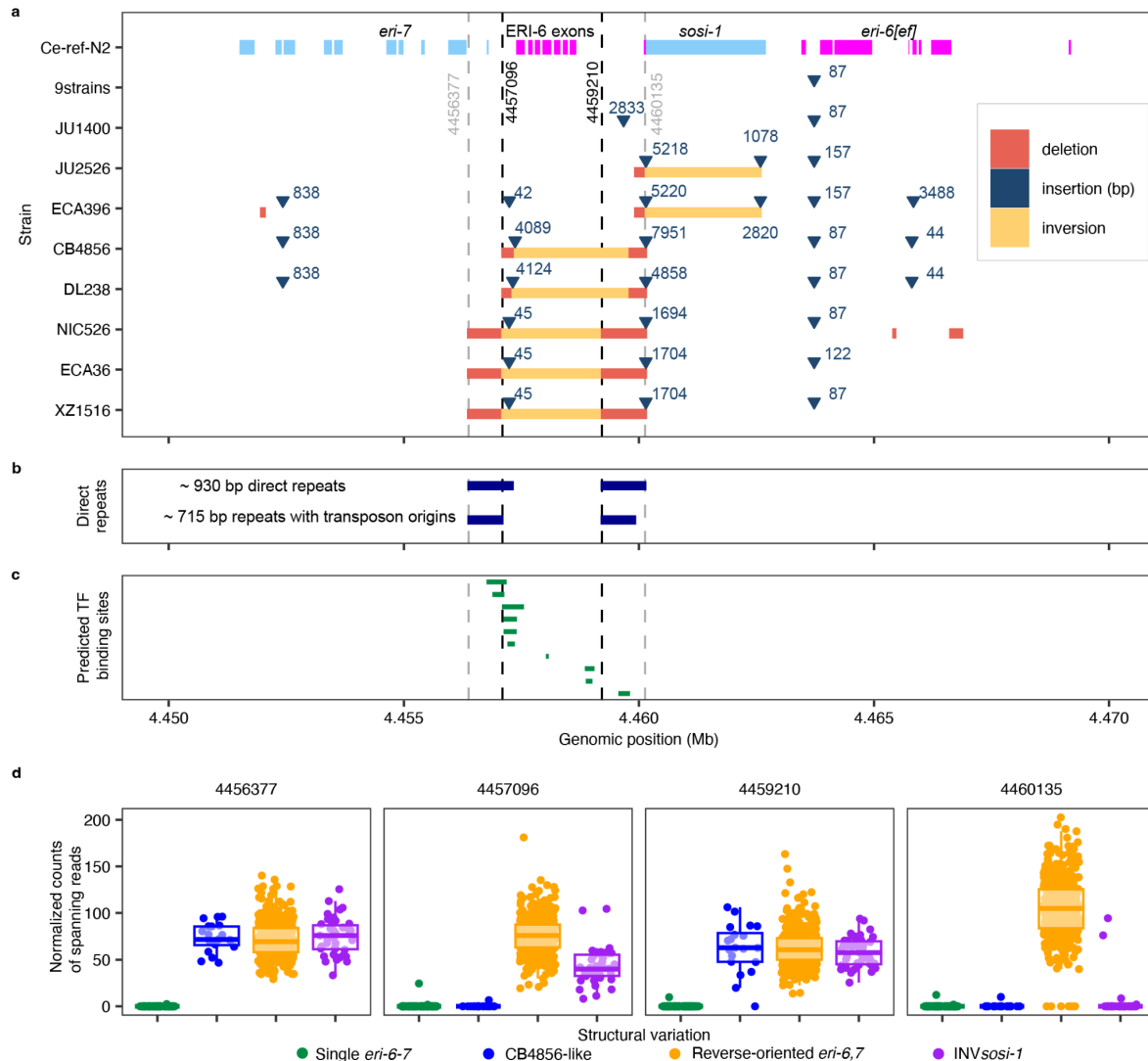


783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793

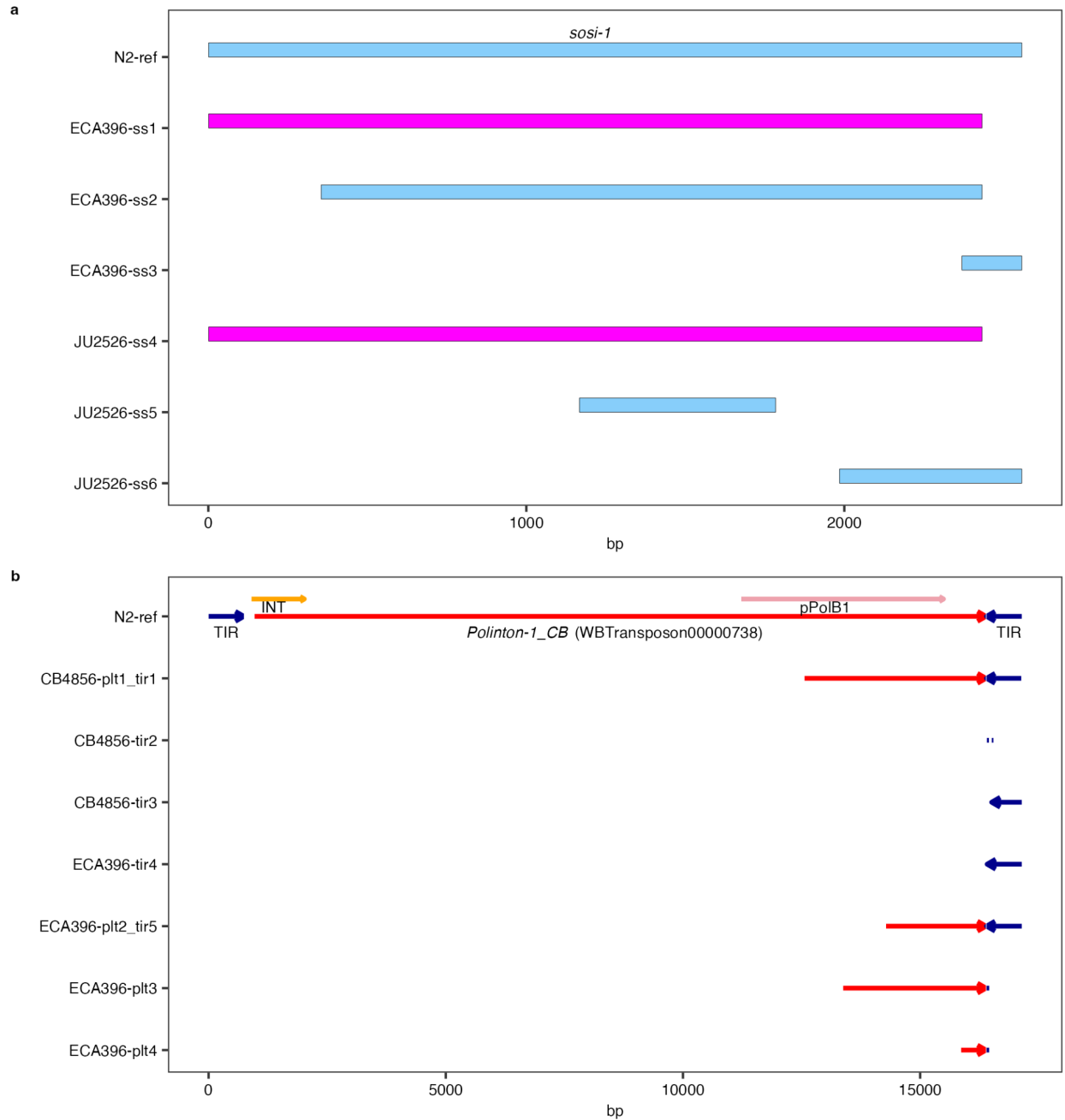
**Extended Data Fig. 1 | Expression QTL maps and expression variation of *eri-6* related transcripts.** **a**, Expression QTL maps using LOCO and INBRED methods. Each point represents an eQTL with its position on the x-axis and the genomic position of the transcript on the y-axis. Local and distant eQTL are colored gold and purple, respectively. Red points represent distant eQTL associated with the *eri-6/7* locus. **b**, Tukey box plots showing expression ( $-\log_2(\text{normalized TPM}+0.5)$ ) variation of ten transcripts at the *eri-6/7* locus. **c**, Tukey box plots showing expression variation of 13 transcripts across the genome between strains with the reference (REF) or alternative (ALT) alleles at the SNV of 4,464,670 bp on chromosome I. **b,c**, Each point represents a strain. Box edges denote the 25th and 75th quantiles of the data; and whiskers represent 1.5x the interquartile range.



794  
 795 **Extended Data Fig. 2 | The SNV candidate cannot explain expression variation in *eri-6*.**  
 796 Expression variation of *eri-6*[*c.1,e.1,f.1*] among four wild *C. elegans* strains (JU2141, JU3144,  
 797 JU2106, JU642) and their eight mutant strains at *eri-6*[*e*] (I: 4,464,670) using CRISPR-Cas9-  
 798 mediated genome editing as previously described<sup>1-3</sup>. The guide RNAs crECA163  
 799 (GCTGTGCCACGATCGGAGTA) (Synthego, CA, USA) was used for the editing. The homologous  
 800 recombination templates crECA162  
 801 (tgtcatttgatcccgcctcggcattttcaacgatgacgaaaagtcttctaacaatctcgaatTaccctactccgatcgtggcacagctc  
 802 aatagcctcaaagagctgaaactgaaagtagccg) and crECA164  
 803 (tgtcatttgatcccgcctcggcattttcaacgatgacgaaaagtcttctaacaatctcgaatGaccctactccgatcgtggcacagctc  
 804 aatagcctcaaagagctgaaactgaaagtagccg) (IDT, IL, USA) were used for wild strains with the  
 805 reference (REF) and alternative (ALT) alleles at the target site, respectively. Genotypes of F2  
 806 progeny were detected with primers oECA1989 (GGTGGTGGCAGCGCATCTAGTC) and  
 807 oECA1990 (GCTCCCGAATGTAGCCACCGA) using PCR and Sanger sequencing. Edit1 and  
 808 Edit2 are two independent edits in each of the four backgrounds. Each point represents a  
 809 biological replicate. Transcriptomes in synchronized young adult stage animals of each replicate  
 810 were measured by RNA sequencing and quantified as previously described<sup>4</sup> (also see details in  
 811 Methods).

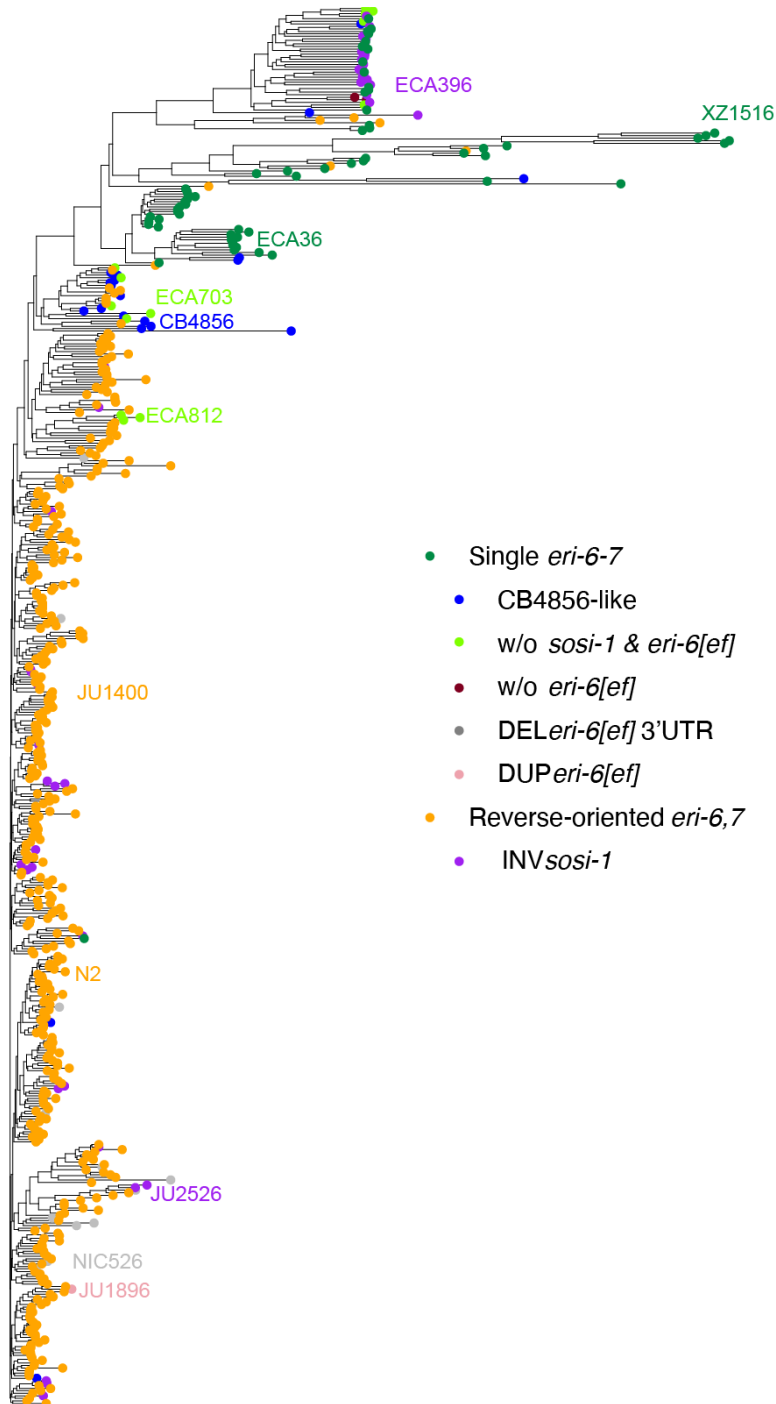


812  
 813 **Extended Data Fig. 3 | Structural variation in the *eri-6/7* region.** **a**, Large deletions ( $\geq 38$  bp),  
 814 insertions ( $\geq 42$  bp), and inversions of 17 wild strains with PacBio assemblies to the reference N2  
 815 genome in the *eri-6/7* region are represented by red rectangles, dark blue triangles, and yellow  
 816 rectangles, respectively. Sizes of the insertions are indicated in bp. Exons of *eri-6/7* and *sosi-1*  
 817 are plotted as rectangles on top and are colored magenta and light blue for plus and minus coding  
 818 strands, respectively. Nine strains (DL226, EG4725, JU310, JU1395, JU2600, MY2147, MY2693,  
 819 NIC2, QX1794) with highly identical sequences in this region were represented together. The only  
 820 local structural difference in these nine strains as compared to N2 is a shared 87-bp insertion  
 821 upstream of *eri-6[e]*. **b**, Ranges of the  $\sim 930$  bp direct repeats<sup>5</sup> and the  $\sim 715$  bp parts with *Polinton*  
 822 origins are indicated as blue rectangles, respectively (Supplementary Table 4). Dashed vertical  
 823 gray and black lines indicate outside boundaries of direct repeats and break points of inversions  
 824 (defined by comparison between the strains XZ1516 and the reference N2). **c**, Predicted  
 825 transcription factor (TF) binding sites<sup>6</sup> within the repeat regions are indicated as green rectangles.  
 826 **d**, Number of reads spanning 20 bp surrounding each boundary/break-point position was counted  
 827 and percent of this count normalized by the mean coverage per bp in the *eri-6/7* locus for each  
 828 strain was plotted on the y-axis against the structural variation on the x-axis. High counts of reads  
 829 spanning the inversion breakpoints indicate Watson ERI-6 exons and direct repeats as in the  
 830 reference genome, except for the CB4856-like strains.



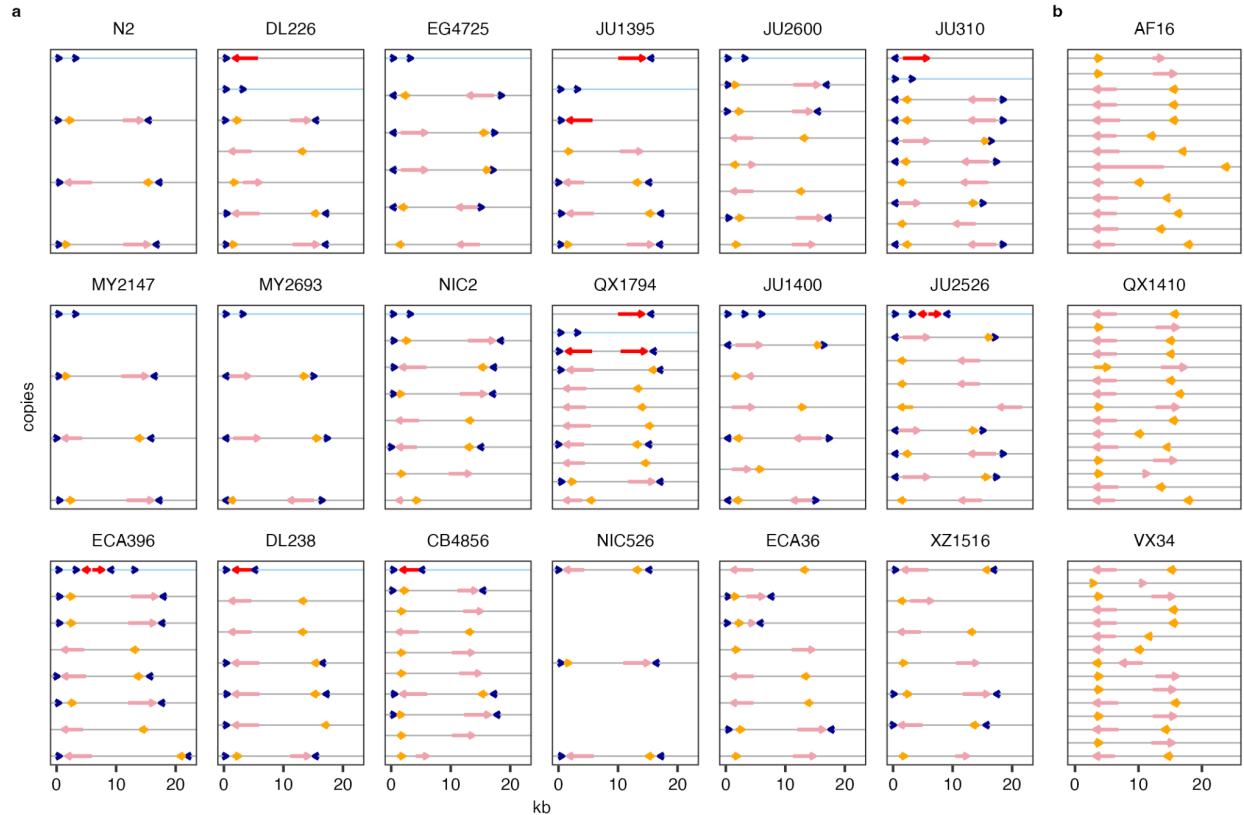
831  
832 **Extended Data Fig. 4 | Alignment of *Polinton* and *sosi-1* fragments to the reference.**  
833 Sequence alignments of fragments (“ss1-6”, “tir1-5”, and “plt1-4”) indicated in **Fig. 2** to *sosi-1*  
834 **(a)** and the largest *Polinton* remnant (*Polinton-1\_CB*, WBTransposon00000738, as a red arrow  
835 on top) **(b)** in the reference N2 genome are shown. Positions of the retroviral-like-element  
836 integrase (INT) and the protein-primed DNA polymerase B genes (pPolB1) are indicated as  
837 orange and pink arrows, respectively. Note that sequences of segments CB4856-tir2, tir3, and  
838 ECA396-tir4 can also be aligned to the terminal inverted repeat (TIR, blue arrows) on the left.



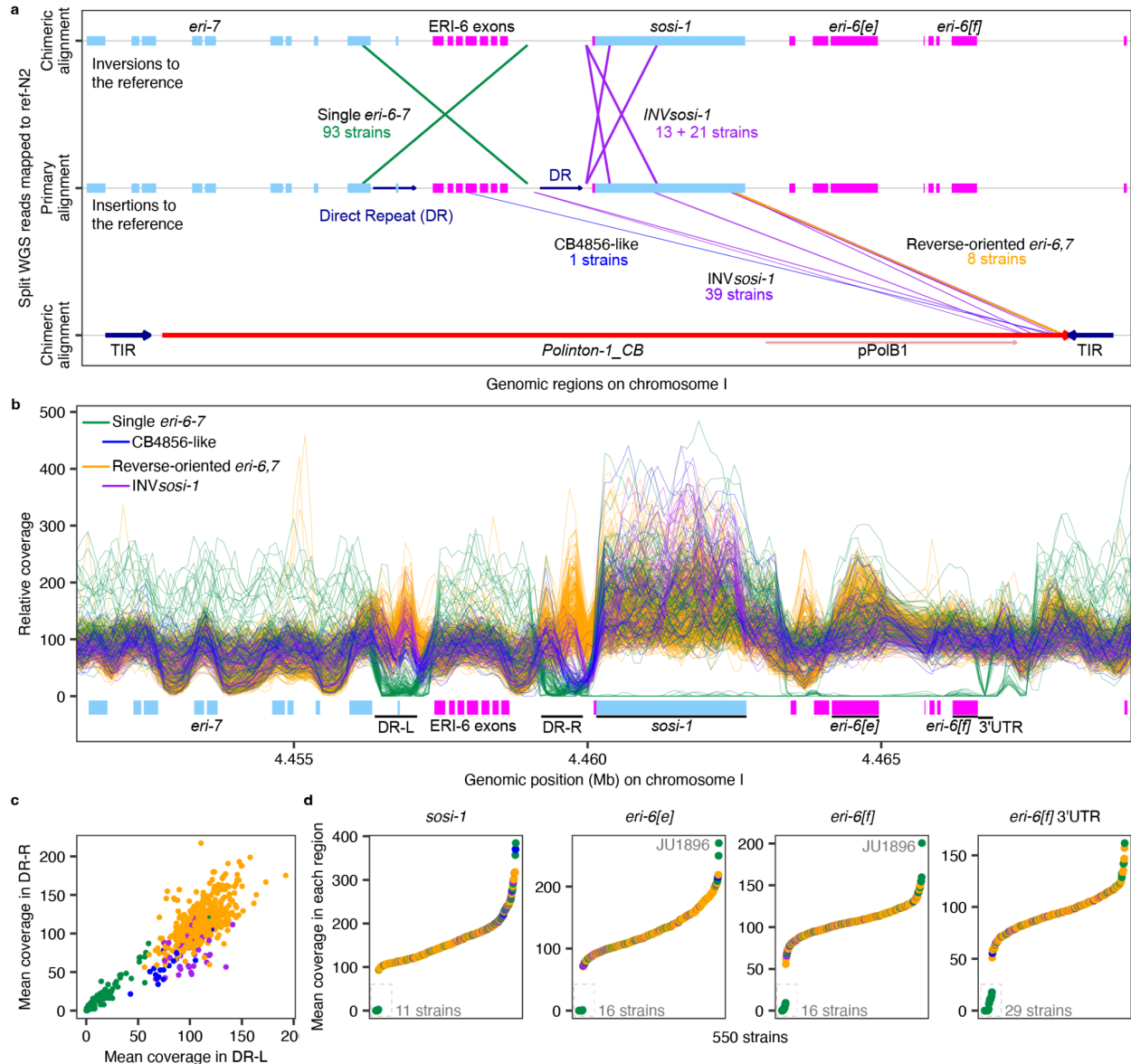


839

840 **Extended Data Fig. 5 | Genetic relatedness of 550 wild *C. elegans* strains.** Genetic  
841 relatedness tree using 1,199,944 biallelic SNVs throughout the genome. Recombination occurs  
842 within the species so this tree only represents overall relatedness. Each point represents a strain  
843 and is colored by its structural variation in the *eri-6/7* region. The 19 strains clustered with  
844 ECA396 in Fig. 3 are also clustered with ECA396 on the same branch.



**Extended Data Fig. 6 | Presence of TIRs and *Polintons* in *C. elegans* and *C. briggsae* strains.** *Polinton* insertions were identified in 18 *C. elegans* strains (a) and three *C. briggsae* strains (b) by requiring the presence of both pPoIB1 (pink arrows) and INT (orange arrows) within 20 kb. Blue arrows represent TIRs. Red arrows represent pPoIB that was found close to TIRs but without nearby INT. Direction of arrows represents orientations. Horizontal lines indicate different copies in the genomes, with the blue lines highlighting the copies in the *eri-6/7* locus. All the identified *Polinton* insertions and TIRs were plotted.



853

854

855

856

857

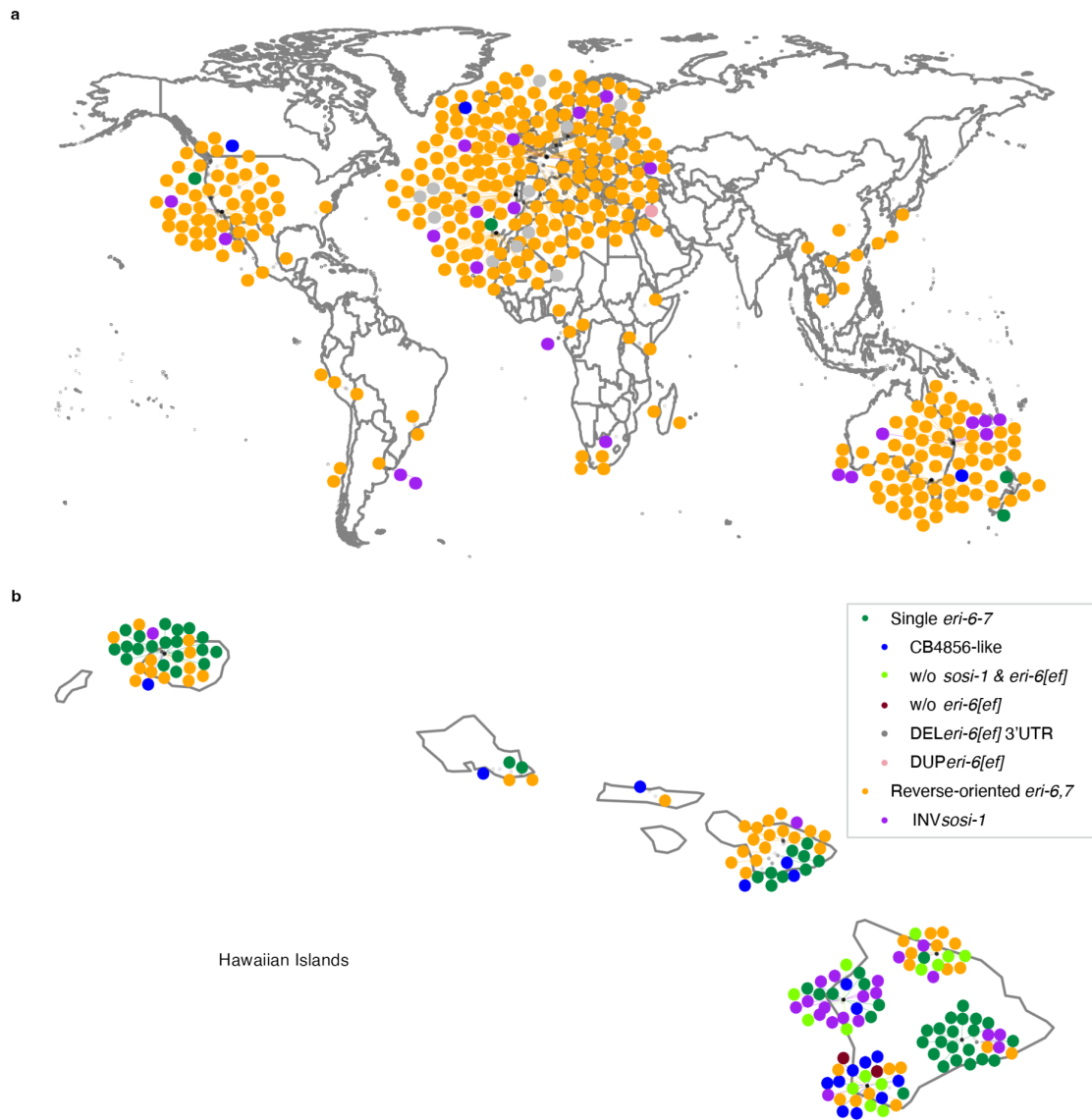
858

859

860

861

**Extended Data Fig. 7 | Inversions and other structural variants within the *eri-6/7* locus. a,** Inversions and *Polinton-1\_CB* insertions within the *eri-6/7* locus are represented by lines that connect positions of the primary and chimeric alignments of split reads (See Methods). **b,** Sliding windows on normalized coverages per bp with a 200-kb window size and a 100-bp step size in the *eri-6/7* locus of 550 strains. **c,d** Mean coverages in each strain within the two direct repeats (**c**) and four other regions (**d**) indicated in **b** were shown. Each line (**b**) / point (**c,d**) represents one strain and is colored by structural variants indicated in **b**. Extreme low coverages indicate deletions and extreme high coverages indicate duplications compared to the reference genome.



862  
863 **Extended Data Fig. 8 | Geographical distribution of wild *C. elegans*.** Geographical distribution  
864 of 550 strains worldwide (a) and detailed on the Hawaiian Islands (b). Each point represents a  
865 strain and is colored by its structural variation in the *eri-6/7* region.

866 **References**

- 867 1. Kim, H. *et al.* A co-CRISPR strategy for efficient genome editing in *Caenorhabditis*  
868 *elegans*. *Genetics* **197**, 1069–1080 (2014).
- 869 2. Paix, A., Folkmann, A., Rasoloson, D. & Seydoux, G. High Efficiency, Homology-Directed  
870 Genome Editing in *Caenorhabditis elegans* Using CRISPR-Cas9 Ribonucleoprotein  
871 Complexes. *Genetics* **201**, 47–54 (2015).
- 872 3. Evans, K. S. *et al.* Natural variation in the sequestosome-related gene, *sqst-5*, underlies  
873 zinc homeostasis in *Caenorhabditis elegans*. *PLoS Genet.* **16**, e1008986 (2020).
- 874 4. Zhang, G., Roberto, N. M., Lee, D., Hahnel, S. R. & Andersen, E. C. The impact of species-  
875 wide gene expression variation on *Caenorhabditis elegans* complex traits. *Nat. Commun.*  
876 **13**, 1–13 (2022).
- 877 5. Fischer, S. E. J., Butler, M. D., Pan, Q. & Ruvkun, G. Trans-splicing in *C. elegans*  
878 generates the negative RNAi regulator ERI-6/7. *Nature* **455**, 491–496 (2008).
- 879 6. Harris, T. W. *et al.* WormBase: a modern Model Organism Information Resource. *Nucleic*  
880 *Acids Res.* **48**, D762–D767 (2020).