

1 **An atlas of gene expression variation across the *Caenorhabditis elegans* species**

2

3 Gaotian Zhang¹, Nicole M. Roberto¹, Daehan Lee¹, Steffen R. Hahnel¹, and Erik C.
4 Andersen^{1,*}

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

7

8 ORCID IDs:

9 0000-0001-6468-1341 (G.Z.)

10 0000-0002-0546-8484 (D.L.)

11 0000-0001-8848-0691 (S.R.H.)

12 0000-0003-0229-9651 (E.C.A.)

13

14 *Corresponding author:

15 Erik C. Andersen

16 Department of Molecular Biosciences

17 Northwestern University

18 4619 Silverman Hall

19 2205 Tech Drive

20 Evanston, IL 60208

21 E-mail: erik.andersen@northwestern.edu

22 **Abstract**

23 Phenotypic variation in diverse organism-level traits have been studied in
24 *Caenorhabditis elegans* wild strains, but differences in gene expression and the
25 underlying variation in regulatory mechanisms are largely unknown. Here, we use natural
26 variation in gene expression to connect genetic variants to differences in organismal-
27 level traits, including drug and toxicant responses. We performed transcriptomic analysis
28 on 207 genetically distinct *C. elegans* wild strains to study natural regulatory variation of
29 gene expression. Using this massive dataset, we performed genome-wide association
30 mappings to investigate the genetic basis underlying gene expression variation and
31 revealed complex genetic architectures. We found a large collection of hotspots
32 enriched for expression quantitative trait loci across the genome. We further used
33 mediation analysis to understand how gene expression variation could underlie
34 organism-level phenotypic variation for a variety of complex traits. These results reveal
35 the natural diversity in gene expression and possible regulatory mechanisms in this
36 keystone model organism, highlighting the promise of gene expression variation in
37 shaping phenotypic diversity.

38

39

40 Introduction

41 Quantitative genetic mapping approaches, such as genome-wide association
42 (GWA) and linkage mapping, have been used in a variety of organisms to disentangle the
43 underlying genetic basis of gene expression variation by considering the expression level
44 of each gene as a quantitative trait^{1–9}. Expression quantitative trait loci (eQTL) affecting
45 gene expression are often classified into local eQTL (located close to the genes that they
46 influence) and distant eQTL (located further away from the genes that they influence)^{10,11}.
47 Local eQTL are abundant in the genome. For example, over half the genes in yeast and
48 94.7% of all protein-coding genes in human tissues are hypothesized to have associated
49 local eQTL^{7,8}. Genetic variants underlying local eQTL might influence the expression of
50 a specific gene by affecting transcription factor binding sites, chromatin accessibility,
51 other promoter elements, enhancers, or other factors at post-transcriptional levels¹².
52 Genes encoding diffusible factors, such as transcription factors, chromatin cofactors,
53 and RNAs, are often considered the most likely genes to underlie distant eQTL. Distant
54 eQTL hotspots in several species have been suggested to account for the variation in
55 expression of many genes located throughout the genome^{2,3,7,9,13}. Although a substantial
56 amount of eQTL have been identified in different species, it is still largely unknown how
57 gene expression variation relates to organism-level phenotypic differences.

58 The nematode *Caenorhabditis elegans* is a powerful model to study the genetic
59 basis of natural variation in diverse quantitative traits^{14–16}. Genome-wide gene expression
60 variation in different developmental stages and various conditions at the whole-organism
61 or cellular resolution have been discovered and thousands of eQTL have been identified
62 in several studies over the past two decades^{3,9,17–23}. However, most of these studies used
63 two-parent recombinant inbred lines derived from crosses of the laboratory-adapted
64 reference strain, N2, and the genetically diverse Hawaiian strain, CB4856. Consequently,
65 the observed variation in gene expression and their identified eQTL were limited to the
66 differences among a small number of *C. elegans* strains and only revealed a tiny fraction
67 of the natural diversity of gene expression and regulatory mechanisms in this species.
68 The *C. elegans* Natural Diversity Resource (CeNDR) has a collection of 540 genetically
69 distinct wild *C. elegans* strains^{16,24,25}. Variation in diverse organism-level phenotypes has

70 been observed among these wild strains, and many underlying QTL, quantitative trait
71 genes (QTGs), and quantitative trait variants (QTVs) have been identified using GWA
72 mappings^{15,16}. Therefore, a genome-wide analysis could improve our understanding of
73 the role of gene regulation in shaping organism-level phenotypic diversity, adaptation,
74 and evolution of *C. elegans*.

75 Here, we investigated the natural variation in gene expression of 207 genetically
76 distinct *C. elegans* wild strains by performing bulk mRNA sequencing on synchronized
77 young adult hermaphrodites. We used GWA mappings to identify 6,545 eQTL associated
78 with variation in expression of 5,291 transcripts of 4,520 genes. We found that local eQTL
79 explained most of the narrow-sense heritability and showed larger effects on expression
80 variation than distant eQTL. We identified 67 hotspots that comprise 1,828 distant eQTL
81 across the *C. elegans* genome. We further found a diverse collection of potential
82 regulatory mechanisms that underlie these distant eQTL hotspots. Additionally, we
83 applied mediation analysis to gene expression and other quantitative trait variation data
84 to elucidate putative mechanisms that can play a role in organism-level trait variation.
85 Our results provide an unprecedented resource of transcriptome profiles and genome-
86 wide regulatory regions that facilitate future studies. Furthermore, we demonstrate
87 efficient methods to locate causal genes that underlie mechanisms of organism-level
88 trait differences across the *C. elegans* species.

89 Results

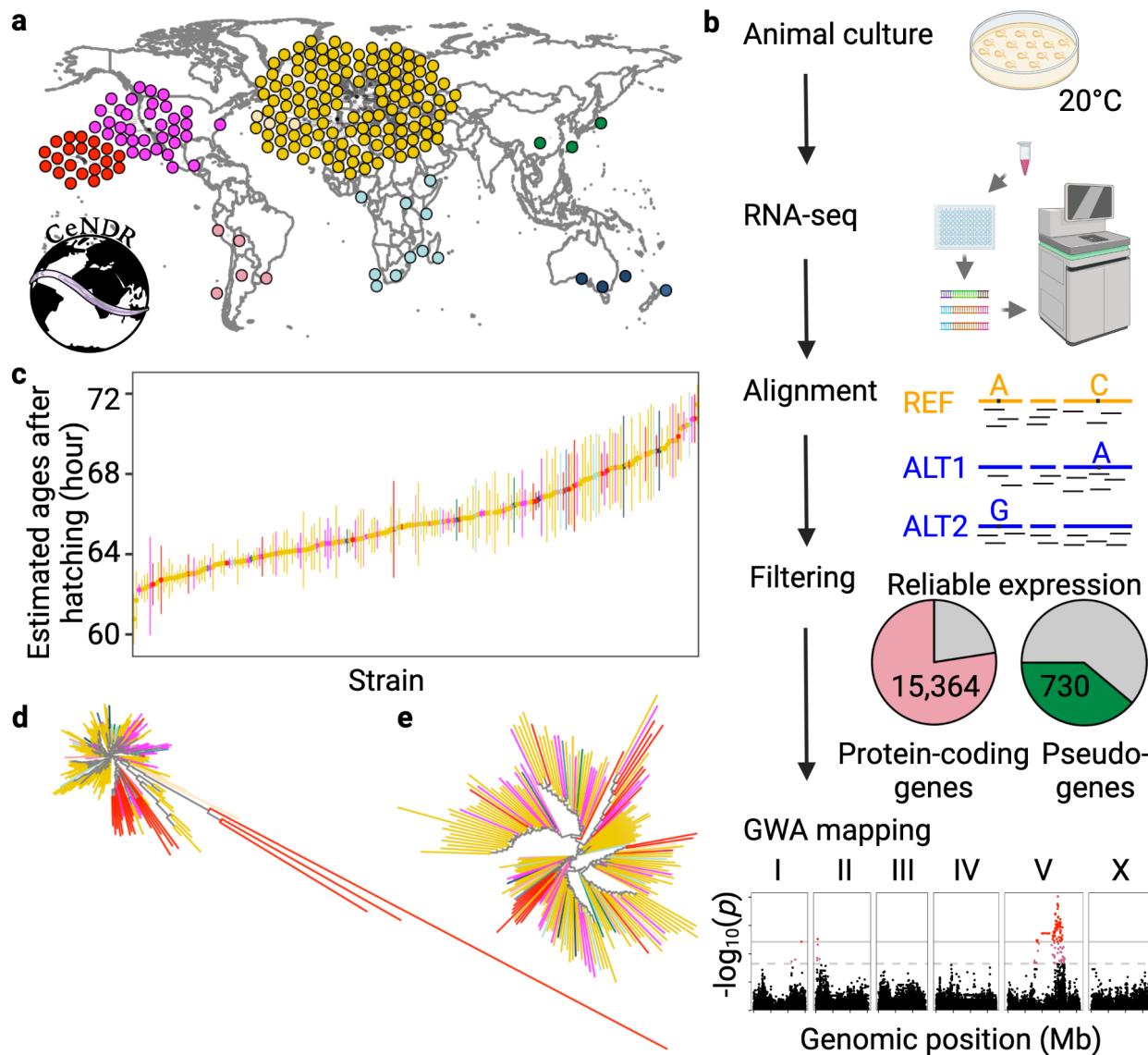
90 Transcriptome profiles of 207 wild *C. elegans* strains

91 We obtained 207 wild *C. elegans* strains from CeNDR²⁵ (Fig. 1a). We grew and
92 harvested synchronized populations of each strain at the young adult stage in
93 independently grown and prepared biological replicates (Fig. 1b). We performed bulk
94 RNA sequencing to measure expression levels and aligned reads to strain-specific
95 transcriptomes (Fig. 1b, Supplementary Fig. 1, Supplementary Data 1). We focused on
96 protein-coding genes and pseudogenes and filtered out those genes with low and/or
97 rarely detected expression (See Methods). Because various hyper-divergent regions with

98 extremely high nucleotide diversity were identified in the genomes of wild *C. elegans*
99 strains^{26,27}, RNA sequencing reads might be poorly aligned and expression abundances
100 might be underestimated for genes in these regions. For each strain, we filtered out
101 transcripts that fell into the known hyper-divergent regions. We also dropped outlier
102 samples by comparing sample-to-sample expression distances (Supplementary Fig. 1).
103 To further verify the homogeneity of developmental stages of our samples, we evaluated
104 the age of each sample when they were harvested using our expression data and
105 published time-series expression data²⁸. We inferred that our animals fit an expected
106 developmental age of 60 to 72 hours post hatching (Fig. 1c), during which time the animal
107 is in the young adult stage. Because we harvested the animals at the first embryo-laying
108 event, the age estimation also reflects natural variation in the duration from hatching to
109 the beginning of embryo-laying of wild *C. elegans*. In summary, we obtained reliable
110 expression abundance measurements for 25,849 transcripts from 16,094 genes (15,364
111 protein-coding genes and 730 pseudogenes) in 561 samples of 207 *C. elegans* strains
112 (Fig. 1b, Supplementary Fig. 1, Supplementary Data 1), which we used for downstream
113 analyses.

114 *C. elegans* geographic population structure has been observed previously^{24,27,29}.
115 Wild strains from Hawaii and other regions in the Pacific Rim harbor high genetic diversity
116 and group into distinct clusters using genetic relatedness and principal component
117 analysis^{24,27,29}. Other strains that were isolated largely from Europe have relatively low
118 genetic diversity because of the recent selective sweeps^{24,27,29}. Similar to the previous
119 results, the 207 strains were classified into three major groups consisting of strains from
120 Hawaii, the Pacific coast of the United States, and Europe, respectively, in the genetic
121 relatedness tree (Fig. 1d). Three Hawaiian strains are extremely divergent from all other
122 strains. However, a tree constructed using transcriptome data only exhibited weak
123 geographic relationships and no highly divergent strains (Fig. 1e), suggesting stabilizing
124 selection has constrained variation in gene expression.

125
126
127



128

129 **Fig. 1: Overview of species-wide expression analysis in wild *C. elegans*.**

130 **a** Global distribution of 205 of the 207 wild *C. elegans* strains that were obtained from
131 CeNDR and used in this study. Strains are colored by their sampling location continent,
132 except for Hawaiian strains (in red). The two strains missing on the map are lacking
133 sampling locations. **b** Graphic illustration of the workflow to acquire *C. elegans*
134 transcriptome data. Created using BioRender.com. **c** Estimated developmental age (y-
135 axis) of 561 well clustered samples of the 207 wild *C. elegans* strains (x-axis). Strains on
136 the x-axis are sorted by their mean estimated age from two to three biological replicates.
137 Error bars show standard deviation of estimated age among replicates of each strain. **d**,
138 **e** Two Neighbor-joining trees of the 207 *C. elegans* strains using 851,105 biallelic
139 segregating sites (**d**) and expression of 22,268 transcripts (**e**) are shown. Strains in **c**, **d**,
140 **e** are colored as in **a**.

141

142 **Complex regulatory genetic architectures in wild *C. elegans* strains**

143 To estimate the association between gene expression differences and genetic
144 variation, we calculated the broad-sense heritability (H^2) and the narrow-sense
145 heritability (h^2) for each of the 25,849 transcript expression traits. We observed a median
146 H^2 of 0.31 and a median h^2 of 0.06 (Fig. 2a, Supplementary Data 1), indicating strong
147 influences from environmental factors, epistasis, or other stochastic factors on transcript
148 expression variation^{7,30,31}. Nearly 4,000 traits have a h^2 higher than 0.18, indicating a
149 substantial heritable genetic component of the population-wide expression differences.

150 We performed marker-based GWA mappings to investigate the genetic basis of
151 expression variation in the 25,849 transcripts (Supplementary Data 1). We determined
152 the 5% false discovery rate (FDR) significance threshold for eQTL detection by mapping
153 40,000 permuted transcript expression traits using the EMMA algorithm³² and the eigen-
154 decomposition significance (EIGEN) threshold³³ (See Methods). In total, we detected
155 6,545 significant eQTL associated with variation in expression of 5,291 transcripts from
156 4,520 genes (Fig. 2b, Supplementary Data 2). The correlation of h^2 and H^2 among traits
157 with eQTL is much higher than among traits without eQTL (Kendall's τ coefficient, 0.45
158 and 0.27, respectively) (Fig. 2a), indicating major roles of additive genetic variation on
159 expression variation than other genetic factors. Likely because GWA mappings mainly
160 detect QTL that contribute additively to trait variance, eQTL were detected for 71% of
161 the traits with $h^2 > 0.18$, but only 11% of the remaining traits (Fig. 2a).

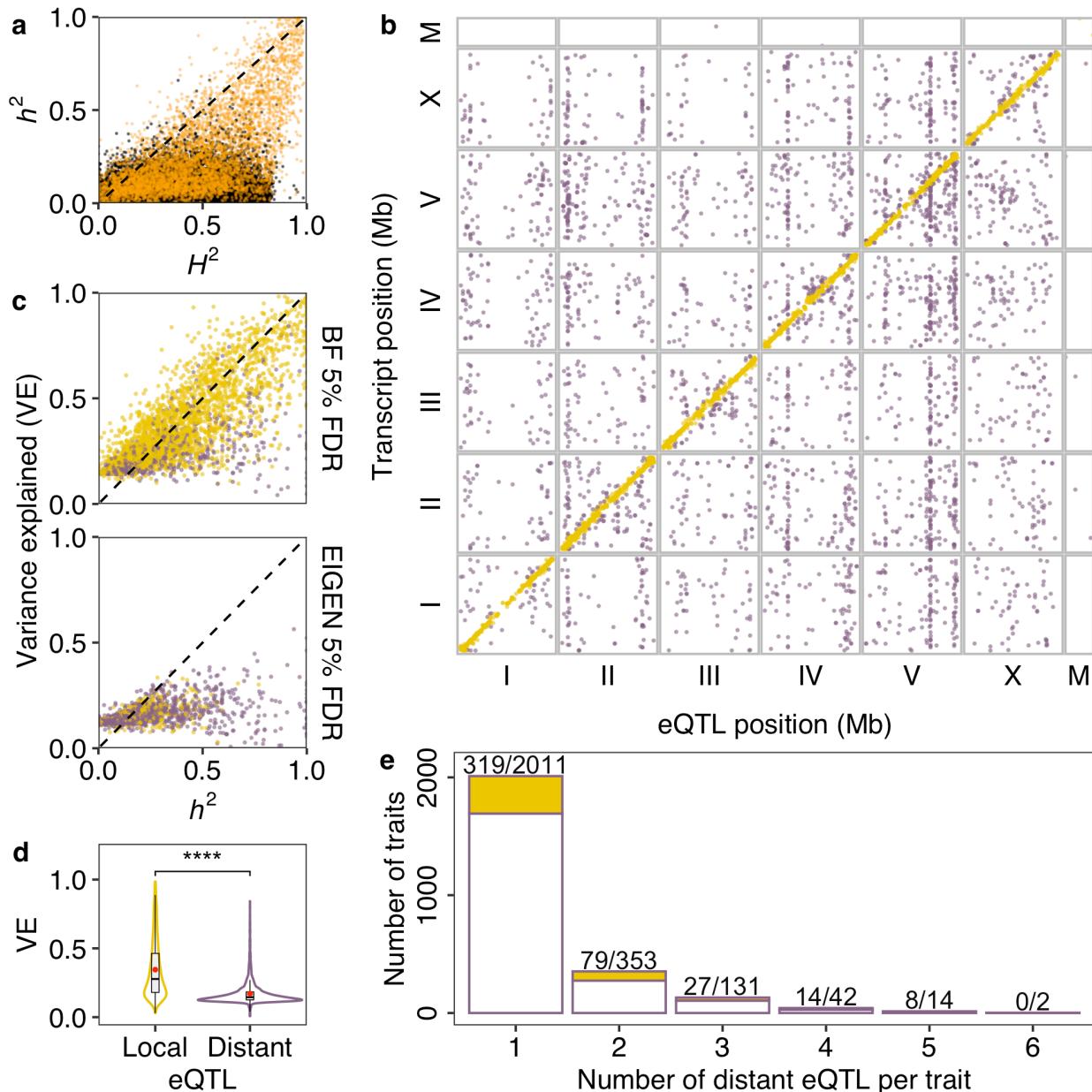
162 In close agreement to previous *C. elegans* eQTL studies using recombinant inbred
163 advanced intercross lines (RIAILs) derived from a cross of the N2 and CB4856 strains^{3,9},
164 eQTL in this study were mostly found on chromosome arms (61%) relative to centers
165 (33%), which is likely related to the genomic distribution of genomic variation (Table 1).
166 Of the 4,520 genes with transcript-level eQTL, we found overrepresentation of
167 nonessential genes (Fisher exact test, odds ratio: 1.18, p -value: 0.001) and
168 underrepresentation of essential genes (Fisher exact test, odds ratio: 0.75, p -value:
169 0.001), suggesting stronger selection against expression variation in essential genes
170 than nonessential genes³⁴. Gene set enrichment analysis (GSEA) on these 4,520 genes
171 showed that proteolysis proteasome-related genes (Fisher Exact Test, Bonferroni FDR

172 corrected $p = 3.76\text{E-}20$), especially genes encoding E3 ligases containing an F-box
173 domain (Fisher Exact Test, Bonferroni FDR corrected $p = 3.73\text{E-}15$), are the most
174 significantly enriched class (Supplementary Fig. 2, Supplementary Data 3). Other
175 significantly enriched gene classes include metabolism (Fisher Exact Test, Bonferroni
176 FDR corrected $p = 2.92\text{E-}12$), stress response (Fisher Exact Test, Bonferroni FDR
177 corrected $p = 7.24\text{E-}12$), and histones (Fisher Exact Test, Bonferroni FDR corrected $p =$
178 $3.23\text{E-}8$). (Supplementary Fig. 2).

179 We classified eQTL located within a two megabase region surrounding each
180 transcript as local eQTL and all other eQTL as distant^{3,9} (Fig. 2b, Table 1, Supplementary
181 Data 2). We identified local eQTL for 3,185 transcripts from 2,655 genes (Fig. 2b, Table
182 1, Supplementary Data 2). The 2,551 local eQTL that passed the Bonferroni 5% FDR
183 threshold explained most of the estimated narrow-sense heritability (Fig. 2c).
184 Additionally, we found 3,360 distant eQTL for 2,553 transcripts from 2,382 genes (Fig.
185 2b, Table 1, Supplementary Data 2). Compared to local eQTL, distant eQTL generally
186 explained significantly lower variance (Fig. 2c, d). We found that local eQTL and up to
187 six distant eQTL could jointly regulate the expression of transcripts (Fig. 2e). Because
188 substantial linkage disequilibrium (LD) is observed within ($r^2 > 0.6$) and between ($r^2 >$
189 0.2) chromosomes in wild *C. elegans* strains^{24,27,35}, we calculated LD among eQTL of
190 each of the 861 transcripts with multiple eQTL. We found low LD among most eQTL,
191 with a median LD of $r^2 = 0.19$ (Supplementary Fig. 3), suggesting complex genetic
192 architectures underlying variation in expression of these transcripts are driven by
193 independent loci.

194

195



196

197 Fig. 2: Expression QTL map of 207 wild *C. elegans* strains.

a Heritability for 25,849 transcript expression traits with (orange) or without (black) detected eQTL. The narrow-sense heritability (h^2 , y-axis) for each trait is plotted against the broad-sense heritability (H^2 , x-axis). **b** The genomic locations of 6,545 eQTL peaks (x-axis) that pass the genome-wide EIGEN 5% FDR threshold are plotted against the genomic locations of the 5,291 transcripts with expression differences (y-axis). Golden points on the diagonal of the map represent local eQTL that colocalize with the transcripts that they influence. Purple points correspond to distant eQTL that are located further away from the transcripts that they influence. **c** The variance explained (VE) by each detected eQTL (y-axis) that passed Bonferroni (BF) 5% FDR or EIGEN 5% FDR threshold for each trait is plotted against the narrow-sense heritability h^2 (x-axis). The

208 dashed lines on the diagonal are shown as visual guides to represent $h^2 = H^2$ (**a**) and $VE = h^2$ (**c**). **d** Comparison of VE between detected local and distant eQTL shown as Violin plots. The mean VE by local or distant eQTL is indicated as red points. Statistical significance was calculated using the Wilcoxon test with $p\text{-value} < 2e-16$ indicated as ***. **e** A histogram showing the number of distant eQTL detected per transcript expression trait. One to six distant eQTL were detected for 2,553 transcript expression traits, of which 447 traits also have one local eQTL. Numbers before slashes (indicated as the golden proportion of each bar) represent the number of traits with a local eQTL in addition to their distant eQTL. Numbers after each slash on top of each bar represent the total number of traits in each category.

218

219 **Table 1: The distribution of eQTL and SNVs.**

220 Genomic domain coordinates were defined previously³⁶. Transcript expression traits
221 and SNVs used for eQTL mappings were listed.

222

Domain	eQTL	Local eQTL	Distant eQTL	Genome	Transcripts	SNVs
Tip	388 (5.93%)	224 (7.03%)	164 (4.88%)	7.37 Mb (7.35%)	1,712 (6.62%)	1,628 (7.76%)
Arm	3,966 (60.60%)	2,027 (63.64%)	1,939 (57.71%)	45.89 Mb (45.76%)	9,503 (36.76%)	12,883 (61.37%)
Center	2,183 (33.35%)	932 (29.26%)	1,251 (37.23%)	47.01 Mb (46.88%)	14,622 (56.57%)	6,429 (30.63%)
MtDNA	8 (0.12%)	2 (0.06%)	6 (0.18%)	0.01 Mb (0.01%)	12 (0.05%)	51 (0.24%)
Total	6,545	3,185	3,360	100.29 Mb	25,849	20,991

223

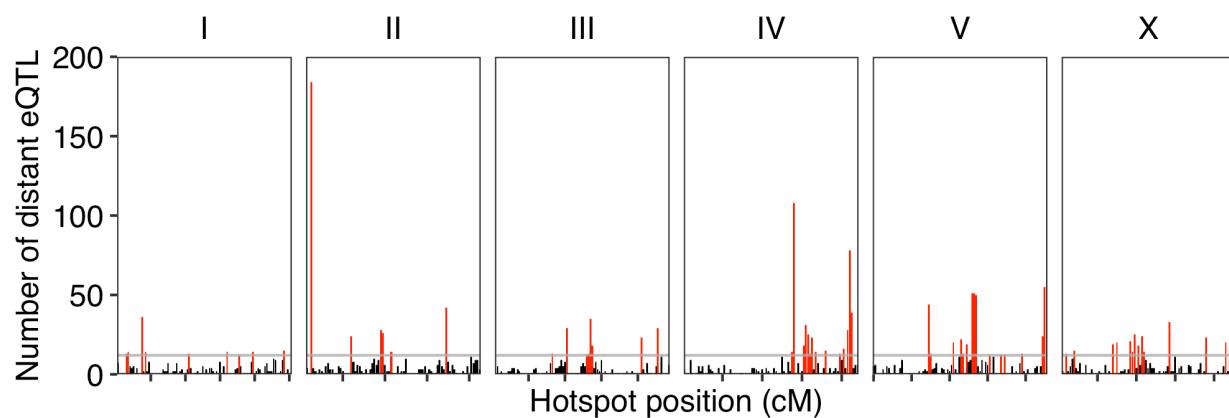
224 **Diverse nature of distant eQTL hotspots**

225 Distant eQTL were not uniformly distributed across the genome. Of the 3,360
226 distant eQTL, 1,828 were clustered into 67 hotspots, each of which affected the
227 expression of 12 to 184 transcripts (Fig. 3). Signatures of selection (Tajima's D values) in

228 hotspots are mostly negative, likely because of the recent selective sweeps
229 (Supplementary Fig. 4)²⁴.

230 GSEA on genes with transcript-level distant eQTL in each hotspot revealed
231 potential shared transcriptional regulatory mechanisms across different genes of the
232 same class (Supplementary Fig. 5, Supplementary Data 3). For example, the hotspot at
233 21.5 cM on chromosome II significantly affected the expression of heat stress related
234 genes (Fisher Exact Test, Bonferroni FDR corrected $p = 7.03E-7$). Our results also
235 showed that a single hotspot could regulate expression of genes in different classes.
236 The hotspot at 2.5 cM on chromosome II significantly affected the expression of genes
237 in three classes, including metallopeptidases (Fisher Exact Test, Bonferroni FDR
238 corrected $p = 1.31E-5$), collagen proteins (Fisher Exact Test, Bonferroni FDR corrected
239 $p = 3.11E-9$), and histones (Fisher Exact Test, Bonferroni FDR corrected $p = 1.26E-6$)
240 (Supplementary Fig. 5, Supplementary Data 3). Furthermore, different hotspots could
241 affect the expression of the same gene class. For example, the hotspot at 45.5 cM on
242 chromosome III was also enriched with distant eQTL of histones (Fisher Exact Test,
243 Bonferroni FDR corrected $p = 8.2E-7$) like the hotspot at 2.5 on chromosome II
244 (Supplementary Fig. 5, Supplementary Data 3). Regulatory genes, such as transcription
245 factors and chromatin cofactors, that are located in each hotspot could underlie the
246 regulation of multiple genes. We found previously known or predicted genes encoding
247 chromatin cofactors and transcription factors³⁷⁻³⁹ in 24 and 59 of the 67 hotspots,
248 respectively (Supplementary Fig. 6).

249



250

251 **Fig. 3: Distant eQTL hotspots.**

252 The number of distant eQTL (y-axis) in each 0.5 cM bin across the genome (x-axis) is
253 shown. Tick marks on the x-axis denote every 10 cM. The horizontal gray line indicates
254 the threshold of 12 eQTL. Bins with 12 or more eQTL were identified as hotspots and
255 are colored red. Bins with fewer than 12 eQTL are colored black.
256

257 To identify causal genes and variants underlying hotspots, we performed fine
258 mapping on distant eQTL in each hotspot and filtered for the most likely candidate
259 variants (see Methods for details) (Supplementary Data 4). Then, we focused on the
260 filtered candidate variants that were mapped for at least four traits in each hotspot and
261 are in genes encoding transcription factors or chromatin cofactors. In total, we identified
262 36 candidate genes encoding transcription factors or chromatin cofactors for 34
263 hotspots. For example, the gene *txx-1*, which encodes a transcription factor necessary
264 for thermosensation in the AFD neurons^{40,41}, might underlie the expression variation of
265 97 transcripts with distant eQTL in three hotspots between 44.5 cM and 45.5 cM on
266 chromosome V. TTX-1 regulates expression of *gcy-8* and *gcy-18* in AFD neurons^{40,41}, but
267 no eQTL were detected for the two genes likely because we measured the expression
268 of whole animals. Additionally, the linker histone gene *hil-2*³⁹ might underlie the
269 expression variation of 46, 10, 17, and four transcripts with distant eQTL in the hotspots
270 at 28 cM, 30.5 cM, 31 cM and 31.5 cM, respectively, on chromosome IV. We also
271 performed GSEA for groups of transcripts whose expression traits were fine mapped to
272 the 36 candidate genes encoding transcription factors or chromatin cofactors. For
273 instance, the 17 traits that fine mapped to *hil-2* in the hotspot at 31 cM on chromosome
274 IV (Supplementary Fig. 7) were enriched in E3 ligases containing an F-box domain (Fisher
275 Exact Test, Bonferroni FDR corrected $p = 0.0003$) and transcription factors of the
276 homeodomain class (Fisher Exact Test, Bonferroni FDR corrected $p = 0.002$). Besides
277 the 36 candidate genes, the hundreds of other fine mapping candidates are not as
278 transcription factors or chromatin cofactors, suggesting other mechanisms underlying
279 distant eQTL. Altogether, as previously implicated in other species^{7,11,42}, our results
280 indicate that a diverse collection of molecular mechanisms likely cause gene expression
281 variation in *C. elegans*.

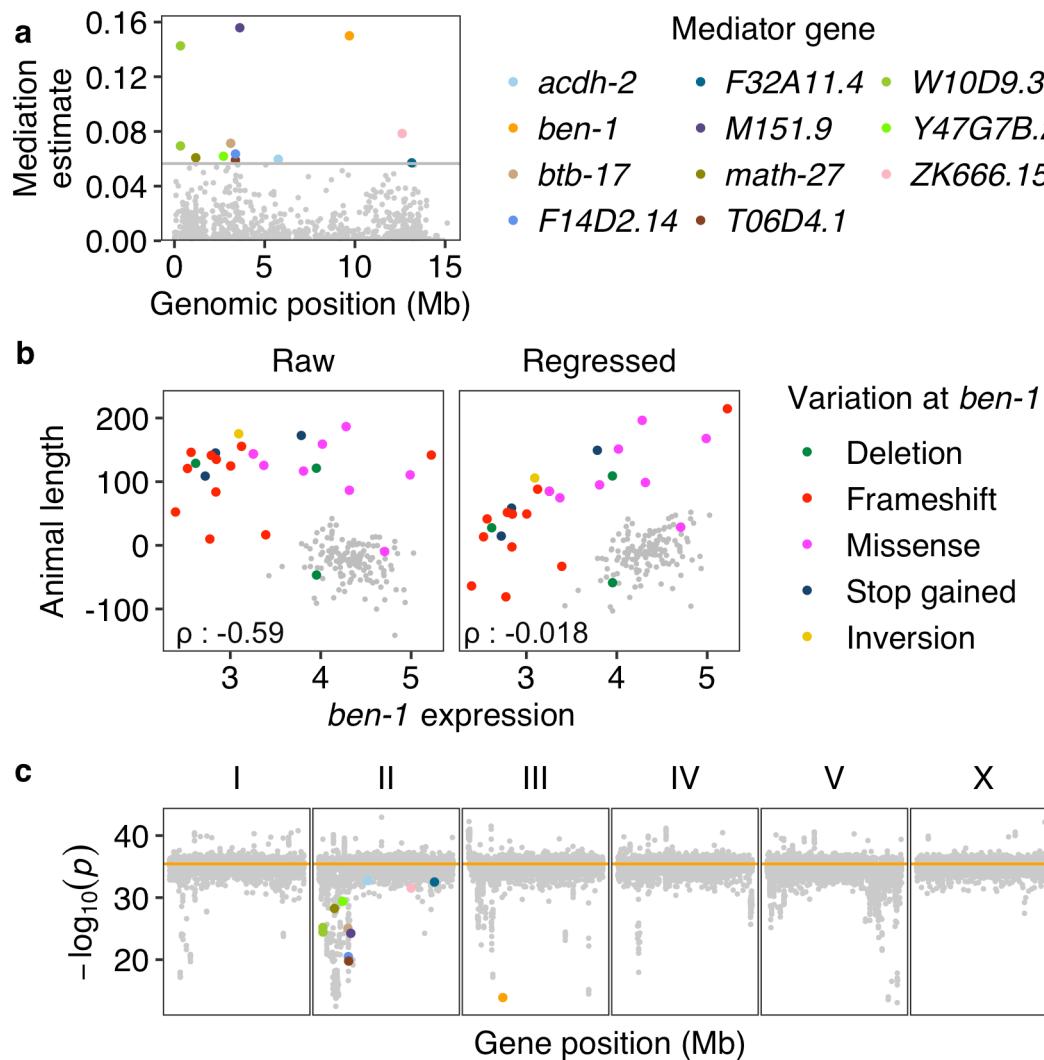
282 **Mediation analysis facilitates candidate gene prioritization**

283 Mediation analysis seeks to identify the mechanism that underlies the relationship
284 between an independent variable and a dependent variable via the inclusion of an
285 intermediary mediating variable. Because gene expression has been found to play an
286 intermediate role between genotypes and phenotypes, it could help to identify the causal
287 mediating genes between genotypes and phenotypes in quantitative genetics mapping
288 studies. We have previously identified mediation effects of *scb-1* expression on
289 responses to several chemotherapeutics and *sqst-5* expression on differential responses
290 to exogenous zinc using linkage mapping experiments^{9,43}. To validate if our expression
291 and eQTL data can be used to identify candidate genes, we first performed mediation
292 analysis on one published GWA study of variation in responses to the commonly used
293 anthelmintic albendazole (ABZ)⁴⁴.

294 Previously, wild *C. elegans* strains were exposed to ABZ and measured for effects
295 on development to identify genomic regions that contribute to variation in ABZ
296 resistance. A single-marker GWA mapping was performed first to detect two QTL on
297 chromosomes II and V, but no putative candidate gene was identified. Using a burden
298 mapping approach, prior knowledge of ABZ resistance in parasitic nematodes, and
299 manually curation of raw sequence read alignment files, the gene *ben-1* was found to
300 underlie natural variation in ABZ resistance variation⁴⁴. The single-marker GWA mapping
301 was not able to detect an association between ABZ resistance and *ben-1* variation
302 because of high allelic heterogeneity caused by rare SNVs and structural variants
303 (Supplementary Fig. 8). However, rare SNVs or structural variants might lead to changes
304 in *ben-1* expression and ABZ resistance. We found two distant eQTL, in regions
305 overlapping with the two organism-level ABZ QTL, for *ben-1* expression variation.
306 Therefore, these results provided an excellent opportunity to test the effectiveness of
307 mediation analysis among organism-level phenotypes, genotype, and gene expression.
308 The mediation estimate for *ben-1* expression was the second strongest hit in the analysis
309 on the phenotype (animal length in response to ABZ), the genotype (GWA QTL of the
310 phenotype), and the expression of 1,157 transcripts (Fig. 4a). We found a moderate
311 negative correlation between the expression of *ben-1* and animal length and almost no

312 correlation after we regressed animal length by the expression of *ben-1* (Fig. 4b),
313 suggesting that expression variation impacts differences in ABZ responses. We further
314 examined genetic variants across strains and found that those strains with relatively low
315 *ben-1* expression and high ABZ resistance all harbor SNVs or structural variants with
316 different predicted effects (Fig. 4b), suggesting that the extreme allelic heterogeneity at
317 the *ben-1* locus might affect ABZ response variation by reducing the abundance of this
318 beta-tubulin. To test the impact of expression variation on phenotypic variation, we
319 regressed animal length by expression of every transcript in our data and performed
320 GWA mappings. Then, we compared the GWA mapping significance value after
321 regression to the original GWA mapping significance value at a pseudo variant marker
322 that represents all the variants in *ben-1* (Fig. 4c, Supplementary Fig. 8)⁴⁵. We found
323 animal length regressed by the expression of *ben-1* showed one of the largest drops in
324 significance, and significance in most of the other mappings was approximately equal to
325 the original significance value (Fig. 4c, Supplementary Fig. 8). These results indicated
326 that increasing *ben-1* expression decreases resistance to ABZ and suggested the
327 applicability of mediation analysis using the expression and eQTL data for other *C.*
328 *elegans* quantitative traits.

329



330

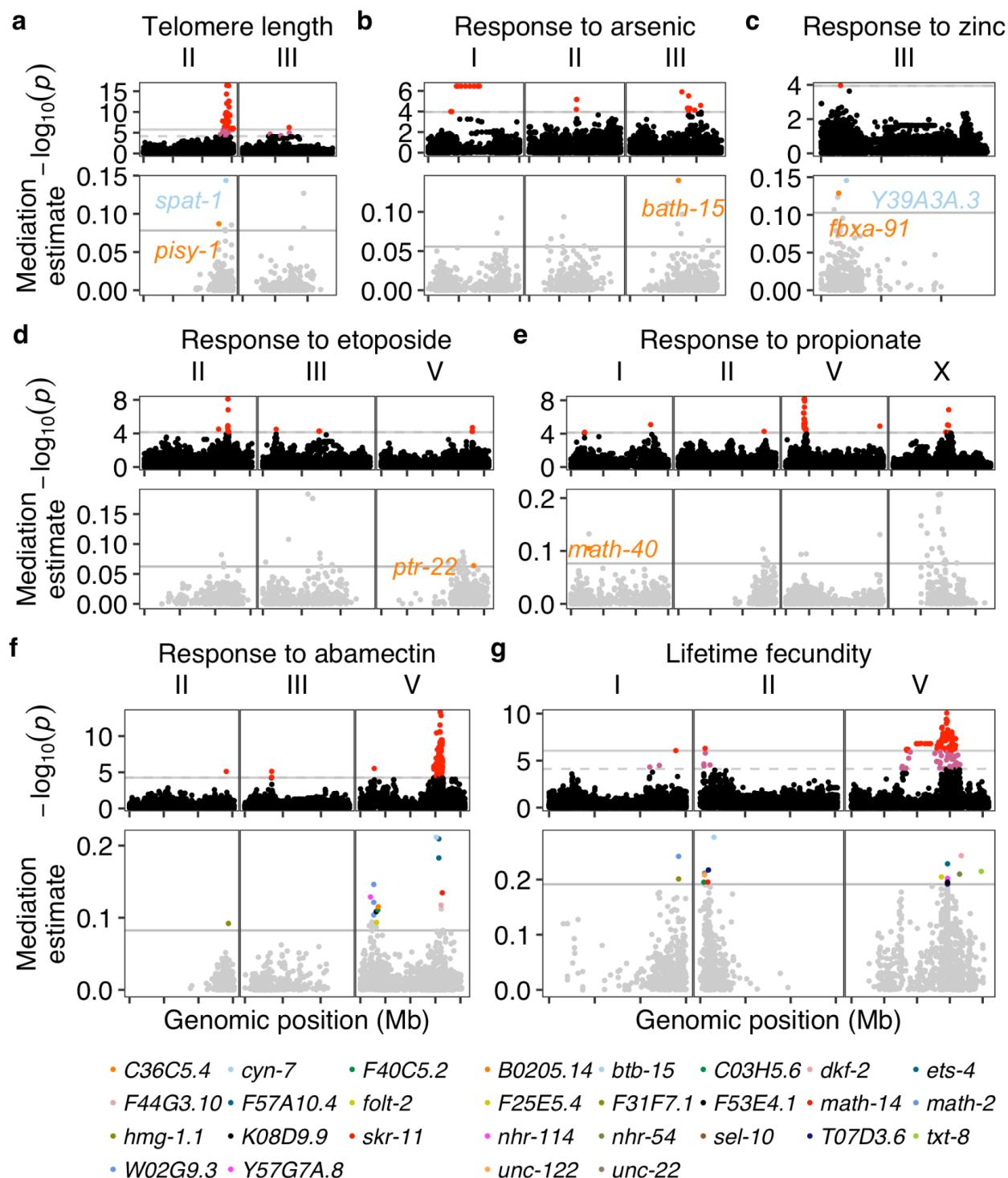
331 **Fig. 4: Mediation effects of *ben-1* expression on *C. elegans* resistance to ABZ.**

332 **a** Mediation estimates (y-axis) calculated as the indirect effect that differences in
 333 expression of each gene play in the overall phenotype are plotted against the genomic
 334 position of the eQTL (x-axis) on chromosome II. The horizontal gray line represents the
 335 99th percentile of the distribution of mediation estimates. Significant mediators are
 336 colored other than gray by their genes as shown in the legend. **b** The correlation of *ben-*
 337 *1* expression (x-axis) to raw animal length and to animal length regressed by *ben-1*
 338 expression on y-axis. The Pearson's coefficient ρ for each correlation was indicated at
 339 bottom left. Strains are colored by the type of their genetic variants in *ben-1*. Strains
 340 without identified variants are colored gray. **c** Significance at the pseudo variant marker
 341 of 25,837 GWA mappings. Each point represents a GWA mapping that is plotted with its
 342 $-\log_{10}(p)$ value (y-axis) at the pseudo variant marker (III: 3,539,640) against the genomic
 343 locations (x-axis) of the transcript of which the expression was used in regression for
 344 animal length. Points for traits regressed by expression of transcripts identified as
 345 significant mediators are colored as in (a). The orange horizontal line represents the
 346 significance at the pseudo variant marker using the raw animal length of 167 strains

347 (Supplementary Fig. 8). GWA mapping results of 12 traits regressed by expression of
348 mitochondrial genes were excluded but all with significance close to the horizontal line.
349

350 We further applied mediation analysis to another eight previously published
351 studies of *C. elegans* natural variation and GWA mappings in diverse traits, including
352 telomere length⁴⁶ (Fig. 5a), responses to arsenic⁴⁷ (Fig. 5b), zinc⁴³ (Fig. 5c), etoposide⁴⁸
353 (Fig. 5d), propionate⁴⁹ (Fig. 5e), abamectin⁵⁰ (Fig. 5f), dauer formation in response to
354 pheromone⁵¹, and lifetime fecundity⁵² (Fig. 5g). Causal variants and genes that partially
355 explained the phenotypic variation in all the eight traits, except for lifetime fecundity,
356 have been identified using fine mappings and genome-editing experiments^{43,46-52}. Only
357 one causal gene, *dbt-1* (for arsenic response variation⁴⁷), has eQTL detected and its
358 expression was tested in mediation analysis for arsenic response variation⁴⁷ (Fig. 5b). No
359 significant mediation effects were found on arsenic response variation by the expression
360 of *dbt-1*. We also did not observe significant differential expression between strains with
361 different alleles at the previously validated causal *dbt-1* QTV (II:7944817)⁴⁷. Therefore,
362 this causal variant possibly causes arsenic response variation only by affecting
363 enzymatic activity⁴⁷ and not the abundance of the *dbt-1* transcript. Instead, we identified
364 *bath-15* as a significant mediator gene for arsenic response variation (Fig. 5b). For the
365 other seven organism-level traits, putative genes whose expression likely mediated the
366 phenotypic variation were detected for six of the traits (Fig. 5). For example, the top
367 mediator gene for the variation in responses to abamectin was *cyn-7*, which is predicted
368 to have peptidyl-prolyl cis-trans isomerase acitivity (Fig. 5f)⁵³. For the variation in lifetime
369 fecundity (Fig. 5g), one of the putative mediator genes was *ets-4*, which is known to
370 affect larval developmental rate, egg-laying rate, and lifespan⁵⁴. Mediator genes suggest
371 candidate genes in addition to those genes identified in fine mappings or linkage
372 mappings. Taken together, we concluded that mediation analysis using the newly
373 generated expression and eQTL data facilitates candidate gene prioritization in GWA
374 studies.

375



376

377 **Fig. 5: Mediation effects of gene expression on variation in seven organism-level**
378 **phenotypes of *C. elegans*.**

379 GWA mapping and mediation analysis results of natural variation in *C. elegans* telomere
380 length (a), responses to arsenic (b), zinc (c), etoposide (d), propionate (e), abamectin (f),

381 and lifetime fecundity (**g**). Top panel: A Manhattan plot indicating the GWA mapping
382 result for each phenotype is shown. Each point represents an SNV that is plotted with
383 its genomic position (x-axis) against its $-\log_{10}(p)$ value (y-axis) in mapping. SNVs that
384 pass the genome-wide EIGEN threshold (the dotted gray horizontal line) and the
385 genome-wide Bonferroni threshold (the solid gray horizontal line) are colored pink and
386 red, respectively. QTL were identified by the EIGEN (**c,d,e,f**) or Bonferroni (**a,b,g**)
387 threshold. Only chromosomes with identified QTL were shown. Bottom panel: Mediation
388 estimates (y-axis) calculated as the indirect effect that differences in expression of each
389 gene plays in the overall phenotype are plotted against the genomic position (x-axis) of
390 the eQTL. The horizontal gray line represents the 99th percentile of the distribution of
391 mediation estimates. The mediator genes with adjusted $p < 0.05$ and interpretable
392 mediation estimate $>$ the 99th percentile estimates threshold are colored other than gray
393 and labeled in the panel (**a-e**) or below the panel (**f, g**). Tick marks on x-axes denote
394 every 5 Mb.
395

396 Discussion

397 *C. elegans* was the first metazoan to have its genome sequenced and has been
398 subjected to numerous genetic screens to identify the genes that underlie diverse traits,
399 including programmed cell death, drug responses, development, and behaviors. Despite
400 huge efforts by a large research community, over 60% of its genes have not been
401 curated with functional annotations or associated with defined mutant phenotypes⁵⁵. A
402 likely reason is that most *C. elegans* research uses the reference strain N2 under
403 laboratory conditions, and the functions of many genes might only be revealed in natural
404 environments or in different genetic backgrounds⁵⁶. In the last decade, wild *C. elegans*
405 strains have exhibited diverse phenotypic variation in natural ecology studies^{16,25,29,57-59}.
406 Here, we provide an unprecedentedly large resource of transcriptome profiles from wild
407 *C. elegans* strains. We used these data and GWA mappings to study gene regulation
408 variation and detected 6,545 eQTL associated with variation in expression of 5,291
409 transcripts of 4,520 genes. These genes are enriched in processes, including the
410 proteasome, metabolism, stress response, etc., suggesting gene expression regulation
411 plays an important role in adaptation of natural *C. elegans* strains to various
412 environments^{60,61}. We identified local eQTL that explained most of the narrow-sense
413 heritability (h^2) and significantly larger variance than distant eQTL, likely because of

414 higher possibilities of pleiotropy and thus stronger selection pressures. We also
415 observed lower variation in gene expression than in genome sequence and
416 underrepresentation of essential genes among all of the genes identified with eQTL,
417 suggesting stabilizing selection against gene expression as previously observed in *C.*
418 *elegans* and other species^{5,12,62,63}.

419 Although previous *C. elegans* eQTL studies using recombinant inbred lines have
420 revealed rich information on the genetic basis of gene expression variation, mapping
421 using 207 genetically distinct wild strains has the advantage of much greater genetic
422 diversity. We reanalyzed results of one previous study that used linkage mapping to
423 identify eQTL from the young adult stage of N2xCB4856 recombinant inbred lines^{3,9}. We
424 reclassified 1,208 local eQTL and 1,179 distant eQTL for 2,054 microarray probes of
425 2,003 genes (Supplementary Fig. 9a). Both the eQTL GWA and linkage mappings
426 detected overlapping local eQTL for 454 genes and distant eQTL for 19 genes, indicating
427 that the CB4856 strain carries the common alternative alleles among wild *C. elegans*
428 strains for these 473 loci. However, among the 6,545 eQTL that we detected, the strains
429 N2 and CB4856 shared the same genotypes in 4,476 eQTL, which could not be
430 discovered using N2xCB4856 recombinant inbred lines. Alternatively, RIAILs might have
431 less linkage disequilibrium between nearby variants and thus smaller eQTL regions of
432 interest than eQTL in wild *C. elegans* strains. The GWA eQTL in this study have a median
433 region of interest of 2.1 Mb (ranged from 12 kb to 18 Mb), whereas the N2xCB4856
434 RIAILs eQTL showed a median size of 0.55 Mb (ranged from 149 bp to 6.8 Mb), which
435 might make the identification of underlying causal variants easier. We further found 17
436 distant eQTL hotspots overlapped between the two studies (Supplementary Fig. 9b).
437 However, these shared hotspots comprise different genes between the two studies,
438 indicating that variation in regulatory factors is not common between the linkage and
439 association mapping studies. Future research should leverage both types of mapping
440 studies to identify common regulatory mechanisms, focusing on local eQTL.

441 In addition to the high linkage disequilibrium across the *C. elegans* genome, the
442 recently discovered hyper-divergent genomic regions made this eQTL study challenging.
443 Approximately 20% of the genomes in some wild *C. elegans* strains were found to have

444 extremely high diversity compared to the N2 reference genome²⁷. Short-sequence reads
445 of wild *C. elegans* strains often fail to align to the N2 reference genome in these regions
446 and showed lower coverage than in other regions²⁷. Similarly, expression levels of genes
447 in hyper-divergent regions could be underestimated because of the poor alignment of
448 RNA-seq reads. Therefore, we only used expression of transcripts in non-divergent
449 regions to map eQTL and flagged the loci that are in common hyper-divergent regions,
450 where we are less confident in the genotypes of wild strains (Supplementary Data 2).
451 Furthermore, we only used distant eQTL that are not in common hyper-divergent regions
452 to identify hotspots. Because hyper-divergent regions were suggested to be under long-
453 term balancing selection, our estimates of Tajima's *D* in hotspots are probably biased
454 towards lower values. Future efforts using long-read sequencing are necessary to study
455 the sequence, expression, natural selection, and evolution of genes in hyper-divergent
456 regions.

457 Variation in gene expression was suggested to impact organism-level phenotypic
458 variation^{7,64–66}. Combining previous GWA studies in *C. elegans* with expression of genes
459 with eQTL, we used mediation analysis to search for organism-level phenotypic variation
460 that can be explained by variation in gene expression. Compared to previous studies
461 using mediation analysis on gene expression and eQTL data from the N2xCB4856
462 recombinant inbred lines^{9,43}, we added a multiple testing correction procedure to our
463 mediation analysis. We performed mediation analysis on ABZ response variation⁴⁴. The
464 causal gene *ben-1* underlying the trait was identified using a burden mapping approach⁴⁴
465 along with prior knowledge^{67,68} about the role of beta-tubulin in this drug response.
466 Although two GWA QTL on chromosomes II and V were found, they were identified likely
467 because of their interchromosomal linkage disequilibrium to variants in the *ben-1* locus⁴⁴
468 (Supplementary Fig. 8). The single-marker GWA mapping could not associate ABZ
469 response variation because of the extreme allelic heterogeneity at the *ben-1* locus.
470 However, we used mediation analysis to identify *ben-1* without consideration of prior
471 knowledge or burden mapping results, demonstrating the power of the approach (Fig.
472 4a). We further identified significant mediators for seven other organism-level traits (Fig.

473 5). The expression of these mediator genes could affect the corresponding phenotypic
474 variation, which should be validated in the future.

475 Mediation analysis provides an efficient hypothesis-generating approach to be
476 performed in parallel to fine mappings. Additionally, mediator genes could contribute to
477 organism-level phenotypic variation in addition to causal genes identified using fine
478 mappings. One limitation of fine mappings is that searching for causal genes and variants
479 is restricted to the QTL region of interest. Mediation analysis can make statistical
480 connections between the organism-level phenotypes and expression of genes far away
481 from the QTL. As mentioned above, large GWA QTL regions of interest make it difficult
482 to identify causal genes, which require validation using genome editing. Future *C.
483 elegans* GWA studies should use both fine mappings and mediation analysis to prioritize
484 candidate genes. If the candidate genes overlap between the two approaches, then
485 validation approaches can be initiated using genome editing. In cases where the two
486 approaches identify different candidate genes, prioritization using prior knowledge
487 across all genes identified by both approaches can inform which genes should be tested
488 for validation using genome editing. Previous studies using fine mappings prioritized
489 candidate genes harboring coding variants predicted to have strong functional impacts.
490 In mediation analysis, noncoding variants that likely affect expression of mediator genes
491 could also be nominated as candidates. For example, upstream variants were suggested
492 to underlie expression variation of the gene *scb-1*, which mediated differences in
493 responses to bleomycin and three other chemotherapeutics^{9,69}. To summarize, we
494 recommend using both fine mappings and mediation analyses to nominate candidate
495 genes and variants.

496 The goal of quantitative genetics is to understand the genetic basis and
497 mechanisms underlying phenotypic variation. Here, we showed that mediation analysis,
498 which uses expression and eQTL data to search connections between genetic variants
499 and complex traits, provides additional loci that might further explain phenotypic
500 variation. The framework we developed for mediation analysis complements marker-
501 based GWA mappings and is also applicable using various other intermediate traits,
502 such as small RNAs, proteins, and metabolites. Any genes and variants underlying

503 variation in these factors can be nominated as candidates for phenotypic validation.
504 Furthermore, we could measure all of these data and complex traits from the exact same
505 samples using *C. elegans*, which can be easily grown at large scale to have synchronized
506 isogenic populations. Analyses using measurements of mRNAs, small RNAs, proteins,
507 and metabolites could strengthen conclusions about causal genes and mechanisms
508 underlying complex traits using a more holistic perspective of organismal phenotypic
509 variation. We foresee this strategy will greatly improve the powers of quantitative genetic
510 mappings in the future.

511 **Methods**

512 **C. elegans strains**

513 We obtained 207 wild *C. elegans* strains from *C. elegans* Natural Diversity Resource
514 (CeNDR)²⁵. Animals were cultured at 20°C on modified nematode growth medium
515 (NGMA) containing 1% agar and 0.7% agarose to prevent burrowing and fed *Escherichia*
516 *coli* (*E. coli*) strain OP50⁷⁰. Prior to each assay, strains were grown for three generations
517 without starvation or encountering dauer-inducing conditions⁷⁰.

518

519 **Animal growth and harvest**

520 We grew and harvested synchronized populations of each strain at the young adult stage
521 with independently grown and prepared biological replicates. Specifically, L4 larval stage
522 hermaphrodites were grown to the gravid adult stage on 6 cm plates and were bleached
523 to obtain synchronized embryos. Approximately 1,000 embryos were grown on each 10
524 cm plate to the young adult stage and were harvested after the first embryo was
525 observed. M9 solution was used to wash harvested animals twice to remove *E. coli*.
526 Animals were then pelleted by centrifugation (2000 rpm for one minute) and Trizol
527 reagent (Ambion) was added to maintain RNA integrity before storage at -80°C.

528

529 **RNA extraction**

530 Frozen samples in Trizol were thawed at room temperature and 100 µL acid-washed
531 sand (Sigma, catalog no. 274739) was added to help to disrupt animal tissues. Then
532 chloroform, isopropanol, and ethanol were used for phase separation, precipitation, and
533 washing steps, respectively. Total RNA pellets were resuspended in nuclease-free water.
534 The concentration of total RNA was determined using the Qubit RNA XR Assay Kit
535 (Invitrogen, catalog no. Q33224). RNA quality was measured using the 2100 Bioanalyzer
536 (Agilent). RNA samples with a minimum RNA integrity number (RIN) of 7 were used to
537 construct Illumina sequencing libraries.

538

539 **RNA library construction and sequencing**

540 Illumina RNA-seq libraries were prepared in 96-well plates. Replicates of the same strain
541 were prepared in different 96-well plates. For each sample, mRNA was purified and
542 enriched from 1 µg of total RNA using the NEBNext Poly(A) mRNA Magnetic Isolation
543 Module (New England Biolabs, catalog no. E7490L). RNA fragmentation, first and
544 second strand cDNA synthesis, and end-repair processing were performed with the
545 NEBNext Ultra II RNA Library Prep with Sample Purification Beads (New England
546 Biolabs, catalog no. E7775L). The cDNA libraries were adapter-ligated using adapters
547 and unique dual indexes in the NEBNext Multiplex Oligos for Illumina (New England
548 Biolabs, catalog no. E6440, E6442) and amplified using 12 PCR cycles. All procedures
549 were performed according to the manufacturer's protocols. The concentration of each
550 RNA-seq library was determined using Qubit dsDNA BR Assay Kit (Invitrogen, catalog
551 no. Q32853). Approximately 96 RNA-seq libraries were pooled and quantified with the
552 2100 Bioanalyzer (Agilent) at Novogene, CA, USA. Each of the pools of libraries were
553 sequenced on a single lane of an Illumina NovaSeq 6000 platform, yielding 150-bp
554 paired-end (PE150) reads.

555

556 In total, RNA-seq data of 608 samples from 207 wild *C. elegans* strains in seven pooled
557 libraries were obtained with an average of 32.6 million reads per sample and a minimum
558 of 16.6 million reads. Of the 207 strains, 194 strains with three replicates and 13 strains
559 with two replicates.

560

561 **Sequence processing and expression abundance quantification**

562 Adapter sequences and low-quality reads in raw sequencing data were removed using
563 *fastp* (v0.20.0)⁷¹. *FastQC* (v0.11.8) analysis
564 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc>) was performed on trimmed
565 FASTQ files to assess read quality (adapter content, read-length distribution, per read
566 GC content, etc.). For RNA-seq mapping, SNV-substituted reference transcriptomes for
567 each of the wild *C. elegans* strains were generated using *BCFtools* (v.1.9)⁷², *gffread*
568 (v0.11.6)⁷³, the N2 reference genome (WS276), a GTF file (WS276)⁵³, and the hard-filtered
569 isotype variant call format (VCF) 20200815 CeNDR release (Supplementary Fig. 1).

570 Transposable element (TE) consensus sequences of *C. elegans* were also extracted from
571 Dfam (release 3.3)⁷⁴ using scripts (<https://github.com/fansalon/TEconsensus>). We used
572 *Kallisto* (v0.44.0) to (1) pseudoalign trimmed RNA-seq reads from each sample to the
573 transcriptome index built from the strain-specific SNV-substituted reference
574 transcriptome (65,173 transcripts) and TE consensus sequences (157 TEs) and (2)
575 quantify expression abundance at the transcript level⁷⁵. On average, 31.3 million reads
576 pseudoaligned to the transcriptome index per sample with a minimum of 15.5 million
577 reads, which were sufficient to capture the expression of more than 70% of the *C.*
578 *elegans* reference genome genes. We used the 608 samples of 207 strains and 39,008
579 transcripts of protein-coding genes and pseudogenes in our analysis.
580

581 **Selection of reliably expressed transcripts**

582 We first normalized the raw counts of transcript expression abundances without the
583 default filtering of low abundance transcripts using the R package *sleuth* (v0.30.0)⁷⁶.
584 Then, we filtered reliably expressed transcripts (26,043) by requiring at least five
585 normalized counts in all the replicates of at least ten strains (Supplementary Fig. 1). We
586 also filtered out 3,775 transcripts of 2,597 genes that are in hyper-divergent genomic
587 regions of at least one strain. We further excluded 194 transcripts in hyper-divergent
588 regions of more than 107 of the 207 strains. In summary, we collected reliable expression
589 abundance for 25,849 transcripts of 16,094 genes (15,364 protein-coding genes and 730
590 pseudogenes).
591

592 **Selection of well clustered samples**

593 We used sample-to-sample distance to select well clustered samples (Supplementary
594 Fig. 1). We first summarized raw counts of reliably expressed transcripts into gene-level
595 abundances using the R package *tximport* (v1.10.1)⁷⁷. Then, we performed variance
596 stabilizing transformations on the gene expression profile using the *vst()* function in the
597 R package *DESeq2* (v1.26.0), which generated log₂ scale normalized expression data⁷⁸.
598 Sample-to-sample pairwise Euclidean distances among the 608 samples were
599 calculated using the generic function *dist()* in R⁷⁹. Our basic assumption is that intra-

600 strain distances among replicates should be smaller than inter-strain distances. Because
601 the majority of the 207 wild *C. elegans* strains exhibit low overall genetic diversity (Fig.
602 1d)^{24,35,80}, we required that the intra-strain distances of replicates be smaller than the
603 median of inter-strain distances of the strain to other strains. Specifically, for each strain,
604 if all of its intra-strain distances were smaller than the median of its inter-strain distances,
605 then all of its replicates were kept. If none of its intra-strain distances were smaller than
606 the median of its inter-strain distances, then all samples of the strain were removed. For
607 strains with three replicates, if one or two of its three intra-strain distances were smaller
608 than the median of its inter-strain distances, then the two replicates with the minimum
609 distances were kept. After removal of some outlier samples, the median of inter-strain
610 distances would change. Therefore, we repeatedly performed the procedures of data
611 transformation, sample-to-sample distance calculation, and filtering by comparing inter-
612 and intra-strain distances until no more samples were removed. Eventually, 561 samples
613 of 207 strains were selected as well clustered samples, which comprised 147 strains
614 with three replicates and 60 strains with two replicates.

615

616 **Transcript expression abundance normalization**

617 We used the function *norm_factors()* in the R package *sleuth* (v0.30.0)⁷⁶ to compute the
618 normalization factors for each sample using the raw transcripts per million reads (TPM)
619 of 22,268 reliably expressed transcripts in non-divergent regions of the 207 strains and
620 their well clustered samples. Then, we normalized the raw TPM of all the 25,849 reliably
621 expressed transcripts of each sample with the normalization factors and used
622 $\log_2(\text{normalized TPM} + 0.5)$ for downstream analysis unless indicated otherwise.

623

624 **Sample age estimation**

625 To further verify the homogeneous developmental stage of our samples, we evaluated
626 the age of each sample when they were harvested using the R package *RAPToR*
627 (v1.1.3)²⁸ (Supplementary Fig. 1). As the requirement of the package, we first generated
628 gene-level expression abundances. Raw TPM of 22,268 reliably expressed transcripts in
629 non-divergent regions were summarized into abundances of 13,637 genes using the R

630 package *tximport* (v1.10.1)⁷⁷. Normalization factors for each sample using gene-level
631 abundances were calculated as described for transcript level and were used to normalize
632 gene level TPM. Correlation of $\log_2(\text{normalized TPM} + 0.5)$ of our data against the
633 reference gene expression time series (Cel_YA_2) in *RAPToR* was computed using the
634 function *ae()* in *RAPToR* with 10,489 intersected genes and default parameters.

635

636 **Genetic and expression relatedness**

637 Genetic variation data for 207 *C. elegans* isotypes were acquired from the hard-filtered
638 isotype variant call format (VCF) 20200815 CeNDR release. These variants were pruned
639 to the 851,105 biallelic single nucleotide variants (SNVs) without missing genotypes. We
640 converted this pruned VCF file to a PHYLIP file using the *vcf2phylip.py* script⁸¹.
641 Expression distance among the 207 wild strains was calculated based on the mean
642 expression of 22,268 transcripts without missing data using the function *dist()* in R. The
643 unrooted neighbor-joining trees for genetic and expression relatedness were made using
644 the R packages *phangorn* (v2.5.5)⁸², *ape* (v5.6)⁸³ and *ggtree* (v1.14.6)⁸⁴.

645

646 **eQTL mapping**

647 *Input phenotype and genotype data*

648 For the 25,849 transcripts, we summarized the expression abundance of replicates to
649 have the mean expression for each transcript of each strain as phenotypes used in GWA
650 mapping (Supplementary Data 1). Genotype data for each of the 207 strains were
651 acquired from the hard-filtered isotype VCF (20200815 CeNDR release).

652

653 *Permutation-based FDR threshold*

654 We performed GWA mapping using the pipeline *cegwas2-nf*
655 (<https://github.com/AndersenLab/cegwas2-nf>). The pipeline uses the eigen-
656 decomposition significance (EIGEN) threshold or the more stringent Bonferroni-
657 corrected significance (BF) threshold to correct for multiple testing because of the large
658 number of genetic markers (SNVs). To further correct for false positive QTL because of
659 the large number of transcript expression traits, we computed a permutation-based

660 False Discovery Rate (FDR) at 5%. We randomly selected 200 traits from our input
661 phenotype file and permuted each of them 200 times. These 40,000 permuted
662 phenotypes were used as input to call QTL using *cegwas2-nf* with EIGEN and BF
663 threshold, respectively, as previously described^{47,49,52}. Briefly, we used *BCFtools*⁷² to filter
664 variants that had any missing genotype calls and variants that were below the 5% minor
665 allele frequency. Then, we used *-indep-pairwise 50 10 0.8* in *PLINK* v1.9^{85,86} to prune the
666 genotypes to 20,991 markers with a linkage disequilibrium (LD) threshold of $r^2 < 0.8$ and
667 then generated the kinship matrix using the *A.mat()* function in the R package *rrBLUP*
668 (v4.6.1)⁸⁷. The number of independent tests (N_{test}) within the genotype matrix was
669 estimated using the R package *RSpectra* (v0.16.0) (<https://github.com/yixuan/RSpectra>)
670 and *correlateR* (0.1) (<https://github.com/AEBilgrau/correlateR>). The eigen-
671 decomposition significance (EIGEN) threshold was calculated as $-\log_{10}(0.05/N_{test})$. We
672 used the *GWAS()* function in the *rrBLUP* package to perform the genome-wide mapping
673 with the EMMA algorithm³². QTL were defined by at least one marker that was above the
674 EIGEN or BF threshold. The EIGEN and BF %5 FDR was calculated as the 95 percentile
675 of the significance of all the detected QTL under each threshold. The EIGEN and BF 5%
676 FDR thresholds were 6.11 and 7.76, respectively.

677

678 eQTL mapping

679 We performed GWA mapping on the expression traits of the 25,849 transcripts as for
680 permuted expression traits but using the EIGEN 5% FDR (6.11) as the threshold. We
681 identified QTL with significance that also passed the Bonferroni 5% FDR threshold to
682 locate the best estimate of QTL positions with the highest significance. We used the
683 generic function *cor()* in R and Pearson correlation coefficient to calculate the phenotypic
684 variance explained by each QTL. We used the *LD()* function from the R package *genetics*
685 (v1.3.8.1.2) (<https://cran.r-project.org/package=genetics>) to calculate the LD correlation
686 coefficient r^2 among QTL for traits with multiple eQTL.

687

688 eQTL classification

689 Local eQTL were classified if the QTL was within a 2 Mb region surrounding the
690 transcript. All other QTL were classified as distant.

691

692 **Heritability calculation**

693 Heritability estimates were calculated for each of the 25,849 traits used for eQTL
694 mapping as previously described⁸⁸. Narrow-sense heritability (h^2) was calculated with the
695 phenotype file and pruned genotypes in eQTL mapping using the functions *mmer()* and
696 *pin()* in the R package *sommer* (v4.1.2)⁸⁹. Broad-sense heritability (H^2) was calculated
697 using expression of replicates of each strain and the *lmer* function in the R package *lme4*
698 (v1.1.21) with the model phenotype $\sim 1 + (1|\text{strain})^{90}$.

699

700 **Hotspot identification**

701 We first filtered out distant eQTL in common hyper-divergent genomic regions of wild *C.*
702 elegans strains. Common hyper-divergent regions were defined among our 206 strains
703 (reference N2 excluded) as described previously²⁷. Briefly, we divided the genome into 1
704 kb bins and calculated the percentage of 206 strains that are classified as hyper-
705 divergent in each bin. Common hyper-divergent regions were defined as bins $\geq 5\%$ ²⁷.

706

707 Distant eQTL hotspots were identified by dividing the genome into 0.5 cM bins and
708 counting the number of non-divergent distant eQTL that mapped to each bin.
709 Significance was determined as bins with more eQTL than the 99th percentile of a
710 Poisson distribution using the maximum likelihood method and the function *eqpois()* in
711 the R package EnvStats^{1,3,9,91}.

712

713 **Reanalysis of RIAILs eQTL**

714 We reclassified eQTL detected in a previous study using microarray expression data
715 from synchronized young adult populations of 208 recombinant inbred advanced
716 intercross lines (RIAILs) derived from N2 and CB4856^{9,36}. A total of 2,540 eQTL from
717 2,196 probes were identified using linkage mappings⁹. We selected 2,387 eQTL of 2,054
718 probes that are in 2,003 live genes based on the probe-gene list in the R package

719 *linkagemapping* (<https://github.com/AndersenLab/linkagemapping>) and the GTF file
720 (WS276)⁵³. We classified 1,208 local eQTL and 1,179 distant eQTL as described above.
721 We further identified hotspots as above for 1,124 distant eQTL that are not in the hyper-
722 divergent regions of CB4856.

723

724 **Population genetics**

725 We use 851,105 biallelic SNVs with no missing calls among the 207 strains from the
726 hard-filtered VCF 20200815 CeNDR release to calculate population genomic statistics.
727 Tajima's D, Watterson's θ, and pi were all calculated using *scikit-allel*⁹². Each of these
728 statistics was calculated for non-overlapping 1,000-bp windows across the genome.

729

730 **Fine mapping for causal genes underlying hotspots**

731 For transcript expression traits with distant eQTL in hotspots, we performed fine
732 mapping using the pipeline *cegwas2-nf* as previously described⁴⁷. Briefly, we defined
733 QTL regions of interest from the GWA mapping as +/- 100 SNVs from the rightmost and
734 leftmost markers above the EIGEN 5% FDR significance threshold. Then, using
735 genotype data from the imputed hard-filtered isotype VCF (20200815 CeNDR release),
736 we generated a QTL region of interest genotype matrix that was filtered as described
737 above, with the one exception that we did not perform LD pruning. We used *PLINK*
738 v1.9^{85,86} to extract the LD between the markers used for fine mapping and the QTL peak
739 marker identified from GWA mappings. We used the same command as above to
740 perform fine mappings. To identify causal genes and variants that affect expression of
741 several transcripts underlying hotspots, we retained the fine-mapped candidate variants
742 that passed the following per trait filters: top 5% most significant markers; out of
743 common hyper-divergent genomic regions; with negative BLOSUM⁹³ scores as
744 characterized and annotated in CeNDR²⁵.

745

746 **Enrichment analysis**

747 Gene set enrichment analyses were carried out for all genes found with transcript-level
748 eQTL and for genes with distant eQTL in each hotspot using the web-based tool
749 *WormCat*³⁹.

750

751 **Mediation analysis**

752 *GWA mapping of diverse C. elegans phenotypes*

753 We obtained nine different phenotype data used in previous *C.elegans* natural variation
754 and GWA studies^{43,44,46–52}. We filtered genetically distinct isotype strains for each trait
755 based on CeNDR (20200815 release) and performed GWA mapping as for permuted
756 expression traits but mostly using EIGEN or BF as the threshold according to the original
757 studies. GWA was performed under EIGEN for two studies originally using BF as the
758 threshold^{48,49}.

759

760 *Mediation analysis*

761 For each QTL of the above phenotypes, we used the genotype (*Exposure*) at the
762 phenotype QTL peak, transcript expression traits (*Mediator*) that have eQTL overlapped
763 with the phenotype QTL, and the phenotype (*Outcome*) as input to perform mediation
764 analysis using the *medTest()* function and 1,000 permutations for *p*-value correction in
765 the R package MultiMed (v2.6.0)
766 (<https://bioconductor.org/packages/release/bioc/html/MultiMed.html>). For mediation,
767 we used only strains with all of the three input data types available and where variation
768 was found. For instance, between the 202 strains used in the study of ABZ resistance⁴⁴
769 and the 207 strains used in this study, 167 strains overlapped. Although we searched
770 overlapped eQTL against QTL in the GWA mapping for ABZ resistance using 202 strains
771 (Supplementary Fig. 8), 167 strains at maximum were used in mediation analysis.
772 Furthermore, because some transcripts were found in hyper-divergent regions in certain
773 strains and their expression data were filtered out, the rest of the strains with all of the
774 data types available might contain no variation in one or all of the three data types and
775 were not used in mediation analysis. For example, we found 1,193 eQTL overlapped with

776 the QTL on chromosome II in the GWA mapping for ABZ resistance, but only 1,157
777 mediation analyses were performed.

778
779 For mediators with adjusted $p < 0.05$ or mediation estimate greater than the 99th
780 percentile of the distribution of mediation estimates, we performed a second mediation
781 analysis as described previously⁹ using the *mediate()* function from the R package
782 *mediation* (version 4.5.0)⁹⁴ to filter out the uninterpretable results where the proportion of
783 the total effect (the estimated effect of genotype on phenotype, ignoring expression) that
784 can be explained by the mediation effect (the estimated effect of expression on
785 phenotype) is negative or larger than 100%.

786
787 *GWA of traits regressed by transcript expression*
788 We regressed the trait animal length (q90.TOF)⁴⁴ by expression of every transcript using
789 the generic function *residuals()* in R, which fits a linear model with the formula (*phenotype*
790 ~ *expression*) to account for any differences in phenotype parameters present in
791 transcript expression. Then GWA was performed for each regressed trait as for permuted
792 expression traits using BF as the threshold.

793
794
795

796 **Acknowledgements**

797 We would like to thank Stefan Zdraljevic and Samuel J. Widmayer for helpful
798 suggestions. G.Z. is supported by the NSF-Simons Center for Quantitative Biology at
799 Northwestern University (awards Simons Foundation/SFARI 597491-RWC and the
800 National Science Foundation 1764421). S.R.H. was funded by a DFG fellowship (HA
801 8449/1-1) from the Deutsche Forschungsgemeinschaft (www.dfg.de). E.C.A. is
802 supported by a National Science Foundation CAREER Award (IOS-1751035) and a grant
803 from the National Institutes of Health R01 DK115690. The *C. elegans* Natural Diversity
804 Resource is supported by a National Science Foundation Living Collections Award to
805 R.E.T. and E.C.A. (1930382). We also like to thank WormBase for which these analyses
806 would not have been possible.

807 **Author contributions**

808 E.C.A. conceived of and designed the study. D.L., S.R.H., and G.Z. prepared *C. elegans*
809 cultures. G.Z. and N.M.R. performed RNA-seq experiments. G.Z. analyzed the data. G.Z.
810 and E.C.A. wrote the manuscript.

811 **Competing interests**

812 The authors declare no competing interests.

813 **Data and code availability**

814 The raw RNA-seq data generated in this study are available at the NCBI Sequence Read
815 Archive (Project PRJNA669810). The raw expression counts and TPM quantified in this
816 study are available at NCBI's Gene Expression Omnibus (Series GSE186719). The RNA-
817 seq mapping pipeline can be found at <https://github.com/AndersenLab/PEmRNA-seq-nf>. The mediation analysis pipeline can be found at
818 https://github.com/AndersenLab/mediation_GWAeQTL. The datasets and code for
819 generating all figures can be found at <https://github.com/AndersenLab/WI-Ce-eQTL>.

821

822 **References**

- 823 1. Brem, R. B., Yvert, G., Clinton, R. & Kruglyak, L. Genetic dissection of
824 transcriptional regulation in budding yeast. *Science* **296**, 752–755 (2002).
- 825 2. West, M. A. L. *et al.* Global eQTL mapping reveals the complex genetic
826 architecture of transcript-level variation in Arabidopsis. *Genetics* **175**, 1441–1450
827 (2007).
- 828 3. Rockman, M. V., Skrovanek, S. S. & Kruglyak, L. Selection at linked sites shapes
829 heritable phenotypic variation in *C. elegans*. *Science* **330**, 372–376 (2010).
- 830 4. Zan, Y., Shen, X., Forsberg, S. K. G. & Carlborg, Ö. Genetic Regulation of
831 Transcriptional Variation in Natural *Arabidopsis thaliana* Accessions. *G3* **6**, 2319–
832 2328 (2016).
- 833 5. Kita, R., Venkataram, S., Zhou, Y. & Fraser, H. B. High-resolution mapping of cis-
834 regulatory variation in budding yeast. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E10736–
835 E10744 (2017).
- 836 6. GTEx Consortium *et al.* Genetic effects on gene expression across human tissues.
837 *Nature* **550**, 204–213 (2017).
- 838 7. Albert, F. W., Bloom, J. S., Siegel, J., Day, L. & Kruglyak, L. Genetics of trans-
839 regulatory variation in gene expression. *Elife* **7**, 1–39 (2018).
- 840 8. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across
841 human tissues. *Science* **369**, 1318–1330 (2020).
- 842 9. Evans, K. S. & Andersen, E. C. The Gene scb-1 Underlies Variation in
843 *Caenorhabditis elegans* Chemotherapeutic Responses. *G3* **10**, 2353–2364 (2020).
- 844 10. Rockman, M. V. & Kruglyak, L. Genetics of global gene expression. *Nat. Rev.*
845 *Genet.* **7**, 862–872 (2006).
- 846 11. Albert, F. W. & Kruglyak, L. The role of regulatory variation in complex traits and
847 disease. *Nat. Rev. Genet.* **16**, 197–212 (2015).
- 848 12. Hill, M. S., Vande Zande, P. & Wittkopp, P. J. Molecular and evolutionary
849 processes generating variation in gene expression. *Nat. Rev. Genet.* **22**, 203–215
850 (2021).

- 851 13. Brem, R. B., Yvert, G., Clinton, R. & Kruglyak, L. Genetic dissection of
852 transcriptional regulation in budding yeast. *Science* **296**, 752–755 (2002).
- 853 14. Gaertner, B. E. & Phillips, P. C. *Caenorhabditis elegans* as a platform for molecular
854 quantitative genetics and the systems biology of natural variation. *Genet. Res.* **92**,
855 331–348 (2010).
- 856 15. Snoek, B. L. *et al.* WormQTL2: an interactive platform for systems genetics in
857 *Caenorhabditis elegans*. *Database* **2020**, (2020).
- 858 16. Evans, K. S., van Wijk, M. H., McGrath, P. T., Andersen, E. C. & Sterken, M. G.
859 From QTL to gene: *C. elegans* facilitates discoveries of the genetic mechanisms
860 underlying natural variation. *Trends Genet.* **0**, (2021).
- 861 17. Li, Y. *et al.* Mapping determinants of gene expression plasticity by genetical
862 genomics in *C. elegans*. *PLoS Genet.* **2**, e222 (2006).
- 863 18. Viñuela, A., Snoek, L. B., Riksen, J. A. G. & Kammenga, J. E. Genome-wide gene
864 expression regulation as a function of genotype and age in *C. elegans*. *Genome*
865 *Res.* **20**, 929–937 (2010).
- 866 19. Li, Y. *et al.* Global genetic robustness of the alternative splicing machinery in
867 *Caenorhabditis elegans*. *Genetics* **186**, 405–410 (2010).
- 868 20. Sterken, M. G. *et al.* Ras/MAPK Modifier Loci Revealed by eQTL in *Caenorhabditis*
869 *elegans*. *G3* **7**, 3185–3193 (2017).
- 870 21. Snoek, B. L. *et al.* Contribution of trans regulatory eQTL to cryptic genetic variation
871 in *C. elegans*. *BMC Genomics* **18**, 500 (2017).
- 872 22. Ben-David, E. *et al.* Whole-organism eQTL mapping at cellular resolution with
873 single-cell sequencing. *Elife* **10**, (2021).
- 874 23. Snoek, B. L. *et al.* The genetics of gene expression in a *Caenorhabditis elegans*
875 multiparental recombinant inbred line population. *G3* **11**, (2021).
- 876 24. Andersen, E. C. *et al.* Chromosome-scale selective sweeps shape *Caenorhabditis*
877 *elegans* genomic diversity. *Nat. Genet.* **44**, 285–290 (2012).
- 878 25. Cook, D. E., Zdraljevic, S., Roberts, J. P. & Andersen, E. C. CeNDR, the
879 *Caenorhabditis elegans* natural diversity resource. *Nucleic Acids Res.* **45**, D650–
880 D657 (2017).

- 881 26. Thompson, O. A. *et al.* Remarkably Divergent Regions Punctuate the Genome
882 Assembly of the *Caenorhabditis elegans* Hawaiian Strain CB4856. *Genetics* **200**,
883 975–989 (2015).
- 884 27. Lee, D. *et al.* Balancing selection maintains hyper-divergent haplotypes in
885 *Caenorhabditis elegans*. *Nat Ecol Evol* 1–14 (2021).
- 886 28. Bulteau, R. & Francesconi, M. Real age prediction from the transcriptome with
887 RAPToR. *bioRxiv* 2021.09.07.459270 (2021) doi:10.1101/2021.09.07.459270.
- 888 29. Crombie, T. A. *et al.* Deep sampling of Hawaiian *Caenorhabditis elegans* reveals
889 high genetic diversity and admixture with global populations. *Elife* **8**, e50465
890 (2019).
- 891 30. Bloom, J. S., Ehrenreich, I. M., Loo, W. T., Lite, T.-L. V. & Kruglyak, L. Finding the
892 sources of missing heritability in a yeast cross. *Nature* **494**, 234–237 (2013).
- 893 31. Chen, A., Liu, Y., Williams, S. M., Morris, N. & Buchner, D. A. Widespread epistasis
894 regulates glucose homeostasis and gene expression. *PLoS Genet.* **13**, e1007025
895 (2017).
- 896 32. Kang, H. M. *et al.* Efficient control of population structure in model organism
897 association mapping. *Genetics* **178**, 1709–1723 (2008).
- 898 33. Li, J. & Ji, L. Adjusting multiple testing in multilocus analyses using the eigenvalues
899 of a correlation matrix. *Heredity* **95**, 221–227 (2005).
- 900 34. Campos, T. L., Korhonen, P. K., Sternberg, P. W., Gasser, R. B. & Young, N. D.
901 Predicting gene essentiality in *Caenorhabditis elegans* by feature engineering and
902 machine-learning. *Comput. Struct. Biotechnol. J.* **18**, 1093–1102 (2020).
- 903 35. Cutter, A. D. Nucleotide polymorphism and linkage disequilibrium in wild
904 populations of the partial selfer *Caenorhabditis elegans*. *Genetics* **172**, 171–184
905 (2006).
- 906 36. Rockman, M. V. & Kruglyak, L. Recombinational landscape and population
907 genomics of *Caenorhabditis elegans*. *PLoS Genet.* **5**, e1000419 (2009).
- 908 37. Araya, C. L. *et al.* Regulatory analysis of the *C. elegans* genome with
909 spatiotemporal resolution. *Nature* **512**, 400–405 (2014).
- 910 38. Kudron, M. M. *et al.* The ModERN Resource: Genome-Wide Binding Profiles for

- 911 Hundreds of *Drosophila* and *Caenorhabditis elegans* Transcription Factors.
912 *Genetics* **208**, 937–949 (2018).
- 913 39. Holdorf, A. D. *et al.* WormCat: An Online Tool for Annotation and Visualization of
914 *Caenorhabditis elegans* Genome-Scale Data. *Genetics* **214**, 279–294 (2020).
- 915 40. Satterlee, J. S. *et al.* Specification of thermosensory neuron fate in *C. elegans*
916 requires *tx-1*, a homolog of *otd/Otx*. *Neuron* **31**, 943–956 (2001).
- 917 41. Kagoshima, H. & Kohara, Y. Co-expression of the transcription factors CEH-14
918 and TTX-1 regulates AFD neuron-specific genes *gcy-8* and *gcy-18* in *C. elegans*.
919 *Dev. Biol.* **399**, 325–336 (2015).
- 920 42. Fairfax, B. P. *et al.* Innate immune activity conditions the effect of regulatory
921 variants upon monocyte gene expression. *Science* **343**, 1246949 (2014).
- 922 43. Evans, K. S. *et al.* Natural variation in the sequestosome-related gene, *sqst-5*,
923 underlies zinc homeostasis in *Caenorhabditis elegans*. *PLoS Genet.* **16**, e1008986
924 (2020).
- 925 44. Hahnel, S. R. *et al.* Extreme allelic heterogeneity at a *Caenorhabditis elegans* beta-
926 tubulin locus explains natural resistance to benzimidazoles. *PLoS Pathog.* **14**,
927 e1007226 (2018).
- 928 45. Churchill, G. A., Gatti, D. M., Munger, S. C. & Svenson, K. L. The Diversity Outbred
929 mouse population. *Mamm. Genome* **23**, 713–718 (2012).
- 930 46. Cook, D. E. *et al.* The Genetic Basis of Natural Variation in *Caenorhabditis elegans*
931 Telomere Length. *Genetics* **204**, 371–383 (2016).
- 932 47. Zdraljevic, S. *et al.* Natural variation in *C. elegans* arsenic toxicity is explained by
933 differences in branched chain amino acid metabolism. *Elife* **8**, e40260 (2019).
- 934 48. Zdraljevic, S. *et al.* Natural variation in a single amino acid substitution underlies
935 physiological responses to topoisomerase II poisons. *PLoS Genet.* **13**, e1006891
936 (2017).
- 937 49. Na, H., Zdraljevic, S., Tanny, R. E., Walhout, A. J. M. & Andersen, E. C. Natural
938 variation in a glucuronosyltransferase modulates propionate sensitivity in a *C.*
939 *elegans* propionic acidemia model. *PLoS Genet.* **16**, e1008984 (2020).
- 940 50. Evans, K. S. *et al.* Two novel loci underlie natural differences in *Caenorhabditis*

- 941 elegans abamectin responses. *PLoS Pathog.* **17**, e1009297 (2021).
- 942 51. Lee, D. *et al.* Selection and gene flow shape niche-associated variation in
943 pheromone response. *Nat Ecol Evol* **3**, 1455–1463 (2019).
- 944 52. Zhang, G., Mostad, J. D. & Andersen, E. C. Natural variation in fecundity is
945 correlated with species-wide levels of divergence in *Caenorhabditis elegans*. *G3*
946 (2021) doi:10.1093/g3journal/jkab168.
- 947 53. Harris, T. W. *et al.* WormBase: a modern Model Organism Information Resource.
948 *Nucleic Acids Res.* **48**, D762–D767 (2020).
- 949 54. Thyagarajan, B. *et al.* ETS-4 is a transcriptional regulator of life span in
950 *Caenorhabditis elegans*. *PLoS Genet.* **6**, e1001125 (2010).
- 951 55. Petersen, C., Dirksen, P. & Schulenburg, H. Why we need more ecology for
952 genetic models such as *C. elegans*. *Trends Genet.* **31**, 120–127 (2015).
- 953 56. Sterken, M. G., Snoek, L. B., Kammenga, J. E. & Andersen, E. C. The laboratory
954 domestication of *Caenorhabditis elegans*. *Trends Genet.* **31**, 224–231 (2015).
- 955 57. Frézal, L. & Félix, M.-A. The natural history of model organisms: *C. elegans* outside
956 the Petri dish. *Elife* **4**, e05849 (2015).
- 957 58. Schulenburg, H. & Félix, M.-A. The Natural Biotic Environment of *Caenorhabditis*
958 *elegans*. *Genetics* **206**, 55–86 (2017).
- 959 59. Crombie, T. A. *et al.* Local adaptation and spatiotemporal patterns of genetic
960 diversity revealed by repeated sampling of *Caenorhabditis elegans* across the
961 Hawaiian Islands. *bioRxiv* 2021.10.11.463952 (2021)
962 doi:10.1101/2021.10.11.463952.
- 963 60. Thomas, J. H. Adaptive evolution in two large families of ubiquitin-ligase adapters
964 in nematodes and plants. *Genome Res.* **16**, 1017–1030 (2006).
- 965 61. de Nadal, E., Ammerer, G. & Posas, F. Controlling gene expression in response to
966 stress. *Nat. Rev. Genet.* **12**, 833–845 (2011).
- 967 62. Denver, D. R. *et al.* The transcriptional consequences of mutation and natural
968 selection in *Caenorhabditis elegans*. *Nat. Genet.* **37**, 544–548 (2005).
- 969 63. Rifkin, S. A., Houle, D., Kim, J. & White, K. P. A mutation accumulation assay
970 reveals a broad capacity for rapid evolution of gene expression. *Nature* **438**, 220–

- 971 223 (2005).
- 972 64. King, M. C. & Wilson, A. C. Evolution at two levels in humans and chimpanzees.
- 973 *Science* **188**, 107–116 (1975).
- 974 65. Oliver, F. *et al.* Regulatory variation at glycan-3 underlies a major growth QTL in
- 975 mice. *PLoS Biol.* **3**, e135 (2005).
- 976 66. Courtier-Orgogozo, V., Arnoult, L., Prigent, S. R., Wiltgen, S. & Martin, A.
- 977 Gephebase, a database of genotype-phenotype relationships for natural and
- 978 domesticated variation in Eukaryotes. *Nucleic Acids Res.* **48**, D696–D703 (2020).
- 979 67. Driscoll, M., Dean, E., Reilly, E., Bergholz, E. & Chalfie, M. Genetic and molecular
- 980 analysis of a *Caenorhabditis elegans* beta-tubulin that conveys benzimidazole
- 981 sensitivity. *J. Cell Biol.* **109**, 2993–3003 (1989).
- 982 68. Kwa, M. S., Veenstra, J. G. & Roos, M. H. Benzimidazole resistance in
- 983 *Haemonchus contortus* is correlated with a conserved mutation at amino acid 200
- 984 in beta-tubulin isotype 1. *Mol. Biochem. Parasitol.* **63**, 299–303 (1994).
- 985 69. Brady, S. C. *et al.* A Novel Gene Underlies Bleomycin-Response Variation in
- 986 *Caenorhabditis elegans*. *Genetics* **212**, 1453–1468 (2019).
- 987 70. Andersen, E. C., Bloom, J. S., Gerke, J. P. & Kruglyak, L. A variant in the
- 988 neuropeptide receptor npr-1 is a major determinant of *Caenorhabditis elegans*
- 989 growth and physiology. *PLoS Genet.* **10**, e1004156 (2014).
- 990 71. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ
- 991 preprocessor. *Bioinformatics* **34**, i884–i890 (2018).
- 992 72. Li, H. A statistical framework for SNP calling, mutation discovery, association
- 993 mapping and population genetical parameter estimation from sequencing data.
- 994 *Bioinformatics* **27**, 2987–2993 (2011).
- 995 73. Pertea, G. & Pertea, M. GFF Utilities: GffRead and GffCompare. *F1000Res.* **9**, 304
- 996 (2020).
- 997 74. Storer, J., Hubley, R., Rosen, J., Wheeler, T. J. & Smit, A. F. The Dfam community
- 998 resource of transposable element families, sequence models, and genome
- 999 annotations. *Mob. DNA* **12**, 2 (2021).
- 1000 75. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-

- 1001 seq quantification. *Nat. Biotechnol.* **34**, 525–527 (2016).
- 1002 76. Pimentel, H., Bray, N. L., Puente, S., Melsted, P. & Pachter, L. Differential analysis
1003 of RNA-seq incorporating quantification uncertainty. *Nat. Methods* **14**, 687–690
1004 (2017).
- 1005 77. Soneson, C., Love, M. I. & Robinson, M. D. Differential analyses for RNA-seq:
1006 transcript-level estimates improve gene-level inferences. *F1000Res.* **4**, 1521
1007 (2015).
- 1008 78. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and
1009 dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 1010 79. Core Team, R. & Others. R: A language and environment for statistical computing.
1011 Vienna, Austria: R Foundation for Statistical Computing. Available (2013).
- 1012 80. Barrière, A. & Félix, M.-A. Natural variation and population genetics of
1013 *Caenorhabditis elegans*. *WormBook* 1–19 (2005).
- 1014 81. Ortiz, E. M. *vcf2phylip v2.0: convert a VCF matrix into several matrix formats for*
1015 *phylogenetic analysis*. (2019). doi:10.5281/zenodo.2540861.
- 1016 82. Schliep, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592–593
1017 (2011).
- 1018 83. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and
1019 evolutionary analyses in R. *Bioinformatics* **35**, 526–528 (2019).
- 1020 84. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T. Ggtree : An r package for
1021 visualization and annotation of phylogenetic trees with their covariates and other
1022 associated data. *Methods Ecol. Evol.* **8**, 28–36 (2017).
- 1023 85. Purcell, S. et al. PLINK: a tool set for whole-genome association and population-
1024 based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 1025 86. Chang, C. C. et al. Second-generation PLINK: rising to the challenge of larger and
1026 richer datasets. *Gigascience* **4**, 7 (2015).
- 1027 87. Endelman, J. B. Ridge regression and other kernels for genomic selection with R
1028 package rrBLUP. *Plant Genome* **4**, 250–255 (2011).
- 1029 88. Wit, J., Rodriguez, B. C. & Andersen, E. C. Natural variation in *Caenorhabditis*
1030 *elegans* responses to the anthelmintic emodepside. *Int. J. Parasitol. Drugs Drug*

- 1031 *Resist.* **16**, 1–8 (2021).
- 1032 89. Covarrubias-Pazaran, G. Genome-Assisted Prediction of Quantitative Traits Using
1033 the R Package *sommer*. *PLoS One* **11**, e0156744 (2016).
- 1034 90. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting Linear Mixed-Effects Models
1035 Using lme4. *Journal of Statistical Software, Articles* **67**, 1–48 (2015).
- 1036 91. Millard, S. P. *EnvStats: An R Package for Environmental Statistics*. (Springer-Verlag
1037 New York, 2013).
- 1038 92. Miles, A., Ralph, P., Rae, S. & Pisupati, R. *cggf/scikit-allel: v1.2.1.* (2019).
1039 doi:10.5281/zenodo.3238280.
- 1040 93. Henikoff, S. & Henikoff, J. G. Amino acid substitution matrices from protein blocks.
1041 *Proc. Natl. Acad. Sci. U. S. A.* **89**, 10915–10919 (1992).
- 1042 94. Tingley, D., Yamamoto, T., Hirose, K., Keele, L. & Imai, K. mediation: R Package
1043 for Causal Mediation Analysis. *Journal of Statistical Software, Articles* **59**, 1–38
1044 (2014).

Supplementary Information

An atlas of gene expression variation across the *Caenorhabditis elegans* species

16

Description of Additional Supplementary Files

17

18 File Name: Supplementary Data 1

19 Description: Expression abundances and heritabilities of 25,849 transcripts in 207 wild
20 *C. elegans* strains

21

22 File Name: Supplementary Data 2

23 Description: GWA mapping results of 5,291 transcript expression traits, eQTL
24 classification, and distant eQTL hotspots

25

26 File Name: Supplementary Data 3

27 Description: Enrichment of genes with eQTL, distant eQTL in hotspots

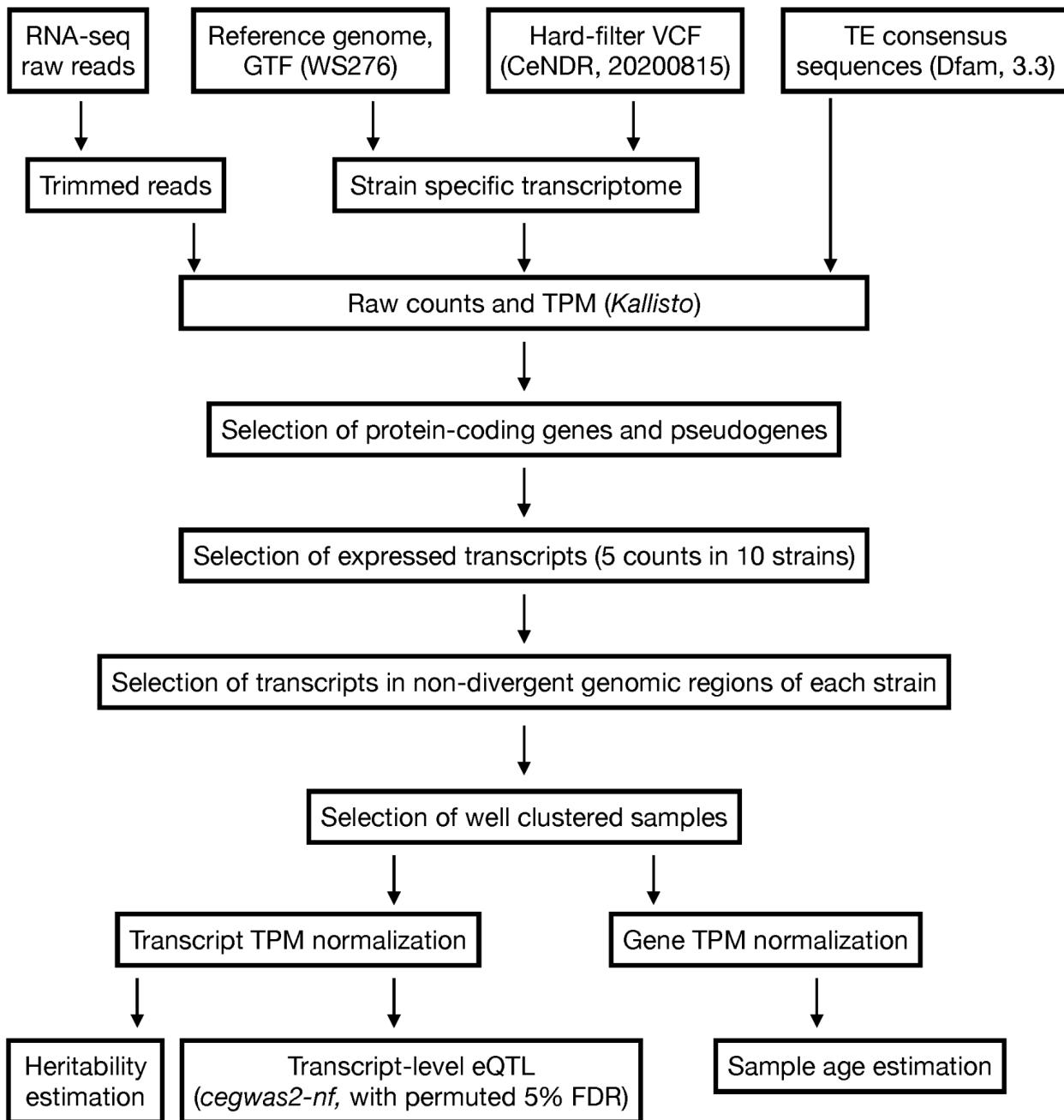
28

29 File Name: Supplementary Data 4

30 Description: Fine mappings for distant eQTL in hotspots

31 **Supplementary Figures**

Workflow of RNA-seq data analysis



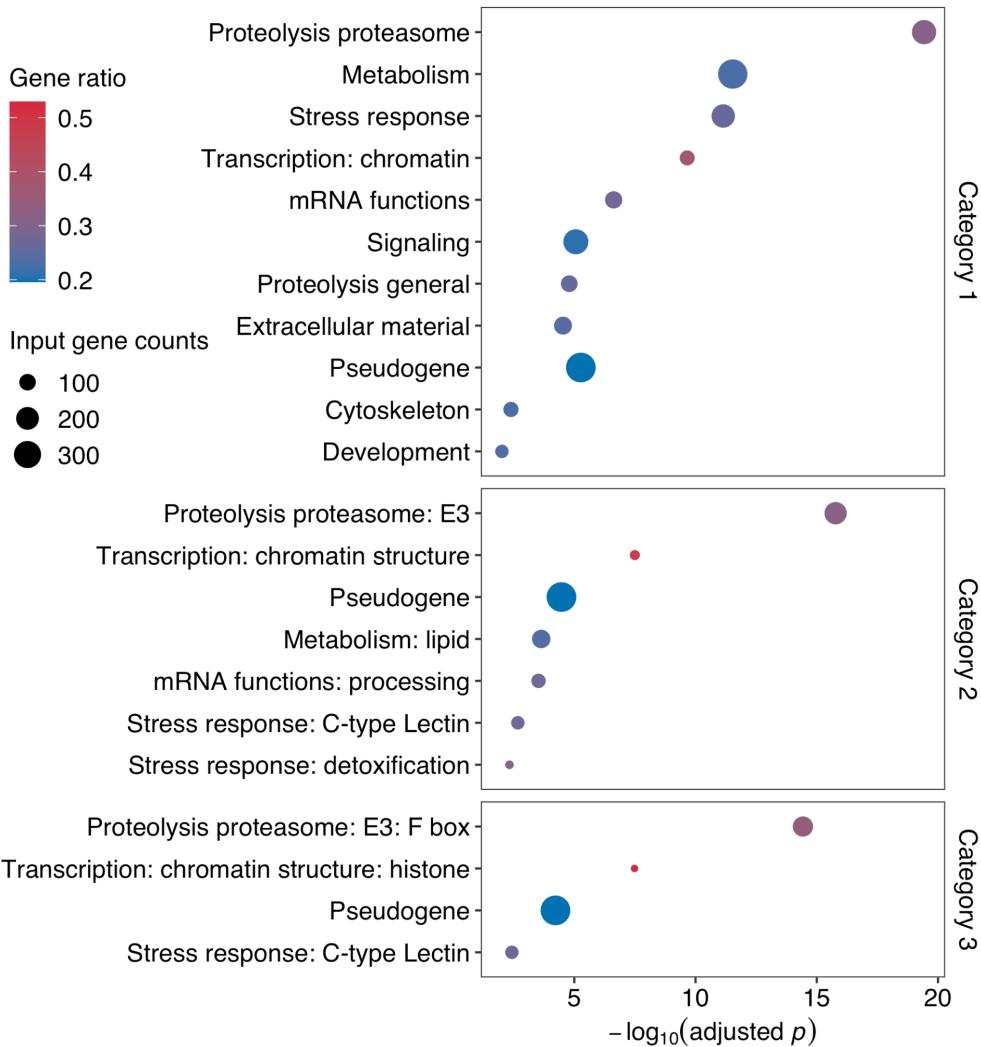
32

33 **Supplementary Fig. 1**

34 Workflow of RNA-seq data processing.

35

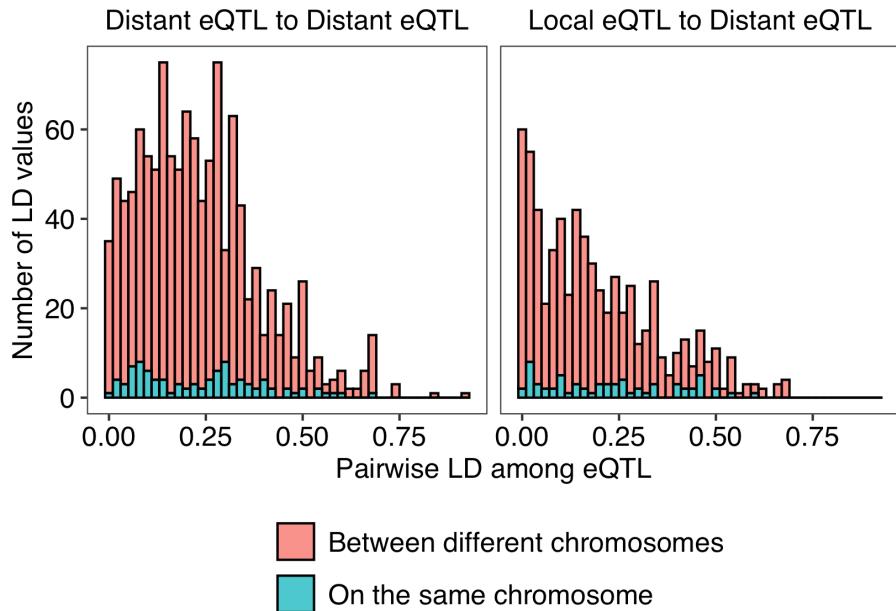
36



37

38 **Supplementary Fig. 2**

39 Gene set enrichment analysis for genes with transcript level eQTL. Enriched gene
40 classes of broad and specific categories (Category 1 to 3)³⁹ are shown on the y axis.
41 Bonferroni FDR corrected significance values using Fisher Exact Test for gene set
42 enrichment analysis are shown on the x axis. The sizes of the circles correspond to the
43 input gene counts of the annotation and the colors of the circles correspond to the gene
44 ratio of input gene counts to total gene counts of the annotation.

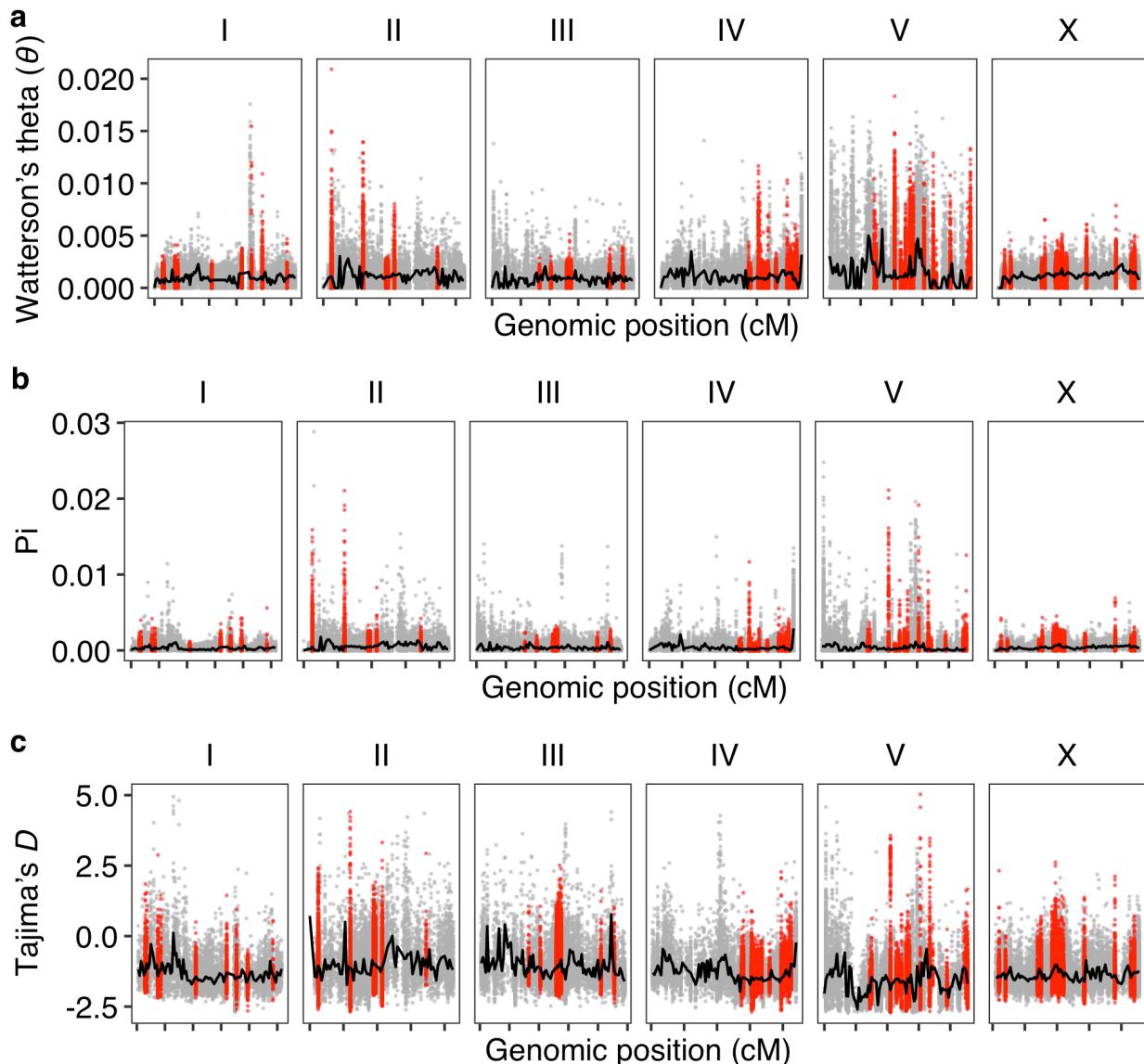


45

46 **Supplementary Fig. 3**

47 A histogram showing the distribution of linkage disequilibrium (LD) values (x-axis) among
48 QTL of multiple eQTL of transcript expression traits. A total of 861 traits were found with
49 multiple eQTL. LD of eQTL from the same chromosome and different chromosomes are
50 colored salmon and blue, respectively.

51

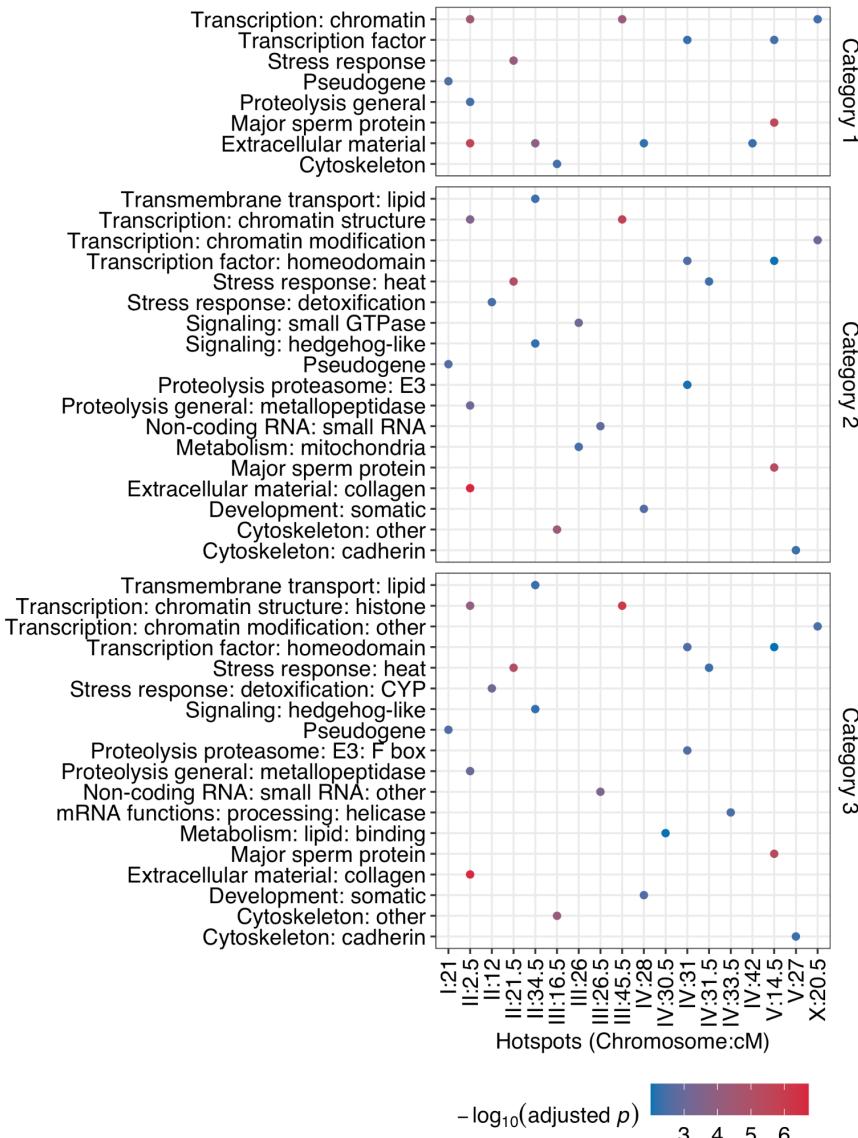


52

53 **Supplementary Fig. 4**

54 Genome-wide pattern of (a) Watterson's theta (θ), (b) nucleotide diversity (π), and (c)
55 Tajima's D . Each point represents the value (y-axis) for a non-overlapping 1 kb genomic
56 window and is plotted against the genome position (x-axis) with tick marks denoting
57 every 10 cM. Points for genomic windows in distant eQTL hotspots are colored red.
58 Other points are colored gray. Median values of each statistic in each 0.5 cM bin were
59 colored black. Tajima's D values that suggest purifying selection are outliers for most
60 values within a hotspot.

61

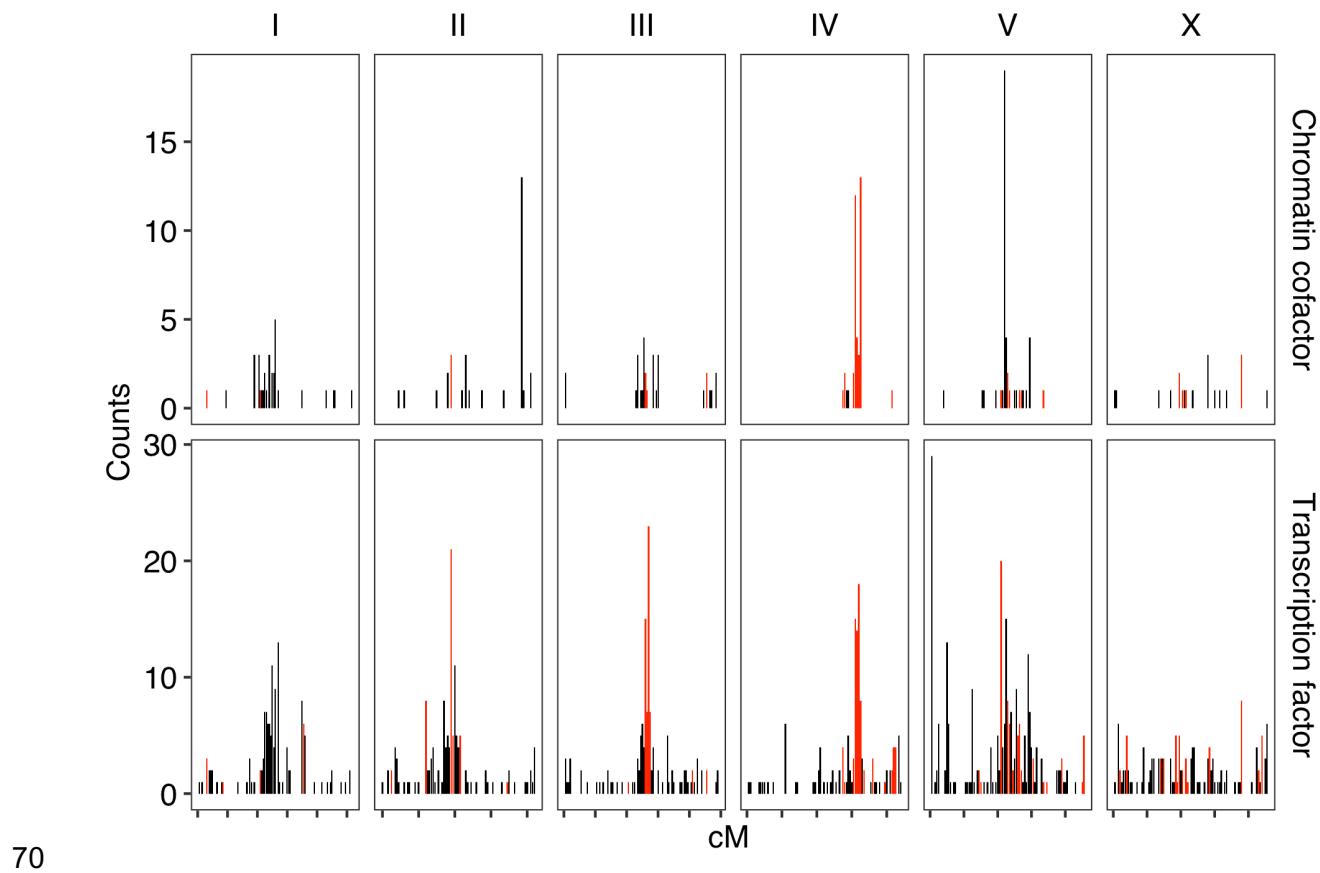


62

63

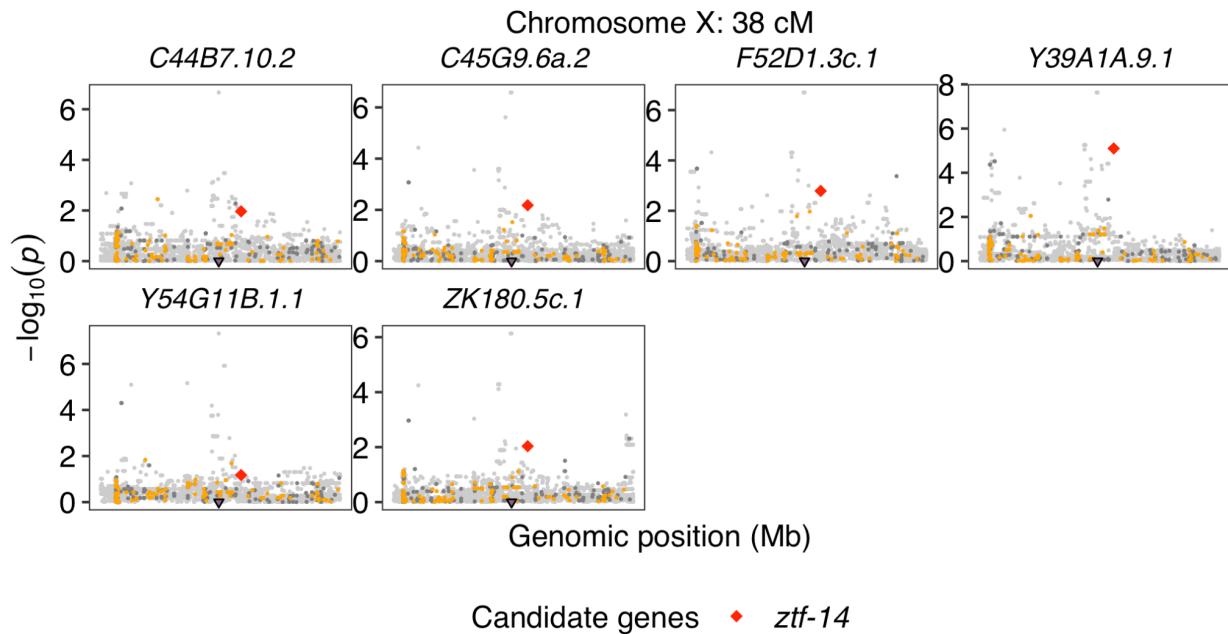
64 **Supplementary Fig. 5**

65 Gene set enrichment analysis for genes with transcript level distant eQTL in each
 66 hotspot. Broad and specific categories of enriched gene (Category 1 to 3)³⁹ are shown
 67 on the y axis. Distant eQTL hotspots with significant gene set enrichment are shown on
 68 the x axis. The colors of the circles correspond to Bonferroni FDR corrected significance
 69 values using Fisher Exact Test.



71 **Supplementary Fig. 6**

72 Number of genes encoding chromatin cofactors and transcription factors in each 0.5 cM
73 bins of the *C. elegans* genome. Bins that were identified as distant eQTL hotspots are
74 colored red. Other bins are colored black.
75

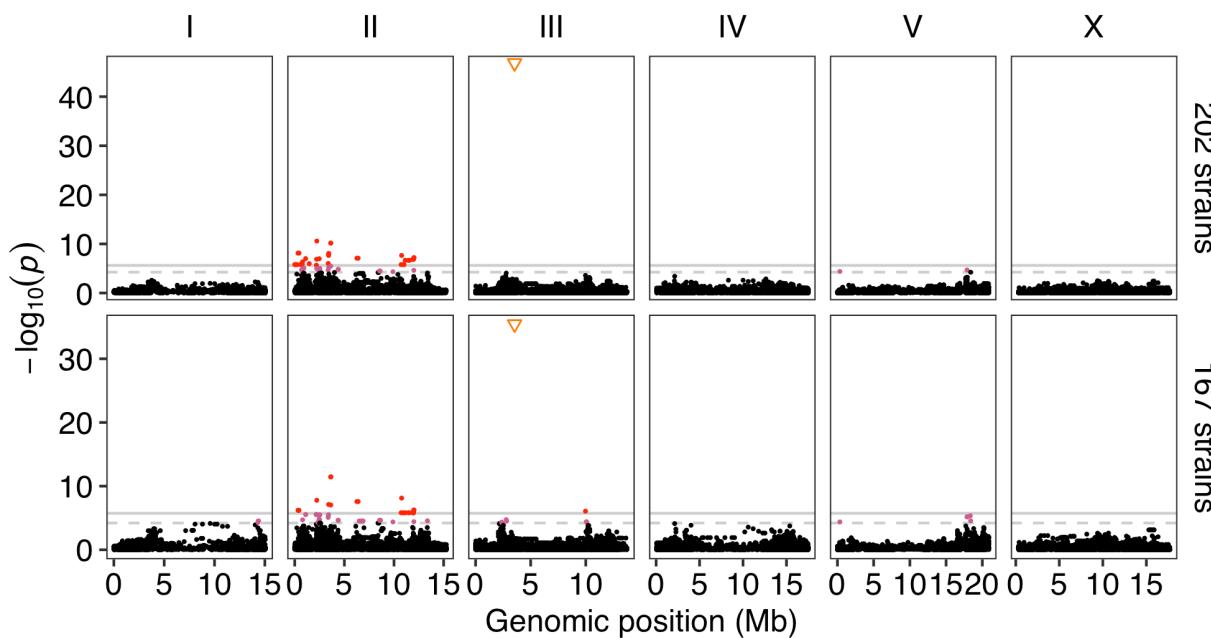


132

133 **Supplementary Fig. 7**

134 Fine mapping of transcript expression traits with distant eQTL in different hotspots is
135 shown. Genomic position (x-axis) is plotted against the $-\log_{10}(p)$ values (y-axis) for each
136 variant. Purple triangles on the x-axis represent eQTL positions. Candidate variants with
137 negative BLOSUM scores in genes encoding transcription factors or chromatin cofactors
138 are indicated as red diamonds. Other variants that are with negative BLOSUM scores,
139 with non-negative BLOSUM scores or intergenic are colored orange, dark gray, and light
140 gray, respectively. Transcript names of each trait are indicated above each panel.

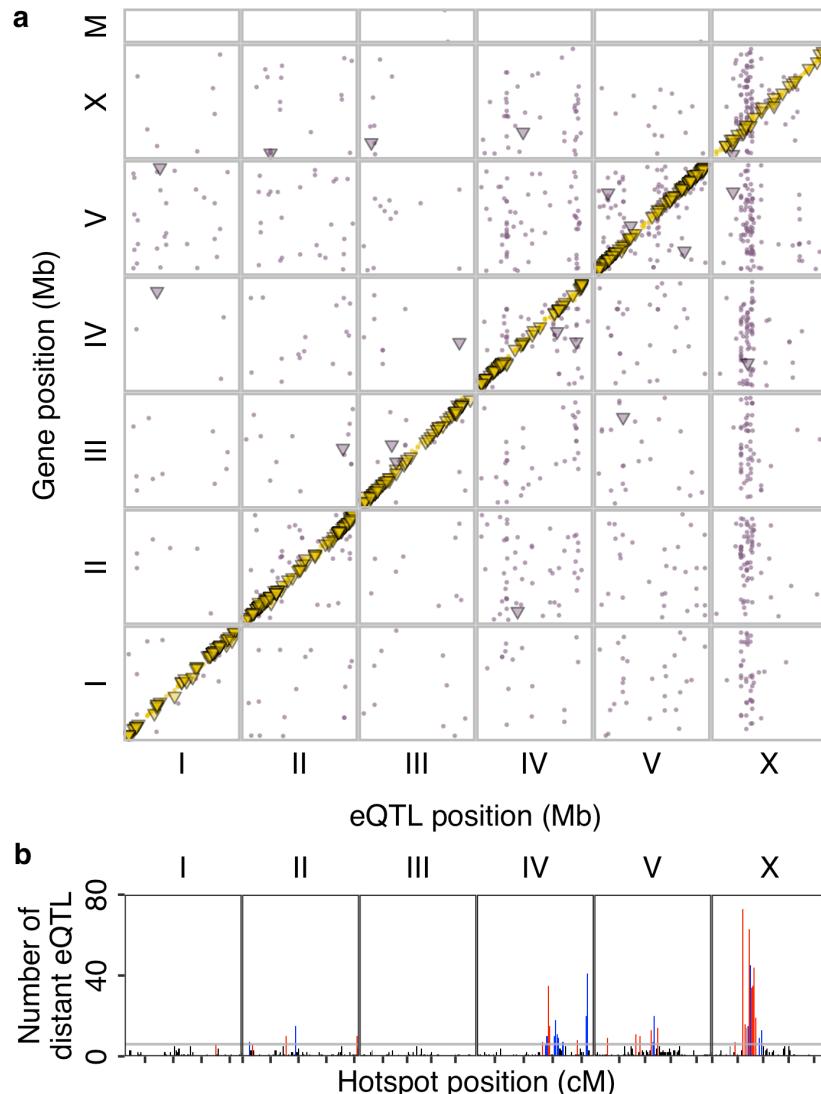
141



142

143 **Supplementary Fig. 8**

144 Manhattan plots indicating the GWA mapping result for animal length (q90.TOF) of 202
145 (top panel) and 167 (bottom panel) *C. elegans* wild strains in response to ABZ⁴⁴ are
146 shown. Each point represents an SNV that is plotted with its genomic position (x-axis)
147 against its $-\log_{10}(p)$ value (y-axis) from the GWA mapping. Real SNVs that pass the
148 genome-wide EIGEN threshold (the dotted gray horizontal line) and the genome-wide
149 Bonferroni threshold (the solid gray horizontal line) are colored pink and red, respectively.
150 The pseudo SNV marker representing high allelic heterogeneity in the gene *ben-1* at
151 position 3,539,640 on chromosome III is indicated as an orange inverted triangle.



152

153 **Supplementary Fig. 9**

154 **RIAILs eQTL.**

155 **a**, The genomic locations of 2,387 eQTL peaks (x-axis) in the RIAILs eQTL studies^{3,9} are
156 plotted against the genomic locations of the 2,003 genes with expression differences (y-
157 axis). Golden points or triangles on the diagonal of the map represent local eQTL. Purple
158 points or triangles correspond to distant eQTL. Triangles represent eQTL that were also
159 found in our study. **b**, The number of distant eQTL (y-axis) in each 0.5 cM bin across the
160 genome (x-axis) is shown. Tick marks on the x-axis denote every 10 cM. The horizontal
161 gray line indicates the threshold of 6 eQTL. Bins with 6 or more eQTL were identified as
162 hotspots and are colored red or blue. Bins with fewer than 6 eQTL are colored black.
163 Blue bins represent hotspots that were also found in our study.

164

165