1

2

3

4

5

6 **Comprehensive detection of structural variation and transposable element differences**

7 **between wild type laboratory lineages of *C. elegans***

8

9 Zachary D. Bush[1], Alice F. S. Naftaly[1], Devin Dinwiddie[1], Cora Albers[1], Kenneth J. Hillers[2], and

10 Diana E. Libuda[1]*

11

12 [1] Institute of Molecular Biology, Department of Biology, University of Oregon, 1229 Franklin Blvd
13 Eugene, OR 97403, USA

14 [2] Biological Sciences Department, California Polytechnic State University, San Luis Obispo,
15 California, USA
16
17 *Corresponding author

18

19 Corresponding Author and Lead Contact Information:

20 Diana E. Libuda, PhD
21 University of Oregon
22 Institute of Molecular Biology
23 1229 Franklin Blvd
24 Eugene, OR 97403
25 541-346-5092 (phone)
26 541-346-4854 (fax)
27 dlibuda@uoregon.edu

28

29 Running title: Structural variations in *C. elegans* genomes

30 Keywords: genome stability; sequence variation; reference genomes; genetic drift;

31 transposons; whole genome sequencing; *C. elegans*

**Abstract**

Genomic structural variations (SVs) and transposable elements (TEs) can be significant

contributors to genome evolution, altered gene expression, and risk of genetic diseases. Recent

advancements in long-read sequencing have greatly improved the quality of *de novo* genome

assemblies and enhanced the detection of sequence variants at the scale of hundreds or

thousands of bases. Comparisons between two diverged wild isolates of *Caenorhabditis*

*elegans*, the Bristol and Hawaiian strains, have been widely utilized in the analysis of small

genetic variations. Genetic drift, including SVs and rearrangements of repeated sequences such

as TEs, can occur over time from long-term maintenance of wild type isolates within the

laboratory.  To comprehensively detect both large and small structural variations as well as TEs

due to genetic drift, we generated *de novo* genome assemblies and annotations for each strain

from our lab collection using both long- and short-read sequencing and compared our

assemblies and annotations with that of other lab wild type strains. Within our lab assemblies,

we annotate over 3.1Mb of sequence divergence between the Bristol and Hawaiian isolates:

337,584 SNPs, 94,503 small insertion-deletions (<50bp), and 4,334 structural variations

(>50bp). Further, we define the location and movement of specific DNA TEs between N2 Bristol

and CB4856 Hawaiian wild type isolates.  Specifically, we find the N2 Bristol genome has 20.6%

more TEs from the *Tc1/mariner* family than the CB4856 Hawaiian genome. Moreover, we

identified Zator elements as the most abundant and mobile TE family in the genome.  Using

specific TE sequences with unique SNPs, we also identify 38 TEs that moved

intrachromosomally and 9 TEs that moved interchromosomally between the N2 Bristol and

CB4856 Hawaiian genomes.  By comparing the *de novo* genome assembly of our lab collection

Bristol isolate to the VC2010 Bristol assembly, we also reveal that lab lineages display over 2

Mb of total variation: 1,162 SNPs, 1,528 indels, and 897 SVs with 95% of the variation due to

SVs. Overall, our work demonstrates the unique contribution of SVs and TEs to variation and

57    genetic drift between wild type laboratory strains assumed to be isogenic despite growing

58    evidence of genetic drift and phenotypic variation.

59  **Author Summary**

60  For multiple model organisms, propagation of wild type strains in independent labs can lead to

61  multiple phenotypic differences over time. To assess recombination, map mutations, and

62  understand genomic changes during speciation, *Caenorhabditis elegans* researchers primarily

63  use the wild type isolates Bristol and Hawaiian. Here, we map structural variations,

64  transposable elements, and sequence divergence between the Bristol and Hawaiian natural

65  isolates and between genomes of different lab lineages of these same strains.

66

67 **Introduction**

68      Genomic variants, through mutation and recombination, in individuals and genetic drift in

69 populations underly the core process of evolution. Functional characterization of sequence

70 variants guides our understanding of phenotypic variances within species while also being

71 critical to identifying heritable disease-causing mutations (Haraksingh and Snyder 2013).

72 Genomic variation has been reported at multiple scales, from single nucleotide polymorphisms

73 (SNPs) to short insertions/deletions (indels) to much larger structural variants (SVs). SVs are

74 defined as insertions, deletions, or chromosomal rearrangements at least 50bp in length. SVs

75 can cause loss of function mutations through large gene deletions or alter gene expression by

76 disrupting spatial interactions between regulatory sequences (Stranger et al. 2007; Hurles,

77 Dermitzakis, and Tyler-Smith 2008). Accurate detection of both sequence variants and

78 chromosome rearrangements is critical for understanding how genomic variation may contribute

79 to phenotypic plasticity in individuals and populations of the same species.

80      Transposable elements (TEs) are a class of repetitive DNA sequences capable of

81 moving to new locations in the genome. TE mobility is a source of genomic structural variation

82 that can also alter gene expression (Girard and Freeling 1999; Slotkin and Martienssen 2007)

83 and drive, sometimes rapid, evolutionary changes within species (Van't Hof et al. 2016;

84 Feschotte and Pritham 2007). Notably, transposons account for a significant fraction of the total

85 DNA sequence in many eukaryotic species (Chalopin et al. 2015; Gilbert, Peccoud, and

86 Cordaux 2021), which provides many opportunities for TE-driven structural rearrangements.

87 The *Tc1/mariner* family of DNA transposons is one of the most abundant TEs across species

88 (Eide and Anderson 1985; Plasterk, Izsvák, and Ivics 1999), and early studies in *C. elegans*

89 found it to be one of the few mobile transposons observed under laboratory conditions (Fischer,

90 Wienholds, and Plasterk 2003). To repress or limit transposon mobilization, transposon

91 silencing is tightly regulated through multiple mechanisms including chromatin modification and

92 RNA interference (Sijen and Plasterk 2003; H.-C. Lee et al. 2012). Despite their ubiquity and

5

93    impact on genomic architecture, the comprehensive annotation and inclusion of TEs in

94    comparative genomic analyses has been challenging. Many studies have incompletely

95    characterized the genomic distribution of TEs because older, short-read based genome

96    assemblies could not accurately map the full content of repetitive sequences. Further, programs

97    that automatically detect TEs based on sequence homology and conserved sequence elements

98    rely heavily on libraries of older reference sequences that may predate the discovery of TE

99    fragments and newer TE families. As new families of transposable elements are discovered

100   (Bao et al. 2009) along with new technology that aids their annotation and tracking (Riehl et al.

101   2022), determining the genomic composition and mobility of new TEs will enable our

102   understanding of their role in genome evolution and genome integrity.

103        Foundational research on genomic variation has utilized next generation short-read

104   sequencing, long-read sequencing, and the direct comparison of reference genome assemblies

105   to identify genomic variants (Mahmoud et al. 2019; Lappalainen et al. 2019). SNPs and indels,

106   ranging in size from 1 bp to 50bp, can be identified with high confidence using short sequencing

107   reads that are 100-150bp (Muzzey, Evans, and Lieber 2015). In contrast, SVs are challenging to

108   annotate using short-read sequencing because the sequencing reads are often smaller than the

109   size of an SV (Sudmant et al. 2015; Mahmoud et al. 2019; Lesack et al. 2022). Similarly, the

110   highly repetitive sequences of TEs present significant challenges to mapping and annotation

111   with traditional short read sequencing methods. With the advent of higher quality long-read

112   sequencing technologies which generate ~10kb-30kb reads with lower genomic coverage, the

113   accurate annotation of large regions of genomic variation such as SVs and transposable

114   elements has become easier (Sakamoto et al. 2021). New tools to identify SVs via assembly-to-

115   assembly alignments (Delcher et al. 1999; Nattestad and Schatz 2016; Li 2018; Goel et al.

116   2019)  are not constrained by read-length to identify SVs and depend on high-quality reference

117   assemblies. Thus, a high-quality reference genome assembly is a critical resource for any

118   model organism. Methods of variant detection that leverage a combined utilization of short- and

119   long-read sequencing can provide more accurate reference sequences to fully address

120   undiscovered genomic variations previously not detected by short-read sequencing alone.

121      *Caenorhabditis elegans* was the first multicellular organism to have its genome fully

122   sequenced (C. elegans Sequencing Consortium 1998) and has been exploited to pioneer many

123   comparative genomic studies. To understand how genetic variation influences phenotypic

124   differences and genomic processes within species, *C. elegans* researchers primarily utilize two

125   highly diverged wild type strains estimated to have diverged 30,000-50,000 generations ago

126   (Thomas et al. 2015): N2 (isolated in Bristol, England) and CB4856 (isolated in Maui, Hawaii)

127   (Nicholas, Dougherty, and Hansen 1959; Sulston and Brenner 1974; Hodgkin and Doniach

128   1997; Crombie et al. 2019). Earlier comparisons of the Bristol and Hawaiian lineages were

129   critical for studying genetic variation, gene families, and evolution of genome structures (Koch et

130   al. 2000; Wicks et al. 2001; Stewart et al. 2005; Maydan et al. 2010). The *C. elegans* genome,

131   comprised of 5 autosomes and the X chromosome, displays a nonuniform distribution of

132   sequence variation when comparing the genomes of wild isolates. Although a large amount of

133   sequence divergence was previously found between the N2 Bristol and CB4856 Hawaiian

134   lineages (Thompson et al. 2015; Andersen et al. 2012), the increased quality of reference

135   genomes, sequencing technology, and variant detection methods enables the identification of

136   additional variations (in particular large structural variations) that previously went undetected in

137   these *C. elegans* genomes.

138      Recently, Bristol and Hawaiian genomes were reassembled *de novo* using a

139   combination of short-read Illumina sequencing as well as long-read sequencing from PacBio

140   and Oxford Nanopore platforms (Yoshimura et al. 2019; Kim et al. 2019). Compared to the

141   previous short-read based assemblies of N2 Bristol, the new assembly of N2 Bristol, called

142   VC2010, identified 53 more predicted genes, 1.8Mb of additional sequence, and eliminated 98%

143   of existing gaps in the N2 Bristol genome. Thus, the VC2010 Bristol genome very likely better

144   represents the genome of Bristol *C. elegans* currently used in laboratories worldwide

145   (Yoshimura et al. 2019). The first CB4856 Hawaiian genome assembly was completed in 2015

146   by iteratively correcting the pre-existing N2 Bristol reference assembly (*C. elegans* Sequencing

147   Consortium 1998) with short-read sequencing data (Thompson et al. 2015). This study identified

148   327,050 single-nucleotide polymorphisms (SNPs) and nearly 80,000 indels relative to N2; a

149   marked increase relative to previous comparisons, which had identified 6,000-17,000 SNPs and

150   small indels (Wicks et al. 2001; Swan et al. 2002) between N2 Bristol and CB4856 Hawaiian.

151   Due to the size of the short read sequences employed in the analysis, the iterative correction

152   method used to assemble the CB4856 Hawaiian genome may not have detected all structural

153   rearrangements and repetitive sequences. In 2019, the first *de novo* CB4856 Hawaiian

154   assembly from long-read sequencing extended the length of the Hawaiian genome, and was

155   further able to characterize over 3,000 previously uncharacterized SVs (Kim et al. 2019). Thus,

156   combining long-read and short-read sequencing in *de novo* genome assembly not only

157   extended the known length of both the N2 Bristol and CB4856 Hawaiian isolate genomes, but

158   broadened our understanding of how much genomic variation exists between these wild-type

159   strains.

160        Many *C. elegans* research labs utilize N2 Bristol and CB4856 Hawaiian as standard wild

161   type strains, but long-term passaging in each lab may lead to the accumulation of many smaller

162   sequence variants and large genomic structural variations. Early assessments of laboratory

163   lineages of the N2 Bristol strain, for example, identified many duplications ranging in size from

164   200bp to 108kb, with some affecting as many as 26 genes (Vergara et al. 2009). To determine

165   the extent of genetic variation between our laboratory lineages of N2 Bristol and CB4856

166   Hawaiian, we generated two high-quality reference assemblies for the N2 and CB4856 strains

167   used in our laboratory to compare to that of other high-quality reference genomes for N2 and

168   CB4856. By leveraging recent technological advancements in sequencing and variant detection,

169   we provide a comprehensive annotation of  SNPs, indels, structural variations, and transposable

170   elements between our lineages of the Bristol and Hawaiian strains. From our comprehensive

171  mapping of TEs in our reference genomes, we report Zator elements to be the most abundant

172  and mobile TE family in the *C. elegans* genome.  Further, by comparing our assembled

173  genomes to recently published VC2010 Bristol and CB4856 Hawaiian long-read assemblies

174  (Yoshimura et al. 2019; C. Kim et al. 2019), we identified SNPs, indels, and SVs unique to

175  different lab wild type strains. These variations were enriched in intergenic regions of the *C.*

176  *elegans* genome, suggesting that variations in regulatory sequences and other non-coding

177  regions may underlie the phenotypic variances previously observed between laboratory strains.

178  Taken together, our systematic analysis of genetic variation between natural and laboratory wild

179  type isolates highlights the impact of large structural variants, TE composition, and other

180  chromosomal rearrangements accumulating in the genomes of laboratory model organisms.

181

182  **Results**

183  ***De novo* genome assembly using combined long and short-read sequencing produces**

184  **high quality genomes**

185  To perform systematic comparisons of multiple wild type genomes from different

186  laboratory isogenic strains, we generated *de novo* assemblies of N2 Bristol and CB4856

187  Hawaiian. The N2 Bristol genome was assembled from PacBio long-reads with 136x coverage

188  producing 121 contigs and a 100.4Mb genome (Figure 1A) The CB4856 Hawaiian genome was

189  generated from PacBio long-reads with 132x coverage from 169 contigs to give a 98.8Mb

190  assembly (Figure 1B). These long-read assemblies were then supplemented with Illumina

191  paired end short-reads with a sequencing depth of 540x and 628x for N2 Bristol and CB4856

192  Hawaiian respectively (Figure 1A-B).

193  To assess the quality of our reference genomes, we examined assembly-to-assembly

194  alignments and the orthologous gene content for each assembly. A strong assembly would

195  show a similar proportion of aligned bases and a high degree of synteny when comparing

196  across assemblies. In concordance with comparisons in earlier studies, 99.2% of bases across

197    our N2 Bristol and CB4856 Hawaiian assembled genomes were aligned (Kim et al. 2019), and

198    more than 92.2% of bases within alignments were syntenic (Table 1).
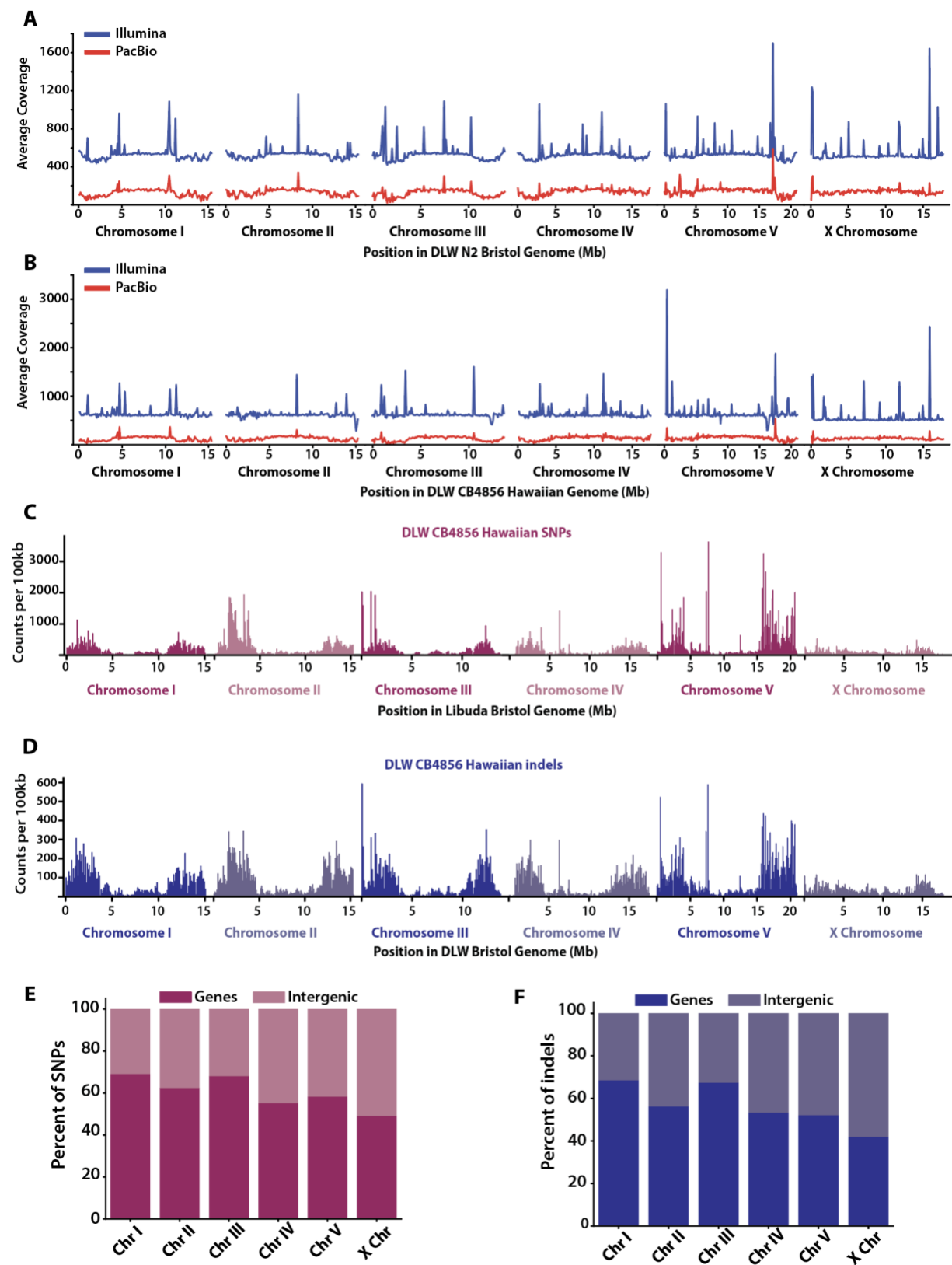
199

**Table 1. Comparisons between the DLW N2 Bristol genome (this study) and DLW CB4856 Hawaiian genome (this study)**

| | Chromosomes | | | | | | Total |
|---|---|---|---|---|---|---|---|
| | **I** | **II** | **III** | **IV** | **V** | **X** | |
| **DLW N2 Bristol Chromosome Length (this study)** | 15,114,068 | 15,311,845 | 13,819,453 | 17,493,838 | 20,953,657 | 17,739,129 | 100,431,990 |
| **DLW CB4856 Hawaiian Chromosome Length (this study)** | 15,045,644 | 15,257,363 | 13,206,755 | 17,183,882 | 20,547,529 | 17,584,915 | 98,826,088 |
| **N2 Bristol Bases Aligned** | 15,100,574 | 15,303,320 | 13,222,676 | 17,330,119 | 20,947,147 | 17,738,394 | 99,642,230 (99.21%) |
| **% Syntenic Aligned Bases** | 93.31 | 88.56 | 90.61 | 95.42 | 87.04 | 98.73 | 92.23 |
| **SNPs*** | 30,394 | 48,365 | 29,881 | 30,497 | 87,300 | 19,861 | 246,298 |
| **Indels*** | 11,460 | 13,716 | 10,530 | 11,221 | 20,063 | 6,799 | 73,789 (275,442 bp) |
| **SVs** | 863 | 808 | 649 | 619 | 925 | 470 | 4,334 (2,654,902 bp) |
| **HDRs** | 185 | 270 | 165 | 138 | 356 | 60 | 1,174 (6,864,884 bp) |

200
201    * All variants listed are only those for which the DLW CB4856 Hawaiian genome was
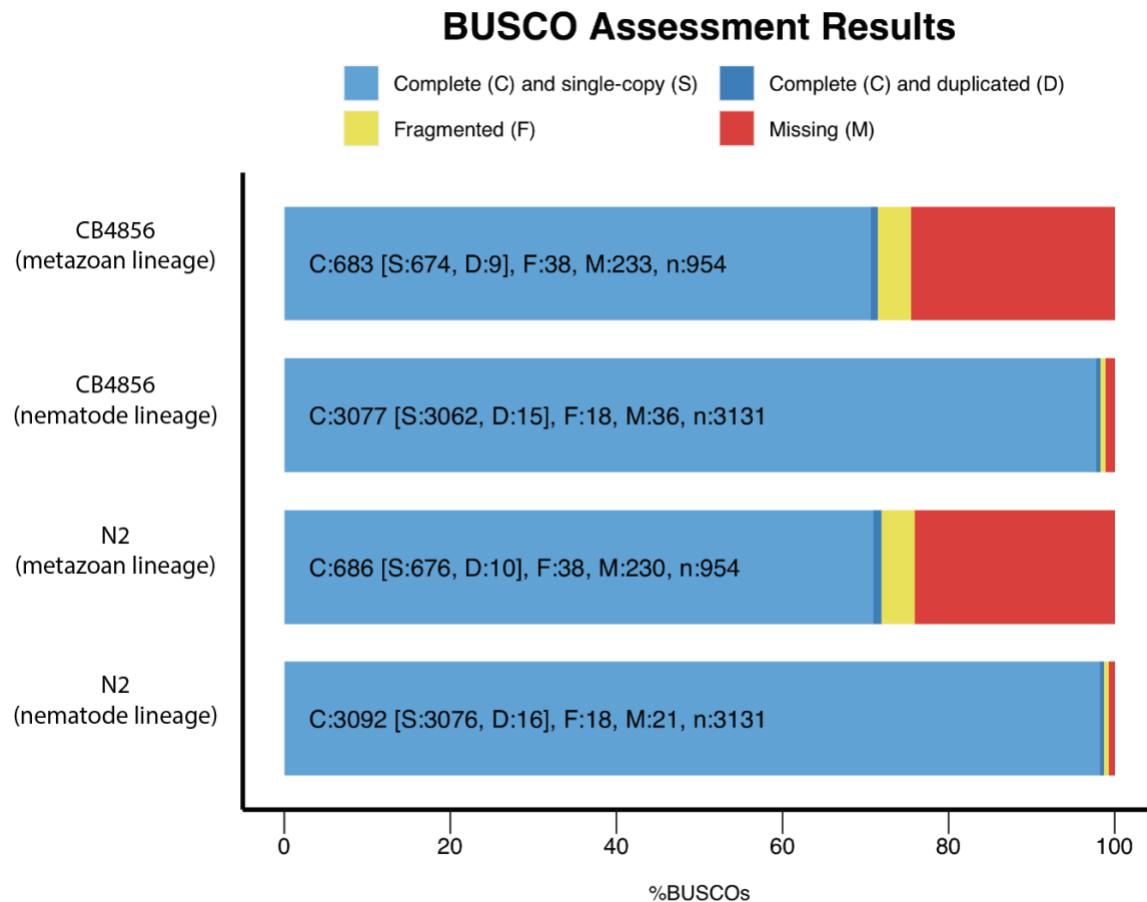202    homozygous

203
204
205

206 **Figure 1. Genomic distribution of SNPs and indels between the DLW N2 Bristol and DLW**
207 **CB4856 Hawaiian genomes. (A)** Line plots showing the average sequencing coverage in
208 100kb bins across each chromosome in the DLW N2 Bristol genome. **(B)** Line plots showing the
209 average sequencing coverage in 100kb bins across each chromosome in the DLW CB4856
210 Hawaiian genome. For each plot in A and B, the coverage for Illumina short-read sequencing is
211 shown in blue, and sequencing coverage for PacBio long-reads is shown in red. **(C)** Histograms
212 depicting the distribution of CB4856 Hawaiian SNPs across each DLW N2 Bristol chromosome
213 in 100kb bins. **(D)** Histograms of the distributions of CB4856 Hawaiian indels across each DLW
214 N2 Bristol chromosome in 100kb bins. **(E)** The proportion of SNPs that overlap with remapped
215 gene annotations versus intergenic regions in the DLW N2 Bristol genome. **(F)** The proportion of
216 indels that overlap with gene versus intergenic regions in the Bristol genome.
217

218 Analysis of universal single-copy orthologs (Simão et al. 2015; Manni et al. 2021) in our *de novo*

219 N2 Bristol and CB4856 Hawaiian genomes revealed greater than 98% completeness

220 (Supplemental Figure S1) and validate that our assemblies are high quality.

221

222 ***De novo* genome assemblies of the N2 Bristol and CB4856 Hawaiian isolates enhance**

223 **detection of genomic variation**

224 Previous comparisons of the genetic variation between N2 Bristol and CB4856 Hawaiian

225 have relied on a short-read N2 Bristol reference genome (Thompson et al. 2015; Kim et al.

226 2019), and the amount of variation has yet to be re-assessed using a modern long-read N2

227 Bristol assembly. Utilizing our N2 Bristol and CB4856 Hawaiian strains, we aligned CB4856

228 Hawaiian short reads to our N2 Bristol assembly. This analysis revealed a total of 246,298

229 homozygous SNPs and 73,789 homozygous indels across the genome (Table 1, Figure 1C-D).

230 While many of these SNPs and indels overlapped with gene annotations, they were under-

231 enriched in gene sequences (Figure 1D-E, Supplemental Figure 2). To identify large sequence

232 variants and chromosome rearrangements, we used whole-genome alignments (see Methods).

233 We identified a total of 4,364 structural variants, which are categorized as insertions, deletions,

234 and other chromosomal rearrangements spanning at least 50bp.

235

236

## BUSCO Assessment Results

**Legend:**
- Complete (C) and single-copy (S)
- Complete (C) and duplicated (D)
- Fragmented (F)
- Missing (M)

CB4856 (metazoan lineage): C:683 [S:674, D:9], F:38, M:233, n:954

CB4856 (nematode lineage): C:3077 [S:3062, D:15], F:18, M:36, n:3131

N2 (metazoan lineage): C:686 [S:676, D:10], F:38, M:230, n:954

N2 (nematode lineage): C:3092 [S:3076, D:16], F:18, M:21, n:3131
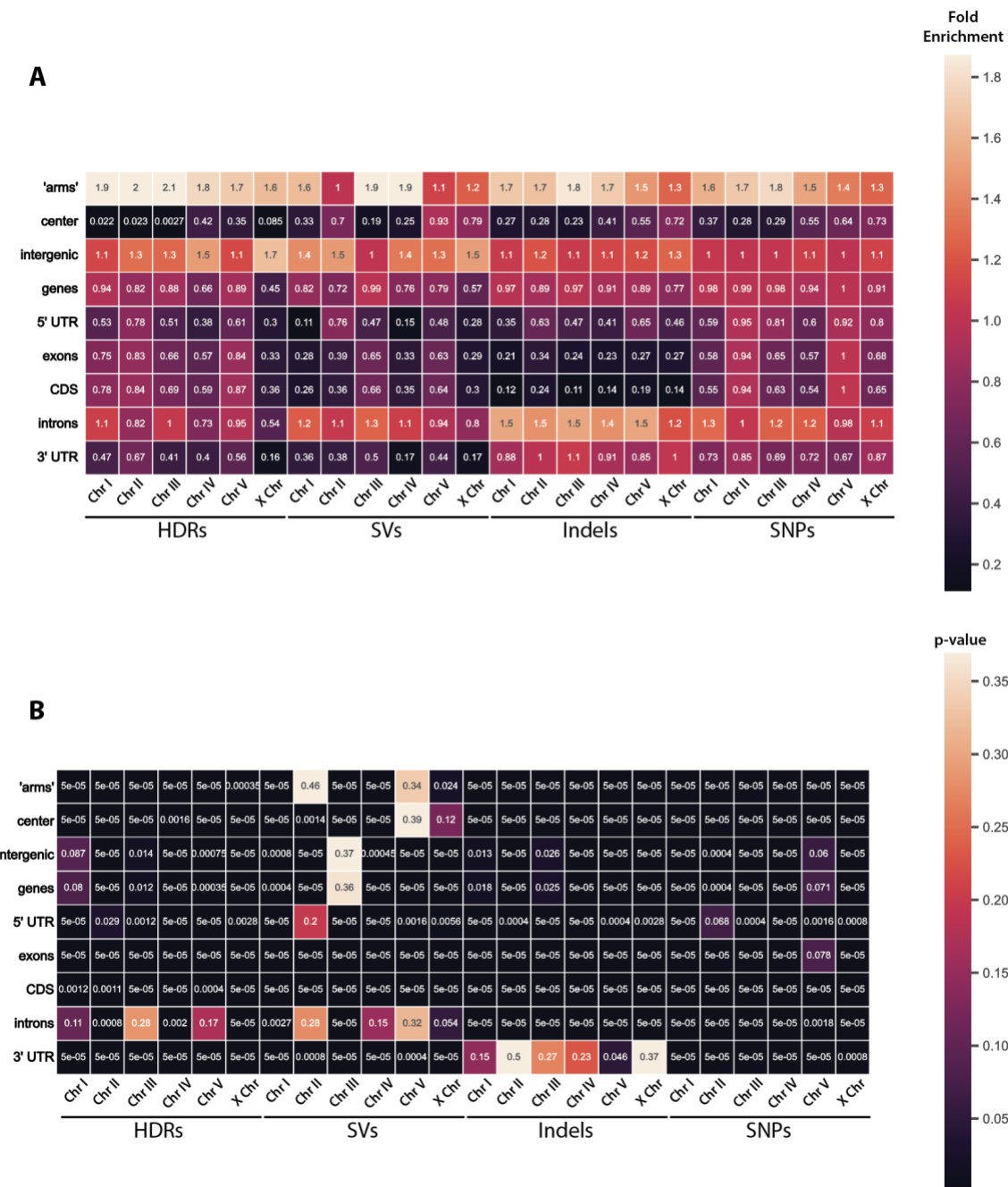
%BUSCOs

**Supplemental Figure S1**. BUSCO analysis of the DLW N2 Bristol and DLW CB4856 Hawaiian genome assemblies. The presence of orthologous genes from metazoan and nematode lineages are shown for each genome assembly. Each orthologous gene analyzed is depicted as either Complete (C, blues), Fragmented (F, yellow), or Missing (M, red). Complete orthologs are then further categorized as single-copy (S, light blue) or duplicated (D, dark blue).

We also identified 1,174 Highly Divergent Regions (HDRs) (Goel et al. 2019) across the genome. HDRs are defined as regions of the genome over 50bp in length that result in low-quality pairwise alignments due to the presence of multiple gaps within these alignments (Goel et al. 2019). Overall, greater than 9.9% of the DLW N2 Bristol genome (~10.0Mb) displayed variation through SNPs, indels, SVs, and HDRs when compared to the DLW CB4856 Hawaiian genome. SVs and HDRs represented only 1.3% and 0.3%, of variant sites between N2 Bristol and CB4856 Hawaiian respectively, but accounted for over 94% (9.5Mb) of sequence variation (Table1). Including heterozygous variants, our short-read analysis detected 3% more SNPs and

13

251     18% more indels than previously discovered using short-read assemblies of N2 Bristol and

252     CB4856 Hawaiian (Thompson et al. 2015). Utilizing whole-genome alignment comparisons (Li

253     2018; Goel et al. 2019), we identified 985 more SV sites than previously reported (Nattestad

254     and Schatz 2016; Kim et al. 2019). This increased sensitivity in variant site detection highlights

255     the power of combining long-read and short-read sequencing to create accurate genome

256     assemblies for comparative genomic studies.

257            Given an enhanced detection of variant sites between our N2 Bristol and CB4856

258     Hawaiian assemblies, we were interested in the genome-wide distribution of all variant sites.

259     Given previous reports (Thompson et al. 2015; Kim et al. 2019), we expected a greater density

260     of variation in the terminal thirds (the "arm-like" regions) of each chromosome. Indeed, there is a

261     significant concentration of SNPs, indels, SVs and HDRs in the arm-like regions relative to the

262     central region of each chromosome (Supplemental Figure S2). Over 78% of all SNPs, indels,

263     SVs, and HDRs are in the arm-like domains of each chromosome (Genome-wide averages:

264     75.12% of SNPs, 78.24% of indels, 71.39% of SVs, 90.77% of HDRs). To determine if the

265     enrichment of SNPs, indels, and SVs in the chromosomal arm-like regions was significant, we

266     compared the observed distribution of each variant category with random permutations of each

267     category of variant (Heger et al. 2013). SNPs, indels and HDRs on the autosomes were

268     significantly enriched in the arm-like regions (SNPs: 1.36-1.77 fold enrichment; Indels: 1.47-1.84

269     fold enrichment; HDRs: 1.70-2.06 fold enrichment; $p < 0.001$ by hypergeometric test). SVs,

270     however, were only significantly enriched on the arm-like regions of autosomes I, III, and IV

271     (1.64-1.92 fold enrichment; $p<.001$ by hypergeometric test). The fold enrichment of all variants

272     on the arm-like regions of the X chromosome was slightly weaker, ranging from 1.23-1.64

273     (SNPs: 1.26 fold enrichment; Indels: 1.26 fold enrichment; SVs: 1.23 fold enrichment; HDRs:

274     1.64 fold enrichment; all p-values $< 0.05$ by hypergeometric test). Similar to previous

275     observations (Thompson et al. 2015), there were a few hyper-variable regions with a greater

276

**A**

Fold Enrichment

| | ChrI | ChrII | ChrIII | ChrIV | ChrV | XChr | ChrI | ChrII | ChrIII | ChrIV | ChrV | XChr | ChrI | ChrII | ChrIII | ChrIV | ChrV | XChr | ChrI | ChrII | ChrIII | ChrIV | ChrV | XChr |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 'arms' | 1.9 | 2 | 2.1 | 1.8 | 1.7 | 1.6 | 1.6 | 1 | 1.9 | 1.9 | 1.1 | 1.2 | 1.7 | 1.7 | 1.8 | 1.7 | 1.5 | 1.3 | 1.6 | 1.7 | 1.8 | 1.5 | 1.4 | 1.3 |
| center | 0.022 | 0.023 | 0.0027 | 0.42 | 0.35 | 0.085 | 0.33 | 0.7 | 0.19 | 0.25 | 0.93 | 0.79 | 0.27 | 0.28 | 0.23 | 0.41 | 0.55 | 0.72 | 0.37 | 0.28 | 0.29 | 0.55 | 0.64 | 0.73 |
| intergenic | 1.1 | 1.3 | 1.3 | 1.5 | 1.1 | 1.7 | 1.4 | 1.5 | 1 | 1.4 | 1.3 | 1.5 | 1.1 | 1.2 | 1.1 | 1.1 | 1.2 | 1.3 | 1 | 1 | 1 | 1.1 | 1 | 1.1 |
| genes | 0.94 | 0.82 | 0.88 | 0.66 | 0.89 | 0.45 | 0.82 | 0.72 | 0.99 | 0.76 | 0.79 | 0.57 | 0.97 | 0.89 | 0.97 | 0.91 | 0.89 | 0.77 | 0.98 | 0.99 | 0.98 | 0.94 | 1 | 0.91 |
| 5' UTR | 0.53 | 0.78 | 0.51 | 0.38 | 0.61 | 0.3 | 0.11 | 0.76 | 0.47 | 0.15 | 0.48 | 0.28 | 0.35 | 0.63 | 0.47 | 0.41 | 0.65 | 0.46 | 0.59 | 0.95 | 0.81 | 0.6 | 0.92 | 0.8 |
| exons | 0.75 | 0.83 | 0.66 | 0.57 | 0.84 | 0.33 | 0.28 | 0.39 | 0.65 | 0.33 | 0.63 | 0.29 | 0.21 | 0.34 | 0.24 | 0.23 | 0.27 | 0.27 | 0.58 | 0.94 | 0.65 | 0.57 | 1 | 0.68 |
| CDS | 0.78 | 0.84 | 0.69 | 0.59 | 0.87 | 0.36 | 0.26 | 0.36 | 0.66 | 0.35 | 0.64 | 0.3 | 0.12 | 0.24 | 0.11 | 0.14 | 0.19 | 0.14 | 0.55 | 0.94 | 0.63 | 0.54 | 1 | 0.65 |
| introns | 1.1 | 0.82 | 1 | 0.73 | 0.95 | 0.54 | 1.2 | 1.1 | 1.3 | 1.1 | 0.94 | 0.8 | 1.5 | 1.5 | 1.5 | 1.4 | 1.5 | 1.2 | 1.3 | 1 | 1.2 | 1.2 | 0.98 | 1.1 |
| 3' UTR | 0.47 | 0.67 | 0.41 | 0.4 | 0.56 | 0.16 | 0.36 | 0.38 | 0.5 | 0.17 | 0.44 | 0.17 | 0.88 | 1 | 1.1 | 0.91 | 0.85 | 1 | 0.73 | 0.85 | 0.69 | 0.72 | 0.67 | 0.87 |
| | | | HDRs | | | | | | SVs | | | | | | Indels | | | | | | SNPs | | | |

Fold Enrichment scale: 0.2 – 1.8

**B**

p-value

| | ChrI | ChrII | ChrIII | ChrIV | ChrV | XChr | ChrI | ChrII | ChrIII | ChrIV | ChrV | XChr | ChrI | ChrII | ChrIII | ChrIV | ChrV | XChr | ChrI | ChrII | ChrIII | ChrIV | ChrV | XChr |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 'arms' | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 0.00035 | 5e-05 | 0.46 | 5e-05 | 5e-05 | 0.34 | 0.024 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 |
| center | 5e-05 | 5e-05 | 5e-05 | 0.0016 | 5e-05 | 5e-05 | 5e-05 | 0.0014 | 5e-05 | 5e-05 | 0.39 | 0.12 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 |
| intergenic | 0.087 | 5e-05 | 0.014 | 5e-05 | 0.00075 | 5e-05 | 0.0008 | 5e-05 | 0.37 | 0.00045 | 5e-05 | 5e-05 | 0.013 | 5e-05 | 0.026 | 5e-05 | 5e-05 | 5e-05 | 0.0004 | 5e-05 | 5e-05 | 5e-05 | 0.06 | 5e-05 |
| genes | 0.08 | 5e-05 | 0.012 | 5e-05 | 0.00035 | 5e-05 | 0.0004 | 5e-05 | 0.36 | 5e-05 | 5e-05 | 5e-05 | 0.018 | 5e-05 | 0.025 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 0.0004 | 5e-05 | 5e-05 | 0.071 | 5e-05 |
| 5' UTR | 5e-05 | 0.029 | 0.0012 | 5e-05 | 5e-05 | 0.0028 | 5e-05 | 0.2 | 5e-05 | 5e-05 | 0.0016 | 0.0056 | 5e-05 | 0.0004 | 5e-05 | 5e-05 | 0.0004 | 0.0028 | 5e-05 | 0.068 | 0.0004 | 5e-05 | 0.0016 | 0.0008 |
| exons | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 0.078 | 5e-05 |
| CDS | 0.0012 | 0.0011 | 5e-05 | 5e-05 | 0.0004 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 |
| introns | 0.11 | 0.0008 | 0.28 | 0.002 | 0.17 | 5e-05 | 0.0027 | 0.28 | 5e-05 | 0.15 | 0.32 | 0.054 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 0.0018 | 5e-05 |
| 3' UTR | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 0.0008 | 5e-05 | 5e-05 | 0.0004 | 5e-05 | 0.15 | 0.5 | 0.27 | 0.23 | 0.046 | 0.37 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 5e-05 | 0.0008 |
| | | | HDRs | | | | | | SVs | | | | | | Indels | | | | | | SNPs | | | |

p-value scale: 0.05 – 0.35

**Supplemental Figure S2.** GAT interval-association test results analyzing the overlap of DLW CB4856 Hawaiian SNPs, indels, and SVs with N2 genome annotations. A) Heatmap showing the fold enrichment of each variant type within gene annotations for each chromosome. B) Heatmap of p-values associated with corresponding fold enrichments shown in panel A calculated by the hypergeometric test.

286 density of SNPs and short indels in the central regions of the autosomes, particularly on

287 chromosomes IV and V (Figure 1 C-D).

288
289 Structural variations and HDRs account for most of the base-pairs affected by sequence

290 divergence between our N2 Bristol and CB4856 Hawaiian lineages. The SVs identified ranged

291 in size from 50bp to 592kb (Figure 2D-E), and HDRs ranged from 50bp to 199kb. Within the

292 SVs detected, we identified 47 non-alignable structures, 2 duplications, 18 inversions, and 2

293 translocations. Non-alignable regions (NOTALs) are highly diverged regions containing many

294 repeats and low-complexity sequences that are inhibitory to whole-genome alignment. From our

295 whole-genome alignments of the DLW N2 Bristol and DLW CB4856 Hawaiian genomes, the

296 non-alignable regions between the two genomes comprise 1.39Mb of sequence, ranged in size

297 from 50-592kb, and comprise <0.5% of coding genes in the Bristol genome. One 156kb

298 translocation was found on the right end of CB4856 Hawaiian chromosome V (V:15,871,614-

299 16,027,614bp), while the other translocation, 38kb, was found to be inverted near a telomere of

300 CB4856 Hawaiian chromosome IV (IV: 176:38,447bp). The largest duplication was found on

301 Hawaiian chromosome III (III: 11,819,363-11,860,261). Together, our analyses provide

302 improved variant site identification in wild isolate genomes and further illuminates previously

303 undetected large structural variations and HDRs.

304

305 **SNPs, indels, and SVs are under-enriched in coding regions**

306 Genes are enriched in the central region of all chromosomes in *C. elegans* (C. elegans

307 Sequencing Consortium 1998), but there are some genes scattered across the chromosome

308 arm-like regions. While much of the sequence variation is enriched in the arm-like regions, we

309 wanted to know whether this variation was affecting coding sequences across the genome.

310 Thus, we tested whether the SNPs, indels, SVs, and HDRs we identified between our N2 Bristol

311 and CB4856 Hawaiian assemblies were enriched in genes versus intergenic space. Based on

312    our remapped annotations (see Methods, LiftOff (Shumate and Salzberg 2021)), approximately

313    61.8% of the DLW N2 Bristol genome is comprised of gene sequences, with exons and introns

314    representing 28.6% and  33.2% of the genome, respectively. Thus, we would expect

315    corresponding proportions of each variant type to overlap within each annotation if variant sites

316    were uniformly distributed across the genome. To determine if SNPs and indels were enriched

317    in genes, we used the Genomic Association Tester (Heger et al. 2013) to compare the observed

318    overlap of our variant sites in each remapped annotation to simulated uniform distributions of

319    SNP and indel intervals. Fold enrichments represent the ratio of observed overlap to simulated

320    overlaps, whereby a fold enrichment of 1.0 means there is no difference between the observed

321    and simulated datasets. The greatest overlap of SNPs and indels in gene regions were

322    observed on the autosomes (SNPs: 55.2-69.1%; indels: 52.1-68.5%), while only 49.1% of SNPs

323    and 42.0% of indels were found in genes on the X chromosome (Figure 1 E-F). Across the

324    genome, SNPs were slightly under-enriched in gene regions with an average fold enrichment of

325    0.96 (hypergeometric test, p-value <0.05). The average fold enrichment of indels in gene

326    regions was lower than observed with SNPs (fold enrichment of 0.90, p-value <0.05), which

327    could be due to selection against indels within coding regions.  For SNPs and indels that did

328    overlap with genes, intron sequences harbored the greatest amount of each variant type (SNPs:

329    fold enrichment 1.14; indels: fold enrichment 1.45; Supplemental Figure 2). In conclusion, SNPs

330    and indels are slightly overrepresented in intergenic regions of the *C. elegans* genome.

331        The distribution of SVs and HDRs across each chromosome resembles the genomic

332    distribution of SNPs and indels (Figure 2A-B).  To determine whether these large variant regions

333    were enriched in intergenic versus coding regions, we compared the enrichment of simulated

334    uniform distributions of SVs to those we identified. On the autosomes, 44.5-68.5% of SVs

335    overlapped with gene regions compared 31.1% on the X chromosome (Figure 2C). Compared

336    to SNPs and indels, structural variations on each chromosome, except chromosome III,

17

**A**



**B**



**C**



**D**



**E**



337

338

**Figure 2. Genomic distribution and size of SVs between the DLW N2 Bristol and DLW CB4856 Hawaiian genomes. (A)** Histograms depicting the distribution of SVs across each chromosome in 100kb bins. Black dashes above each histogram correspond to the genomic locations of SVs that are greater than 20kb in size. **(B)** Chromosome alignment plot depicting syntenic regions between N2 Bristol and CB4856 Hawaiian, structural variants, and highly divergent regions (HDRs). The width of lines showing SVs are proportional to their size. Only rearrangements 1kb or greater in size are shown. **(C)** Stacked bar plots showing the percentage of CB4856 Hawaiian SVs that overlap with intergenic and gene-coding regions of the DLW N2 Bristol genome. **(D)** Bar plots showing the number of each type of SV identified. **(E)** Strip plots showing the log-scaled size distribution of SVs separated by type. For SV types: NOTAL = non-aligned regions, DEL =deletion, INS = insertion, CPG = copy gain in query genome, CPL = copy loss in query genome, TDM = tandem repeat region, INV = inversion, DUP = duplication, TRANS = translocation, and INVTR = inverted translocation. For D and E, different colors only correspond to the different types of SV identified.

displayed significant fold enrichments in intergenic regions (fold enrichments 1.3-1.5, p-values < 0.001; Supplementary Figure S2). Similar to SVs, highly divergent regions overlapped with 38.7-66.6% of genes on the autosomes and 24.3% on the X chromosome. HDRs were significantly enriched in intergenic regions of all chromosomes except chromosome I (fold enrichments 1.15-1.65; p-values < .05). Taken together, our data demonstrate that non-coding regions on the chromosome arm-like regions harbor most of the sequence variation between N2 Bristol and CB4856 Hawaiian lineages.

**Minimal movement of DNA transposons between the N2 Bristol and CB4856 Hawaiian lineages**

Early analyses of the *C. elegans* genome indicated that approximately 12-16% of the genome is comprised of transposable elements (TEs) (*C. elegans* Sequencing Consortium 1998; Bessereau 2006), with *Tc1/mariner* elements as one of the most widely studied DNA transposons that can be active in laboratory strains (Emmons et al. 1983; Liao, Rosenzweig, and Hirsh 1983). While transposable element distributions have been assessed in wild *C. elegans* strains using older reference genomes and Illumina short-read sequencing (Laricchia et al. 2017)*,* the complete TE composition has not yet been reassessed in a *de novo* assembly built from long-read sequencing. Further, new families of eukaryotic Class II transposons, have

373    been discovered (Bao et al. 2009), and it remains unclear if these emerging families of DNA

374    transposable elements comprise a significant proportion of the *C. elegans* genome.

375       To identify and locate known transposable element sequences in our N2 Bristol and

376    CB4856 Hawaiian assembled genomes, we used a transposable element identification pipeline

377    that applies an ensemble of programs to find all known RNA and DNA transposable element

378    families (Riehl et al. 2022). We found that approximately 14.7 and 14.3% of our N2 Bristol and

379    CB4856 Hawaiian assemblies, respectively, are composed of transposable element sequences

380    (Supplemental Table 1). For both genome assemblies, the distribution of TEs was concentrated

381    in the terminal third, arm-like regions of each chromosome (Figure 3A-B). Class II DNA TEs

382    represented 96% of all TEs identified in each genome, and Zator elements are 52% of these

383    Class II DNA TEs present in each genome (Supplemental Table 1, Figure 3C-D). To our

384    knowledge, movement of Zator elements and other recently identified TE families has not yet

385    been analyzed in *C. elegans* laboratory strains.  Further, we also found that N2 Bristol genome

386    has 20.6% more TEs from the *Tc1/mariner* family than the CB4856 Hawaiian genome

387    (Supplemental Table 1).

388

389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406

**Supplemental Table 1. Transposable Elements identified in DLW N2 Bristol genome (this study) vs DLW CB4856 Hawaiian genome (this study)**

| | DLW N2 Bristol | DLW CB4856 Hawaiian |
|---|---|---|
| **Class I Transposable Elements (Retrotransposons)** | 710 (2,688,730 bp) | 776 (2,522,357 bp) |
| **Gypsy** | 557 (2,195,895 bp) | 592 (2,031,038 bp) |
| **Copia** | 134 (472,195 bp) | 161 (465,785 bp) |
| **SINE** | 9 (2,146 bp) | 9 (1,945 bp) |
| **ERV** | 7 (8,280 bp) | 6 (7,569 bp) |
| **LINE** | 3 (10,214 bp) | 8 (16,038 bp) |
| **Class I intrachromosomal transpositions*** | 0 | |
| **Class I interchromosomal transpositions*** | 0 | |
| **Class II Transposable Elements (DNA transposons)** | 17,682 (12,055,357 bp) | 17,310 (11,606,010 bp) |
| **Tc1/Mariner** | 1870 (1,298,386 bp) | 1,550 (1,131,443 bp) |
| **hAT** | 3,999 (3,988,461 bp) | 3,818 (3,725,667 bp) |
| **CMC** | 1,679 (3,138,647 bp) | 2,011 (3,260,455 bp) |
| **Zator** | 9,159 (3,009,341 bp) | 8,980 (2,907,391 bp) |
| **Novosib** | 46 (12,060 bp) | 28 (12,088 bp) |
| **Helitron** | 39 (368,980 bp) | 43 (329,238 bp) |
| **Sola** | 821 (226,645 bp) | 699 (196,797 bp) |
| **MITE** | 69 (12,837 bp) | 181 (42,931 bp) |
| **Class II intrachromosomal transpositions*** | 38 | |
| **Class II interchromosomal transpositions*** | 9 | |

* All TE sequences with predicted transpositions are relative to the DLW N2 Bristol genome

21

421

422

423 **Figure 3. Genomic distributions of transposable elements in the DLW N2 Bristol and DLW**
424 **CB4856 Hawaiian genomes.** Histograms depicting the distributions of transposable elements
425 across the DLW N2 Bristol genome in 100kb bins. **B)** Histograms depicting the distributions of
426 transposable elements across the DLW CB4856 Hawaiian genome in 100kb bins. **C,D)** Stacked
427 bar plot depicting the percent of total DNA transposable elements on DLW N2 Bristol (C) and
428 DLW CB4856  Hawaiian (D) chromosomes accounted for by specific DNA transposon families.
429 For TE families: CMC= CACTA, Mirage and Chapaev families; hAT = hobo and Activator
430 families; Other = MITE, Novosib and Helitron families. **E)** Ideogram depicting the locations of
431 individual DNA transposable elements that moved between the DLW N2 Bristol genome and the
432 DLW CB4856 Hawaiian genome. DLW N2 Bristol chromosomes are represented by the blue
433 boxes on the top, and DLW CB4856 Hawaiian chromosomes by the red boxes on the bottom.
434 Each line represents an individual transposable element sequence, traced from its position on
435 the DLW N2 Bristol genome to its unique position on the DLW CB4856 Hawaiian genome.
436 Transposable elements predicted to have translocated are colored according to transposon
437 class. Arrow heads across the Bristol N2 chromosomes indicate DNA TEs where duplicated
438 copies are found in the Hawaiian CB4856 genome.
439

440          Since the N2 Bristol and CB4856 Hawaiian lineages were geographically isolated for

441 thousands of generations, we sought to utilize our new TE annotation set to identify individual

442 transposition events that occurred over the course of divergence between the two strains. Using

443 whole-genome alignments and the SNPs we previously defined between these two lineages, we

444 identified specific TE sequences with unique polymorphisms that enables individual transposons

445 to be tracked between the N2 Bristol and CB4856 Hawaiian genome assemblies.  Of the 18,392

446 total transposable elements identified in the N2 Bristol genome, 9,377 TEs were uniquely

447 identifiable by sequence polymorphism. Among all N2 Bristol TEs with SNPs, only 1,535

448 elements were detectable in the CB4856 Hawaiian genome. While the vast majority of TEs were

449 found to have not moved within either genome, we did identify 38 Class II DNA TEs that moved

450 intrachromosomally and 9 TEs that moved interchromosomally (Figure 3E). Specifically, we

451 detected 6 Zator elements and one each of *Tc1/mariner*, Sola, and hAT elements at different

452 interchromosomal locations between the two lineages.  In this analysis, we also found several

453 unique copies of Class II DNA transposable elements in the N2 Bristol genome that had

454 duplicated copies in the CB4856 Hawaiian genome (Figure 3E, arrowheads). While we were

455 able to identify transposition events relative to the N2 Bristol genome, we cannot accurately

456 infer the history of each CB4856 Hawaiian copy to determine which resulted from transposition

457    versus duplication. Overall, the landscape of transposable elements remains largely unchanged

458    across the history of divergence between the N2 Bristol and CB4856 Hawaiian lineages.

459

460    **Structural variants predominate the sequence divergence between lab strains**

461    Much of the work exploring *C. elegans* genetic diversity utilizes comparisons of different

462    natural isolates (Koch et al. 2000; Wicks et al. 2001; Thompson et al. 2015; Andersen et al.

463    2012). Work on germline mutation rates in *C. elegans*, however, suggest that considerable

464    genetic variation may have been incurred during the laboratory setting (Denver et al. 2009).

465    Given the rate of mutation accumulation in the germline ($2.7 \times 10^{-9}$ mutations per site per

466    generation (Denver et al. 2009)) and a generation time of approximately three days, each N2

467    lineage alone may have accumulated up to ~1,500 single nucleotide mutations since the 1970s,

468    and nearly 790 potential mutations since the first genome was published in 1998 (C. elegans

469    Sequencing Consortium 1998). Notably, this predicted variation does not include the

470    accumulation of indels and structural variations. Thus, the N2 Bristol and CB4856 Hawaiian

471    genomes present in each lab strain likely carries considerable genomic variation relative to

472    other labs isolates. Previous studies using earlier genome assemblies identified many

473    segmental duplications between lab lineages of wild type strains (Vergara et al. 2009). This

474    variation may underpin phenotypic variation as well as previous work that has shown the

475    lifespans of laboratory N2 Bristol isolates varies between 12-17 days (Gems and Riddle 2000).

476    Taken together, accumulating evidence suggests that inter-lab genetic variation in wild type

477    backgrounds may contribute to differences in experimental outcomes. High-quality lab-specific

478    reference genomes may be an important tool to understand how genetics influences the

479    phenotypes and processes studied by different laboratory groups.

480    To further evaluate the quality and differences of our genome assemblies, we aligned

481    our N2 Bristol genome to the VC2010 Bristol (Yoshimura et al. 2019), as well as aligned our

482    CB4856 Hawaiian genome to the Kim CB4856 Hawaiian genome (Kim et al. 2019). We

483    expected that examining whole-genome alignments to previously validated long-read

484    assemblies would reveal a striking degree of similarity. Comparing our N2 Bristol genome to

485    VC2010, 99.9% of bases were alignable and 99.8% of bases were in syntenic alignments

486    (Table 2). Analysis of our CB4856 Hawaiian genome versus the Kim CB4856 Hawaiian genome

487    showed that 96.1% of bases were alignable, with 92.3% of bases in syntenic alignments (Table

488    3). This high degree of similarity within alignments gives us increased confidence in the quality

489    of our own genome assemblies.

490        To assess how much genetic variation may exist between lab lineages of the most

491    utilized wild-type strain, we first compared our N2 Bristol genome to VC2010 Bristol. We

492    identified 1,162 homozygous SNPs and 1,528 homozygous indels. (Figure 4A-B, Table 2). In

493    total, over 2.07Mb were affected by SNPs, indels and SVs, with 99.7% of this sequence

494    divergence due to structural variations (Figure 4C, Table 2). While highly divergent regions have

495    been observed between wild populations of *C. elegans* (D. Lee et al. 2021)*,* we were also able

496    to identify over 404kb of sequence as HDRs between these two laboratory Bristol lineages

497    (Table 2). These HDRs identified between laboratory strains represent regions with multiple

498    gaps between both genomes within a pairwise alignment in regions of synteny (Goel et al.

499    2019). In addition, we identified two inverted duplications (5.4kb and 12.9kb on chromosomes III

500    and V, respectively) and 39 simple inversions. Four of these inversions are over 29kb in size

501    and account for 11.6% of all structural variation between our N2 Bristol and the VC2010 Bristol

502    genomes. SVs of this nature can be particularly disruptive to genome organization by impairing

503    interactions between regulatory sequences or disrupting gene expression through loss of coding

504    regions (Stranger et al. 2007; Hurles, Dermitzakis, and Tyler-Smith 2008).

505        Examination of our CB4856 Hawaiian lineage compared to the Kim *et al.*, 2019 CB4856

506    Hawaiian assembly (Kim et al. 2019) revealed a greater amount of sequence divergence than

507    comparisons between laboratory lineages of N2 Bristol. We identified 541 homozygous SNPs

508    and 1,298 homozygous indels by aligning our CB4856 Hawaiian short reads to the Kim CB4856

509    Hawaiian genome (Supplementary Figure S3, Table 3, see Methods).

510

511

**Table 2. Comparisons between the DLW N2 Bristol genome (this study) and VC2010 Bristol genome**

| | Chromosomes | | | | | | Total |
|---|---|---|---|---|---|---|---|
| | **I** | **II** | **III** | **IV** | **V** | **X** | |
| **DLW N2 Bristol Chromosome Length (this study)** | 15,114,068 | 15,311,845 | 13,819,453 | 17,493,838 | 20,953,657 | 17,739,129 | 100,431,990 bp |
| **VC2010 Bristol Chromosome Length (Yoshimura et al., 2019)** | 15,331,301 | 15,525,148 | 14,108,536 | 17,759,200 | 21,243,235 | 18,110,855 | 102,078,275 bp |
| **DLW N2 Bristol Bases Aligned** | 15,108,942 | 15,310,622 | 13,819,294 | 17,492,076 | 20,852,291 | 17,738,432 | 100,321,657 bp (99.89%) |
| **% Syntenic Aligned Bases** | 98.83 | 99.83 | 99.47 | 99.12 | 99.05 | 99.50 | 99.28 |
| **SNPs*** | 169 | 124 | 164 | 209 | 280 | 216 | 1,162 |
| **Indels*** | 150 | 261 | 210 | 262 | 378 | 267 | 1,528 (3465 bp) |
| **SVs** | 113 | 134 | 83 | 228 | 175 | 164 | 897 (2,010,282 bp) |
| **HDRs** | 8 | 14 | 10 | 21 | 24 | 11 | 88 (406,737 bp) |

512
513    * All variants listed are only those for which the VC2010 Bristol genome was homozygous
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528

26

529

530

531

**Table 3. Comparisons between the DLW CB4856 Hawaiian genome (this study) and Kim CB4856 Hawaiian genome**

| | Chromosomes | | | | | | Total |
|---|---|---|---|---|---|---|---|
| | **I** | **II** | **III** | **IV** | **V** | **X** | |
| **DLW CB4856 Hawaiian Chromosome Length (this study)** | 15,045,644 | 15,257,363 | 13,206,755 | 17,183,882 | 20,547,529 | 17,584,915 | 98,826,088 |
| **Kim CB4856 Hawaiian Chromosome Length (Kim et al., 2019)** | 15,528,896 | 15,813,191 | 14,110,336 | 17,985,219 | 21,389,866 | 18,073,349 | 102,900,857 |
| **DLW CB4856 Hawaiian Bases Aligned** | 14,620,886 | 14,680,704 | 12,451,582 | 16,482,840 | 19,427,050 | 17,371,269 | 95,034,331 (96.16%) |
| **% Syntenic Aligned Bases** | 94.91 | 92.38 | 93.01 | 89.28 | 89.24 | 96.10 | 92.32 |
| **SNPs*** | 60 | 52 | 71 | 108 | 135 | 115 | 541 |
| **Indels*** | 240 | 190 | 175 | 242 | 238 | 213 | 1,298 (2157 bp) |
| **SVs** | 148 | 274 | 194 | 626 | 660 | 168 | 2,070 (6,923,335 bp) |
| **HDRs** | 19 | 70 | 25 | 100 | 144 | 12 | 370 (3,327,407 bp) |

532

533    * All variants listed are only those for which the Kim CB4856 Hawaiian genome was
534    homozygous

535

536    Notably, analysis of our whole-genome alignments identified over 9.5Mb of structural variation

537    and HDRs between these two genomes. More than 66% of this structural variation, however, is

538    due to unique, non-alignable regions. These non-alignable regions are highly divergent with

539    many gaps in pairwise alignments that contain many repeats and low-complexity sequences.

540    Further, over 3.3Mb in each Hawaiian genome falls within highly divergent regions. Taken

541    together, we detected much more variation than anticipated between laboratory wild-type

542    genomes. In the laboratory isolates of N2 Bristol and CB4856 Hawaiian, SVs affected the

543    greatest number of base pairs, with a large portion of this variation due to large non-alignable



544

**Figure 4. Genomic variation between the DLW N2 Bristol genome and the VC2010 Bristol genome. (A-B)** Histograms depicting the distribution of DLW N2 Bristol SNPs and indels across each VC2010 Bristol chromosome in 500kb bins. **(C)** Scatterplots showing the genomic position of SVs with the log-scaled size of each SV on the y-axis. **(D)** The proportions of DLW N2 Bristol

549 SNPs, indels, and SVs that overlap with intergenic versus gene-coding regions of the VC2010
550 Bristol genome.

551

552



553

554    **Supplementary Figure S3**. Genomic variation between the DLW CB4856 Hawaiian genome
555    and the Kim CB4856 Hawaiian genome. (A-B) Histograms depicting the distribution of SNPs
556    and indels across each Kim CB4856 Hawaiian chromosome in 500kb bins. (C) Scatterplots
557    showing the genomic position of SVs with the log-scaled size of each SV on the y-axis. (D) The
558    proportions of DLW CB4856 Hawaiian SNPs, indels, and SVs that overlap with intergenic
559    versus gene-coding regions of the Kim CB4856 Hawaiian genome.
560

561    regions, duplications, and inversions. Thus, the genomes of wild type strains present in some

562    labs are not only unlike the most widely used reference genome in the *C. elegans* research

563    community, but there are likely many large inter-lab genomic variations that might underlie some

564    of the phenotypic differences observed in laboratory strains.

565

566    **Intergenic enrichment of variant sites between lab lineages of N2 Bristol and CB4856**

567    **Hawaiian**

568            To determine whether specific genomic regions of lab strains are susceptible to

569    sequence variation, we repeated our analysis of the genomic distributions of each variant class

570    and assessed the enrichment of variant sites in gene annotations. We used LiftOff to remap the

571    pre-existing Bristol gene annotations onto both the VC2010 Bristol and Kim CB4856 Hawaiian

572    assemblies with similar success. To address whether the sequence divergence is nonrandomly

573    enriched in each region of interest, we again used GAT to simulate random SNP, indel, SV, and

574    HDR intervals 20,000 times and compared the simulated overlap to what we observed in

575    between genomes. Notably, the stereotypical "arms"-vs-"center" genomic distribution of variants

576    seen when comparing Bristol and Hawaiian genomes is not true for all chromosomes when

577    comparing our Bristol genome to VC2010 (Figure 4A-C), with some chromosomes displaying a

578    concentration of variation in the central region. SNPs were 1.2-1.5 fold enriched in the arm-like

579    regions of chromosomes I, II, III, and V (p-values <0.01). Indels, however, were concentrated in

580    the arm-like regions of all chromosomes with fold enrichments ranging from 1.2-1.5 (p-values <

581    0.05). While SVs were 1.2-2.2 fold enriched in the arm-like regions of each chromosome, this

582    enrichment was only significantly higher than expected by null distributions on chromosomes I,

583    IV, and the X chromosome (Supplemental Figure S4). Further, HDRs between Bristol lineages

584    were 1.6 fold enriched on the arm-like regions of chromosome II and the X chromosome (p-

585    values < 0.05), while displaying significant 1.8-2.1 fold enrichments in the center regions of

586    chromosomes I, IV, and V (p-values < 0.01). Finally, we also wanted to determine whether the

587    variant sites we detected between lab strains impacted gene coding regions. Between the two

588    Bristol lineages, SNPs, indels, SVs, and HDRs were all under-enriched in gene coding regions

589    and displayed significant enrichments in intergenic regions on most chromosomes

590    (Supplemental Figure S4). Thus, variation between laboratory Bristol lineages is largely

591    concentrated in non-coding regions of each chromosome.

592         We next examined the genomic distribution and enrichment of variant sites in gene

593    annotations of the two Hawaiian genomes to see if the patterns of enrichment were similar to

594    the Bristol genomes. After examining the enrichment of all variant types in both the arm-like

595    regions and "centers" of each chromosome, it was clear that most chromosomes were enriched

596    for variant sites in the arm-like regions and intergenic sequences, with a few exceptions as

597    follows. (Supplemental Figure S5). SNPs were only enriched 0.84 fold in the arm-like regions of

598    the X chromosome, and indels were enriched 0.81 fold in the arm-like regions of chromosome

599    IV (p-values < 0.05). SVs were 0.64 fold enriched in the arm-like regions of chromosome II, and

600    HDRs were 0.12 fold enriched in the arm-like regions of chromosome I (p-values < 0.001). We

601    then examined the enrichment of all variants in intergenic versus gene sequences between the

602    two CB4856 genomes. SNPs and indels showed significant 1.7-2.8 fold enrichments in

603    intergenic regions on all chromosomes (p-values < 0.001). SVs displayed a significant 1.2-1.7

604    fold enrichment in the intergenic regions of all chromosomes. HDRs were 1.2-1.7 fold enriched

605    in the intergenic regions of chromosomes II, III, IV, V and the X chromosome (all p-values <

606    0.05).  In conclusion, analysis of the genetic variation between respective lab lineages of N2

607    Bristol and CB4856 Hawaiian revealed a striking amount of variation often present in intergenic

31

608 sequences, with some weak enrichments in the arm-like regions versus the central regions of

609 chromosomes.

610



611

**Supplemental Figure S4.** GAT interval-association test results analyzing the overlap of DLW
N2 Bristol SNPs, indels, and SVs with remapped VC2010 Bristol genome annotations. A)
Heatmap showing the fold enrichment of each variant type within gene annotations for each

615    chromosome. B) Heatmap of p-values associated with corresponding fold enrichments shown in
616    panel A calculated by the hypergeometric test.



617

618    **Supplemental Figure S5.** GAT interval-association test results analyzing the overlap of DLW
619    CB4856 Hawaiian SNPs, indels, and SVs with remapped Kim CB4856 Hawaiian genome
620    annotations. A) Heatmap showing the fold enrichment of each variant type within gene
621    annotations for each chromosome. B) Heatmap of p-values associated with corresponding fold
622    enrichments shown in panel A calculated by the hypergeometric test.
623

624

**Discussion**

626        Detection and characterization of sequence variation between individuals or across

627  species is fundamental to our functional understanding of genomic elements and consequences

628  of variation. Since the first draft of the *C. elegans* genome was released in 1998, the highly

629  divergent strains N2 Bristol and CB4856 Hawaiian have been used extensively for comparative

630  genomics studies(C. elegans Sequencing Consortium 1998; Koch et al. 2000; Wicks et al. 2001;

631  Maydan et al. 2010; Andersen et al. 2012; D. Lee et al. 2021). The combined usage of short and

632  long read sequencing to assemble genomes and to compare them has both increased the

633  quality of our reference genomes as well as enhanced the genome-wide detection of sequence

634  variants, new genes, and new genomic regions (Yoshimura et al. 2019; C. Kim et al. 2019; B. Y.

635  Kim et al. 2021; Sarsani et al. 2019). In this study, we generate *de novo* assemblies for the N2

636  Bristol and CB4856 Hawaiian *C. elegans* isolates from our lab lineage using short-read and

637  long-read sequencing. Our examination of the inter-lab genetic drift among wild-type strains

638  suggests genomic analyses can be improved by resequencing the genomes of labs' wild-type

639  strains or utilizing strains with recently published, accurate genome assemblies. This also

640  presents a strong argument for labs utilizing *C. elegans* in their research to frequently return to

641  cryogenically preserved stocks of their wild type strains. These genomes will serve as additional

642  tools for future comparative genomics studies, especially in the functional characterization of

643  structural variations identified through whole-genome alignments.

644

**Genome assembly and genomic divergence in laboratory isolates**

646        Earlier studies uncovering phenotypic and genetic variations between lab wild-type

647  strains indicated that there are likely many underlying large-scale genomic differences (Denver

648  et al. 2009; Vergara et al. 2009; Gems and Riddle 2000). Here we identify numerous SNPs,

649  indels, SVs, and HDRs between different lab lineages of each wild isolate. The total amount of

34

650    genomic variation is at levels higher than predicted by earlier mutation accumulation studies.

651    Much of this variation, however, is due to SVs and HDRs, which have only recently become a

652    detailed subject of study (Thompson et al. 2015; Kim et al. 2019; Lee et al. 2021). Our genome

653    assemblies of the Bristol and Hawaiian strains corroborate prior results indicating that genomic

654    variation is enriched in the distal arm-like regions of chromosomes between these natural

655    isolates. Evolutionary genomic analysis has shown that recombination in the arm-like regions of

656    each chromosome and balancing selection likely have shaped this landscape of sequence

657    divergence across the 30,000-50,000 generations these strains have been geographically

658    isolated (Thomas et al. 2015; Kern and Hahn 2018). In contrast, we find that the distribution of

659    variant sites across the arm-like regions versus center domains of each chromosome between

660    lab lineages is not as strong or consistent as seen when comparing N2 Bristol to CB4856

661    Hawaiian genomes. This result could indicate that in relatively short timescales (~3,000-5,800

662    generations), selection for the accumulation of mutations in the arm-like regions, particularly in

663    noncoding regions, is not sufficient to consistently eliminate sequence divergence away from

664    the gene-dense chromosome centers. Further, we found that SNPs, indels, and structural

665    variations were highly enriched in intergenic regions when comparing the genomes of laboratory

666    strains. Although many of the sequence variants we identified are not directly disrupting coding

667    sequences, it remains possible that genetic drift in these regions are altering the function of

668    intergenic regulatory sequences such as promoters and enhancers. Thus, the accumulation of

669    disruptive genomic changes within regulatory regions in the gene-dense centers of

670    chromosomes may underpin many of the phenotypic differences observed in laboratory wild-

671    type strains, such as variance in lifespan (Gems and Riddle 2000).

672

673    **Highly variable arm-like domains on *C. elegans* chromosomes**

674        The arm-like regions of *C. elegans* chromosomes exhibit a striking degree of variation

675    that is highly correlated with large domains of increased recombination, which is a pattern

676    observed in many species (Andersen et al. 2012; D. Lee et al. 2021; Kern and Hahn 2018;

677    Rockman and Kruglyak 2009). In *C. elegans,* these divergent autosomal arm-like domains

678    coincide with a disproportionate fraction of newer, rapidly evolving genes as compared to the

679    center regions of each chromosome, which house highly conserved essential genes (C. elegans

680    Sequencing Consortium 1998; Kamath et al. 2003). The development of new tools to detect

681    larger structural variations through alignment of assemblies or long sequencing reads has

682    revealed many SVs on the chromosomal arm-like domains (Mahmoud et al. 2019; C. Kim et al.

683    2019). The fact that SVs are enriched in the arm-like regions, which also display elevated levels

684    of recombination, is notable given the fact that large structural variants such as inversion are

685    typically inhibitory to recombination (Miller, Cook, and Hawley 2019). The arm-like regions of *C.*

686    *elegans* chromosomes are enriched for many repetitive elements, including transposable

687    elements, tandem repeats, and low complexity repeat sequences (C. elegans Sequencing

688    Consortium 1998; Surzycki and Belknap 2000). The presence of many SVs in the arm-like

689    regions could be due to errors in double-strand DNA break repair and heterologous

690    recombination in regions adjacent to highly repetitive sequences, thereby causing chromosomal

691    rearrangements. Similar rearrangement events are known to contribute to many human

692    genomic disorders like Prader-Willy Syndrome or Charcot-Marie-Tooth disease (Carvalho and

693    Lupski 2016; Stankiewicz and Lupski 2010). Future investigations assessing the occurrence of

694    SVs adjacent to highly repetitive regions and sites of homologous recombination will be

695    invaluable in understanding how differences in genomic organization arise between divergent

696    lineages of *C. elegans*.

697        With regard to genomic rearrangements and their impact on genome function, renewed

698    attention must be given to the contribution of transposable elements and their mobility within

699    and between chromosomes. While Sola and Zator elements are relatively recent in their

700    discovery within *C. elegans* and other eukaryotic genomes (Bao et al. 2009; Riehl et al. 2022),

701    our data suggests there may be many active TE copies in these families, particularly Zator

702    elements. Historically, much attention has been given to the impact of *Tc1/Mariner* transposition

703    on genomic architecture, but the contribution of Zator elements to changes in genome structure

704    and gene regulation merits further future investigation. Our analysis of TE mobility only

705    examines two endpoints across the long period of divergence between the Bristol and Hawaiian

706    lineages. It remains unclear, however, whether many of these newly characterized TEs remain

707    active and whether they contribute to the growing catalog phenotypic differences displayed

708    between laboratory lineages of Bristol and Hawaiian *C. elegans*.

709
710    Finally, the generation of multiple independent long read *de novo* genome assemblies

711    for both N2 Bristol and CB4856 Hawaiian isolates provides a powerful toolkit for comparative

712    genomics and evolution studies. Many prior studies assessing the *C. elegans* recombination

713    landscape have relied on mapping recombination in worms heterozygous for Bristol and

714    Hawaiian chromosomes. The high sequence divergence and large structural variations between

715    Bristol and Hawaiian which we describe, however, may have positional impact on the

716    distributions of crossover sites. Our identification of variants in Bristol strains enables

717    polymorphism mapping by crossing different lab-lineages of N2 Bristol, avoiding the potential

718    confounding effects of crosses with other wild isolates. Additionally, further identification and

719    functional characterization of polymorphic sites and structural variations present between lab

720    lineages of N2 Bristol and CB4856 Hawaiian could provide new insights into how pronounced

721    phenotypic differences in the lifespan, feeding behavior, and reproductive fitness arise in

722    modern lab-derived strains (Gems and Riddle 2000; Zhao et al. 2018). To summarize, we

723    demonstrate the importance of using long and short-read sequencing to generate modern

724    reference genome assemblies and maximally detect sequence variation, while highlighting the

725    potential genomic underpinnings of phenotypic variations in laboratory lineages of *C. elegans*.

726

727

**Methods**

***C. elegans* culture and sucrose floatation**

The N2 Bristol and CB4856 Hawaiian strains of *C. elegans* were grown at 20°C on standard

NGM agar plates seeded with the OP50 strain of *E. coli* as a food source. To minimize bacterial

contamination in downstream gDNA sample preps, we performed sucrose floatation on pooled

populations of each isolate. Worms were washed from plates with 8mL cold M9 buffer and

transferred to 15mL glass centrifuge tubes using a glass Pasteur pipette. Collected worms were

centrifuged at 3000rpm at 4°C and washed in 4mL of fresh M9 twice. To separate worms from

bacteria and other debris, 4mL of 60% sucrose solution was added to 4mL of M9 buffer and

worms and vortexed briefly. The mixture was then spun at 5000 rpm at 4°C for 5 minutes. Using

a glass pipette, the floating layer of worms were transferred to a new glass centrifuge tube on

ice and brought up to 4mL in fresh M9. Worms were then incubated at room temp for 30

minutes and gently vortexed every 5 minutes. Worms were washed three times in equal volume

of fresh M9 were performed before storing collected worms in M9 at 20°C before genomic DNA

(gDNA) extraction.

**Long-read and short-read sequencing**

Genomic DNA was extracted from worms using the Qiagen DNeasy Blood and Tissue Kit.

Sequencing was performed on pooled populations of N2 and CB4856 after reducing bacterial

contamination by sucrose float for each strain. For PacBio long-read sequencing, library

preparation was performed on pooled populations of worms for each isolate by the University of

Oregon's Genomics and Cell Characterization Core Facility and sequenced on the Sequel II

system. For Illumina short-read sequencing, library preparation was performed on pooled

populations of worms for each isolate by the University of Oregon's Genomics and Cell

Characterization Core Facility. The short-read libraries were then sequenced on an Illumina

HiSeq4000 (2 x 150bp).

754

755     **Long-read genome assembly and short-read refinement**

756     PacBio long-reads were aligned to the E. coli genome using BWA (Li and Durbin 2009) (version

757     0.7.17), and reads that aligned to the bacterial genome were removed. De novo genome

758     assembly was performed for N2 Bristol and CB4856 Hawaiian using Canu (Koren et al. 2017)

759     (version 1.7). To refine the long-read assemblies, short-reads from each isolate were aligned to

760     their respective long-read assembly using BWA-MEM (version 0.7.17). Aligned reads in SAM

761     format were sorted and converted to BAM format using SAMtools(Li et al. 2009). Using Picard

762     (https://broadinstitute.github.io/picard/), read groups were added via

763     AddOrReplaceReadGroups, and duplicate reads were filtered using MarkDuplicates. Some

764     bases may have been inaccurately called due to lower sequencing coverage, larger error rate in

765     PacBio sequencing, or predominating alleles present in the population of each isolate that could

766     be revealed by greater sequencing depth afforded by Illumina sequencing. GATK's

767     HaplotypeCaller (McKenna et al. 2010) and Freebayes (Garrison and Marth 2012) were utilized

768     to generate VCF files representing potentially inaccurate sites in each initial assembly.

769     Coverage thresholds were manually determined using IGV for each assembly. Sites were

770     filtered according to manual values using VCFtools (Danecek et al. 2011; Danecek and

771     McCarthy 2017)). Error correction was performed on single-nucleotide alleles using BCFtools

772     *consensus* (Danecek and McCarthy 2017) and alternate indel alleles. After filtering potential

773     sites by sequencing depth thresholds determined for each chromosome, this left 4237 and

774     36145 corrections for the N2 Bristol and CB4856 Hawaiian genomes, respectively. Of these

775     sites, less than .7% were unable to be resolved, and all of these were short indels comprising

776     less than .001% of each genome.

777

778     **Assessing genome assembly completeness**

779  To further assess the quality and completeness of our N2 Bristol and CB4856 Hawaiian

780  assemblies, we used BUSCO (Simão et al. 2015; Manni et al. 2021). BUSCO was run in a

781  Docker container (https://busco.ezlab.org/busco_userguide.html) in genome mode. For each

782  assembly, the quality and presence of expected orthologous genes was checked against the

783  nematoda and metazoan lineage databases.

784

785  **SNP and indel Calling in N2 and CB4856 assemblies**

786  Illumina short reads from the DLW N2 Bristol and DLW CB4856 Hawaiian genome were

787  trimmed using Trimmomatic (Bolger, Lohse, and Usadel 2014) to remove adapter and barcode

788  sequences. The trimmed CB4856 reads were then aligned to the DLW N2 Bristol reference

789  genome using BWA-MEM so that SNPs and indels present between N2 Bristol and CB4856

790  Hawaiian could be identified. All resulting variant positions comparing our N2 Bristol and

791  CB4856 Hawaiian genomes are in relation to the N2 Bristol assembly. Aligned reads in SAM

792  format were then sorted using SAMtools (Li et al. 2009) and converted to BAM files. Using

793  Picard read groups were added via AddOrReplaceReadGroups, and duplicate reads were

794  filtered using MarkDuplicates as described above. BAM files with filtered duplicate reads were

795  used to call variants using a combination of GATK HaplotypeCaller, Freebayes, and BCFtools.

796  The three resulting VCF files containing SNPs and indels were then concatenated, further

797  filtered for duplicate sites and low-quality variants, and sorted using BCFtools. SNPs with QUAL

798  scores of 30 or greater, a minimum of 10 variant reads, and a minimum of 30 total, high-quality

799  reads were retained. To draw comparisons between other N2 assemblies, the most recent gene

800  annotations were downloaded from WormBase (*C. elegans* VC2010, PRJEB28388). For

801  comparisons between CB4856 Hawaiian genomes, assemblies were downloaded WormBase

802  (*C. elegans* CB4856, PRJNA275000) and the NCBI BioProject Database under accession

803  number PRJNA523481. To call variants between our N2 Bristol and CB4856 Hawaiian

804   assemblies and those generated by other labs, short reads were aligned to the respective

805   genomes and the SNP and indel calling pipeline was repeated as described above.

806

807   **Calling Structural Variants using whole-genome alignments**

808   All assembly-to-assembly alignments were performed using Minimap2 (Li 2018). SyRI (Goel et

809   al. 2019) was then used to parse the resulting SAM files and call structural variants and highly

810   divergent regions (Structural rearrangements were plotted with the aid of Plotsr within the SyRI

811   package. "NOTAL" or non-alignable regions in each genome were retained as SVs. To acquire

812   NOTAL regions in each query genome, the Minimap2 alignment was repeated with the original

813   reference and query genomes swapped. The sizes of HDRs depicted in Tables 1-3 are sizes

814   relative to the reference genome in each comparison (*i.e*. N2 Bristol in Table 1). When

815   comparing our CB4856 Hawaiian genome to the Kim CB4856 genome, 89% of the size

816   difference in assemblies can be accounted for in the net sequence gained from Kim HDRs and

817   unique NOTAL structures. NOTAL structures and gap-adjacent sequences in the Kim CB4856

818   genome are 1.5 and 1.6-fold enriched for low complexity and repeat sequences, respectively.

819   These regions and sequence features are challenging for genome assembly and likely explain

820   megabase-scale differences in genome assembly sizes.

821

822   **Converting gene annotations between assemblies**

823   We converted gene annotations from the N2 reference assembly (cel235) to our N2 Bristol and

824   CB4856 Hawaiian assemblies, as well as the VC2010 Bristol and Kim CB4856 Hawaiian

825   assemblies. The gene annotations for the WBcel235 genome assembly were downloaded in

826   GFF3 format from Ensembl (http://ftp.ensembl.org/pub/release-

827   105/gff3/caenorhabditis_elegans/). Unlike previously established tools that require pre-

828   generated chain files(James et al. 2003), Liftoff (Shumate and Salzberg 2021) can accurately

829   remap gene annotations onto newly generated assemblies using Minimap2 assembly-to-

41

830    assembly alignments. Rather than aligning whole genomes, Liftoff aligns only regions listed in

831    the annotation files so that genes may be remapped even if there are large structural variations

832    between two genomes. The Liftoff program was then used to remap annotations between the

833    WBcel235 assembly onto each new genome assembly for N2 Bristol and CB4856 Hawaiian).

834

835    **Testing the association between variant sites and gene annotations**

836    For each chromosome, to determine whether SNPs or indels were enriched within gene

837    annotations, fold enrichment analyses were performed using the genomic association tester

838    (GAT) (Heger et al. 2013) tool (https://github.com/AndreasHeger/gat.git). The observed

839    enrichment of each variant type in gene annotations was compared to overlaps in simulated

840    distributions SNPs or indels. Simulated distributions were created using 20,000 iterations

841    whereby each variant type was randomly and uniformly distributed across each chromosome.

842    SNPs and indel distributions were compared against intergenic, gene, intron, exon, and UTR

843    annotations. Comparing the observed enrichment to the simulated distributions, statistical

844    significance was assigned to the observed fold enrichment with p-values calculated from a

845    hypergeometric test calculated within GAT.  Per-chromosome BED files for SNP intervals were

846    created from their original VCF using AWK. Per-chromosome BED files for indel intervals were

847    calculated using a custom script. The GFF3 formatted annotations generated via liftoff were

848    then broken down by chromosome, gene, exon, and UTR regions. Because intron regions were

849    not explicitly written into each GFF3 file, they were calculated using BEDtools (Quinlan and Hall

850    2010). First, a joint BED file containing the UTR and exon regions were made using awk and

851    sorted first by chromosome then by position. Using BEDtools these intervals were combined,

852    and intronic regions were calculated by finding regions in gene intervals not covered by either

853    UTR or exons. Intergenic spaces on each chromosome were calculated with the gene BED files

854    and chromosome sizes as inputs. GAT was then run for each chromosome in each assembly.

855

**Transposable Element Identification and Tracking**

857    The TransposonUltimate pipeline (Riehl et al. 2022) was run for both our N2 Bristol and CB4856

858    Hawaiian genome assemblies. MUST and SINE finder were run independently and integrated

859    into the filtering steps of the pipeline manually. Additionally, we added LTR retriever to the TE

860    identification ensemble to supplement LTR harvest and LTR finder. TE sequences that

861    overlapped with SNPs were identified using BEDtools. CB4856 Hawaiian SNPs were applied to

862    corresponding N2 Bristol TE sequences, and these sequences were cross-referenced with the

863    original TransposonUltimate output for CB4856 Hawaiian for matches. Unique polymorphic TE

864    sequences found in both genomes were then assessed for translocation events by examining

865    genomic start coordinates in each genome. Utilizing whole-genome alignments for each

866    chromosome, TEs were predicted to have moved if starting coordinates for each TE pair were

867    did not correspond to relative changes in coordinates due to alignment.

868

869

870

**Data Availability Statement**

872    The PacBio long-read and the Illumina short-read data generated in this study have been

873    submitted to the NCBI BioProject database (https://www.ncbi.nlm.nih.gov/bioproject/) under

874    accession number PRJNA907379. All custom scripts are available upon request.  Strains are

875    available upon request.

876

**Acknowledgements**

878    We thank N. Kurhanewicz, C. Cahoon, E. Toraason, and J. Conery in the Libuda Lab for

879    thoughtful discussion and comments on the manuscript. We are grateful to the University of

880    Oregon's Genomics and Cell Characterization Core Facility for sample prep and sequencing of

881    our N2 Bristol and CB4856 Hawaiian strains. Over the course of many years beyond the time of

43

882    this study, thanks for inspiration from Walter Rogers and invaluable guidance from Dr. Glenn C.

883    Rowe given to Z.D.B.

884

885    **Funding**

886    This work was supported by the National Institutes of Health T32HD007348 to Z.D.B and

887    National Institutes of Health R35GM128890 and University of Oregon start-up funds to D.E.L.

888

889    **Conflict of Interest**

890    The authors declare no conflicts of interest.

891

892    **Author Contributions**

893    Z.D.B. polished primary genome assemblies, assessed genome quality, annotated the

894    genomes, and analyzed genomic variation and structures. A.F.S.N. helped conceive this study

895    and developed protocols for DNA purification, short-read Illumina sequencing, and variant

896    calling. D.D. developed protocols for long-read sequencing, devised strategies for genome

897    assembly, assembled the contigs and primary assembly, and filled gaps. C.A. identified

898    transposable elements and tracked copies containing SNPs between genome assemblies.

899    K.J.H. helped conceive this study, develop protocols for DNA purification, and purified the DNA

900    for the long-read PacBio sequencing. D.E.L. helped conceive this study, led discussions for the

901    comparison of different genomes, and coordinated work. All authors contributed to the

902    manuscript, with most of the writing by Z.D.B., A.S.F.N., and D.E.L.

903

904

## References

906

Andersen, Erik C, Justin P Gerke, Joshua A Shapiro, Jonathan R Crissman, Rajarshi Ghosh, Joshua S Bloom, Marie-Anne Félix, and Leonid Kruglyak. 2012. "Chromosome-Scale Selective Sweeps Shape Caenorhabditis Elegans Genomic Diversity." *Nature Genetics* 44 (3): 285–90. https://doi.org/10.1038/ng.1050.

Bao, Weidong, Matthew G. Jurka, Vladimir V. Kapitonov, and Jerzy Jurka. 2009. "New Superfamilies of Eukaryotic DNA Transposons and Their Internal Divisions." *Molecular Biology and Evolution* 26 (5): 983–93. https://doi.org/10.1093/molbev/msp013.

Bessereau, Jean-Louis. 2006. "Transposons in C. Elegans." *WormBook : The Online Review of C. Elegans Biology*, January, 1–13. https://doi.org/10.1895/wormbook.1.70.1.

Bolger, Anthony M, Marc Lohse, and Bjoern Usadel. 2014. "Trimmomatic: A Flexible Trimmer for Illumina Sequence Data." *Bioinformatics (Oxford, England)* 30 (15): 2114–20. https://doi.org/10.1093/bioinformatics/btu170.

C. elegans Sequencing Consortium. 1998. "Genome Sequence of the Nematode C. Elegans: A Platform for Investigating Biology." *Science* 282 (5396): 2012–18. https://doi.org/10.1126/science.282.5396.2012.

Carvalho, Claudia M B, and James R Lupski. 2016. "Mechanisms Underlying Structural Variant Formation in Genomic Disorders." *Nature Reviews Genetics*. Nature Publishing Group. https://doi.org/10.1038/nrg.2015.25.

Chalopin, Domitille, Magali Naville, Floriane Plard, Delphine Galiana, and Jean-Nicolas Volff. 2015. "Comparative Analysis of Transposable Elements Highlights Mobilome Diversity and Evolution in Vertebrates." *Genome Biology and Evolution* 7 (2): 567–80. https://doi.org/10.1093/gbe/evv005.

Crombie, Tim A, Stefan Zdraljevic, Daniel E Cook, Robyn E Tanny, Shannon C Brady, Ye Wang, Kathryn S Evans, et al. 2019. "Deep Sampling of Hawaiian Caenorhabditis Elegans Reveals High Genetic Diversity and Admixture with Global Populations." Edited by Graham Coop, Diethard Tautz, and Asher Cutter. *eLife* 8 (December): e50465. https://doi.org/10.7554/eLife.50465.

Danecek, Petr, Adam Auton, Goncalo Abecasis, Cornelis A Albers, Eric Banks, Mark A DePristo, Robert E Handsaker, et al. 2011. "The Variant Call Format and VCFtools." *Bioinformatics (Oxford, England)* 27 (15): 2156–58. https://doi.org/10.1093/bioinformatics/btr330.

Danecek, Petr, and Shane A McCarthy. 2017. "BCFtools/Csq: Haplotype-Aware Variant Consequences." *Bioinformatics (Oxford, England)* 33 (13): 2037–39. https://doi.org/10.1093/bioinformatics/btx100.

Delcher, A L, S Kasif, R D Fleischmann, J Peterson, O White, and S L Salzberg. 1999. "Alignment of Whole Genomes." *Nucleic Acids Research* 27 (11): 2369–76. https://doi.org/10.1093/nar/27.11.2369.

Denver, Dee R, Peter C Dolan, Larry J Wilhelm, Way Sung, J Ignacio Lucas-Lledó, Dana K Howe, Samantha C Lewis, et al. 2009. "A Genome-Wide View of Caenorhabditis Elegans Base-Substitution Mutation Processes." *Proceedings of the National Academy of Sciences of the United States of America* 106 (38): 16310–14. https://doi.org/10.1073/pnas.0904895106.

Eide, D, and P Anderson. 1985. "Transposition of Tc1 in the Nematode Caenorhabditis Elegans." *Proceedings of the National Academy of Sciences* 82 (6): 1756–60. https://doi.org/10.1073/pnas.82.6.1756.

952  Emmons, Scott W., Lewis Yesner, Ke-san Ruan, and Daniel Katzenberg. 1983. "Evidence for a
953       Transposon in Caenorhabditis Elegans." *Cell* 32 (1): 55–65.
954       https://doi.org/10.1016/0092-8674(83)90496-8.
955  Feschotte, Cédric, and Ellen J. Pritham. 2007. "DNA Transposons and the Evolution of
956       Eukaryotic Genomes." *Annual Review of Genetics* 41: 331–68.
957       https://doi.org/10.1146/annurev.genet.40.110405.090448.
958  Fischer, Sylvia E J, Erno Wienholds, and Ronald H A Plasterk. 2003. "Continuous Exchange of
959       Sequence Information Between Dispersed Tc1 Transposons in the *Caenorhabditis*
960       *Elegans* Genome." *Genetics* 164 (1): 127–34.
961       https://doi.org/10.1093/genetics/164.1.127.
962  Garrison, Erik, and Gabor Marth. 2012. "Haplotype-Based Variant Detection from Short-Read
963       Sequencing." arXiv. https://doi.org/10.48550/ARXIV.1207.3907.
964  Gems, D, and D L Riddle. 2000. "Defining Wild-Type Life Span in Caenorhabditis Elegans." *The
965       Journals of Gerontology. Series A, Biological Sciences and Medical Sciences* 55 (5):
966       B215-9. https://doi.org/10.1093/gerona/55.5.b215.
967  Gilbert, Clément, Jean Peccoud, and Richard Cordaux. 2021. "Transposable Elements and the
968       Evolution of Insects." *Annual Review of Entomology* 66 (1): 355–72.
969       https://doi.org/10.1146/annurev-ento-070720-074650.
970  Girard, Lisa, and Michael Freeling. 1999. "Regulatory Changes as a Consequence of
971       Transposon Insertion." *Developmental Genetics* 25 (4): 291–96.
972       https://doi.org/10.1002/(SICI)1520-6408(1999)25:4<291::AID-DVG2>3.0.CO;2-5.
973  Goel, Manish, Hequan Sun, Wen-Biao Jiao, and Korbinian Schneeberger. 2019. "SyRI: Finding
974       Genomic Rearrangements and Local Sequence Differences from Whole-Genome
975       Assemblies." *Genome Biology* 20 (1): 277. https://doi.org/10.1186/s13059-019-1911-0.
976  Haraksingh, Rajini R, and Michael P Snyder. 2013. "Impacts of Variation in the Human Genome
977       on Gene Regulation." *Journal of Molecular Biology* 425 (21): 3970–77.
978       https://doi.org/10.1016/j.jmb.2013.07.015.
979  Heger, Andreas, Caleb Webber, Martin Goodson, Chris P Ponting, and Gerton Lunter. 2013.
980       "GAT: A Simulation Framework for Testing the Association of Genomic Intervals."
981       *Bioinformatics* 29 (16): 2046–48. https://doi.org/10.1093/bioinformatics/btt343.
982  Hodgkin, J, and T Doniach. 1997. "Natural Variation and Copulatory Plug Formation in
983       Caenorhabditis Elegans." *Genetics* 146 (1): 149–64.
984       https://doi.org/10.1093/genetics/146.1.149.
985  Hurles, Matthew E, Emmanouil T Dermitzakis, and Chris Tyler-Smith. 2008. "The Functional
986       Impact of Structural Variation in Humans." *Trends in Genetics : TIG* 24 (5): 238–45.
987       https://doi.org/10.1016/j.tig.2008.03.001.
988  James, Kent W, Baertsch Robert, Hinrichs Angie, Miller Webb, and Haussler David. 2003.
989       "Evolution's Cauldron: Duplication, Deletion, and Rearrangement in the Mouse and
990       Human Genomes." *Proceedings of the National Academy of Sciences* 100 (20): 11484–
991       89. https://doi.org/10.1073/pnas.1932072100.
992  Kamath, Ravi S, Andrew G Fraser, Yan Dong, Gino Poulin, Richard Durbin, Monica Gotta,
993       Alexander Kanapin, et al. 2003. "Systematic Functional Analysis of the Caenorhabditis
994       Elegans Genome Using RNAi." *Nature* 421 (6920): 231–37.
995       https://doi.org/10.1038/nature01278.
996  Kern, Andrew D, and Matthew W Hahn. 2018. "The Neutral Theory in Light of Natural
997       Selection." *Molecular Biology and Evolution* 35 (6): 1366–71.
998       https://doi.org/10.1093/molbev/msy092.
999  Kim, Bernard Y, Jeremy R Wang, Danny E Miller, Olga Barmina, Emily Delaney, Ammon
1000      Thompson, Aaron A Comeault, et al. 2021. "Highly Contiguous Assemblies of 101
1001      Drosophilid Genomes." Edited by Graham Coop, Patricia J Wittkopp, and Timothy B
1002      Sackton. *eLife* 10: e66405. https://doi.org/10.7554/eLife.66405.

Kim, Chuna, Jun Kim, Sunghyun Kim, Daniel E Cook, Kathryn S Evans, Erik C Andersen, and Junho Lee. 2019. "Long-Read Sequencing Reveals Intra-Species Tolerance of Substantial Structural Variations and New Subtelomere Formation in C. Elegans." *Genome Research* 29 (6): 1023–35. https://doi.org/10.1101/gr.246082.118.

Koch, R, H G van Luenen, M van der Horst, K L Thijssen, and R H Plasterk. 2000. "Single Nucleotide Polymorphisms in Wild Isolates of Caenorhabditis Elegans." *Genome Research* 10 (11): 1690–96. https://doi.org/10.1101/gr.gr-1471r.

Koren, Sergey, Brian P Walenz, Konstantin Berlin, Jason R Miller, Nicholas H Bergman, and Adam M Phillippy. 2017. "Canu: Scalable and Accurate Long-Read Assembly via Adaptive k-Mer Weighting and Repeat Separation." *Genome Research* 27 (5): 722–36. https://doi.org/10.1101/gr.215087.116.

Lappalainen, Tuuli, Alexandra J. Scott, Margot Brandt, and Ira M. Hall. 2019. "Genomic Analysis in the Age of Human Genome Sequencing." *Cell* 177 (1): 70–84. https://doi.org/10.1016/j.cell.2019.02.032.

Laricchia, K.M., S. Zdraljevic, D.E. Cook, and E.C. Andersen. 2017. "Natural Variation in the Distribution and Abundance of Transposable Elements Across the Caenorhabditis Elegans Species." *Molecular Biology and Evolution* 34 (9): 2187–2202. https://doi.org/10.1093/molbev/msx155.

Lee, Daehan, Stefan Zdraljevic, Lewis Stevens, Ye Wang, Robyn E. Tanny, Timothy A. Crombie, Daniel E. Cook, et al. 2021. "Balancing Selection Maintains Hyper-Divergent Haplotypes in Caenorhabditis Elegans." *Nature Ecology & Evolution* 5 (6): 794–807. https://doi.org/10.1038/s41559-021-01435-x.

Lee, Heng-Chi, Weifeng Gu, Masaki Shirayama, Elaine Youngman, Darryl Conte, and Craig C. Mello. 2012. "C. Elegans piRNAs Mediate the Genome-Wide Surveillance of Germline Transcripts." *Cell* 150 (1): 78–87. https://doi.org/10.1016/j.cell.2012.06.016.

Lesack, Kyle, Grace M. Mariene, Erik C. Andersen, and James D. Wasmuth. 2022. "Different Structural Variant Prediction Tools Yield Considerably Different Results in Caenorhabditis Elegans." *PLOS ONE* 17 (12): e0278424. https://doi.org/10.1371/journal.pone.0278424.

Li, Heng. 2018. "Minimap2: Pairwise Alignment for Nucleotide Sequences." *Bioinformatics* 34 (18): 3094–3100. https://doi.org/10.1093/bioinformatics/bty191.

Li, Heng, and Richard Durbin. 2009. "Fast and Accurate Short Read Alignment with Burrows–Wheeler Transform." *Bioinformatics* 25 (14): 1754–60. https://doi.org/10.1093/bioinformatics/btp324.

Li, Heng, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup. 2009. "The Sequence Alignment/Map Format and SAMtools." *Bioinformatics (Oxford, England)* 25 (16): 2078–79. https://doi.org/10.1093/bioinformatics/btp352.

Liao, L. W., B. Rosenzweig, and D. Hirsh. 1983. "Analysis of a Transposable Element in Caenorhabditis Elegans." *Proceedings of the National Academy of Sciences of the United States of America* 80 (12): 3585–89. https://doi.org/10.1073/pnas.80.12.3585.

Mahmoud, Medhat, Nastassia Gobet, Diana Ivette Cruz-Dávalos, Ninon Mounier, Christophe Dessimoz, and Fritz J Sedlazeck. 2019. "Structural Variant Calling: The Long and the Short of It." *Genome Biology* 20 (1): 246. https://doi.org/10.1186/s13059-019-1828-7.

Manni, Mosè, Matthew R Berkeley, Mathieu Seppey, Felipe A Simão, and Evgeny M Zdobnov. 2021. "BUSCO Update: Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes." *Molecular Biology and Evolution* 38 (10): 4647–54. https://doi.org/10.1093/molbev/msab199.

Maydan, Jason S, Adam Lorch, Mark L Edgley, Stephane Flibotte, and Donald G Moerman. 2010. "Copy Number Variation in the Genomes of Twelve Natural Isolates of

Caenorhabditis Elegans." *BMC Genomics* 11 (January): 62. https://doi.org/10.1186/1471-2164-11-62.

McKenna, Aaron, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis, Andrew Kernytsky, Kiran Garimella, et al. 2010. "The Genome Analysis Toolkit: A MapReduce Framework for Analyzing next-Generation DNA Sequencing Data." *Genome Research* 20 (9): 1297–1303. https://doi.org/10.1101/gr.107524.110.

Miller, Danny E, Kevin R Cook, and R Scott Hawley. 2019. "The Joy of Balancers." *PLOS Genetics* 15 (11): e1008421.

Muzzey, Dale, Eric A Evans, and Caroline Lieber. 2015. "Understanding the Basics of NGS: From Mechanism to Variant Calling." *Current Genetic Medicine Reports* 3 (4): 158–65. https://doi.org/10.1007/s40142-015-0076-8.

Nattestad, Maria, and Michael C Schatz. 2016. "Assemblytics: A Web Analytics Tool for the Detection of Variants from an Assembly." *Bioinformatics (Oxford, England)* 32 (19): 3021–23. https://doi.org/10.1093/bioinformatics/btw369.

Nicholas, W L, E C Dougherty, and E L Hansen. 1959. "Axenic Cultivation of Caenorhabditis Briggsae (Nematoda: Rhabditidae) with Chemically Undefined Supplements; Comparative Studies with Related Nematodes." *Ann. N.Y. Acad. Sci* 77: 218–36.

Plasterk, Ronald H.A, Zsuzsanna Izsvák, and Zoltán Ivics. 1999. "Resident Aliens: The Tc1/ Mariner Superfamily of Transposable Elements." *Trends in Genetics* 15 (8): 326–32. https://doi.org/10.1016/S0168-9525(99)01777-1.

Quinlan, Aaron R, and Ira M Hall. 2010. "BEDTools: A Flexible Suite of Utilities for Comparing Genomic Features." *Bioinformatics* 26 (6): 841–42. https://doi.org/10.1093/bioinformatics/btq033.

Riehl, Kevin, Cristian Riccio, Eric A Miska, and Martin Hemberg. 2022. "TransposonUltimate: Software for Transposon Classification, Annotation and Detection." *Nucleic Acids Research* 50 (11): e64–e64. https://doi.org/10.1093/nar/gkac136.

Rockman, Matthew V, and Leonid Kruglyak. 2009. "Recombinational Landscape and Population Genomics of Caenorhabditis Elegans." *PLOS Genetics* 5 (3): e1000419.

Sakamoto, Yoshitaka, Suzuko Zaha, Yutaka Suzuki, Masahide Seki, and Ayako Suzuki. 2021. "Application of Long-Read Sequencing to the Detection of Structural Variants in Human Cancer Genomes." *Computational and Structural Biotechnology Journal* 19: 4207–16. https://doi.org/10.1016/j.csbj.2021.07.030.

Sarsani, Vishal Kumar, Narayanan Raghupathy, Ian T Fiddes, Joel Armstrong, Francoise Thibaud-Nissen, Oraya Zinder, Mohan Bolisetty, et al. 2019. "The Genome of C57BL/6J 'Eve', the Mother of the Laboratory Mouse Genome Reference Strain." *G3 (Bethesda, Md.)* 9 (6): 1795–1805. https://doi.org/10.1534/g3.119.400071.

Shumate, Alaina, and Steven L Salzberg. 2021. "Liftoff: Accurate Mapping of Gene Annotations." *Bioinformatics* 37 (12): 1639–43. https://doi.org/10.1093/bioinformatics/btaa1016.

Sijen, Titia, and Ronald H. A. Plasterk. 2003. "Transposon Silencing in the Caenorhabditis Elegans Germ Line by Natural RNAi." *Nature* 426 (6964): 310–14. https://doi.org/10.1038/nature02107.

Simão, Felipe A, Robert M Waterhouse, Panagiotis Ioannidis, Evgenia V Kriventseva, and Evgeny M Zdobnov. 2015. "BUSCO: Assessing Genome Assembly and Annotation Completeness with Single-Copy Orthologs." *Bioinformatics* 31 (19): 3210–12. https://doi.org/10.1093/bioinformatics/btv351.

Slotkin, R. Keith, and Robert Martienssen. 2007. "Transposable Elements and the Epigenetic Regulation of the Genome." *Nature Reviews Genetics* 8 (4): 272–85. https://doi.org/10.1038/nrg2072.

Stankiewicz, Pawe\l, and James R Lupski. 2010. "Structural Variation in the Human Genome and Its Role in Disease." *Annual Review of Medicine* 61 (1): 437–55. https://doi.org/10.1146/annurev-med-100708-204735.

Stewart, Mary K, Nathaniel L Clark, Gennifer Merrihew, Evan M Galloway, and James H Thomas. 2005. "High Genetic Diversity in the Chemoreceptor Superfamily of Caenorhabditis Elegans." *Genetics* 169 (4): 1985–96. https://doi.org/10.1534/genetics.104.035329.

Stranger, Barbara E., Forrest Matthew S., Dunning Mark, Ingle Catherine E., Beazley Claude, Thorne Natalie, Redon Richard, et al. 2007. "Relative Impact of Nucleotide and Copy Number Variation on Gene Expression Phenotypes." *Science* 315 (5813): 848–53. https://doi.org/10.1126/science.1136678.

Sudmant, Peter H, Tobias Rausch, Eugene J Gardner, Robert E Handsaker, Alexej Abyzov, John Huddleston, Yan Zhang, et al. 2015. "An Integrated Map of Structural Variation in 2,504 Human Genomes." *Nature* 526 (7571): 75–81. https://doi.org/10.1038/nature15394.

Sulston, J E, and S Brenner. 1974. "The DNA of Caenorhabditis Elegans." *Genetics* 77 (1): 95–104. https://doi.org/10.1093/genetics/77.1.95.

Surzycki, S A, and W R Belknap. 2000. "Repetitive-DNA Elements Are Similarly Distributed on Caenorhabditis Elegans Autosomes." *Proceedings of the National Academy of Sciences of the United States of America* 97 (1): 245–49. https://doi.org/10.1073/pnas.97.1.245.

Swan, Kathryn A, Damian E Curtis, Kathleen B McKusick, Alexander V Voinov, Felipa A Mapa, and Michael R Cancilla. 2002. "High-Throughput Gene Mapping in Caenorhabditis Elegans." *Genome Research* 12 (7): 1100–1105. https://doi.org/10.1101/gr.208902.

Thomas, Cristel G, Wei Wang, Richard Jovelin, Rajarshi Ghosh, Tatiana Lomasko, Quang Trinh, Leonid Kruglyak, Lincoln D Stein, and Asher D Cutter. 2015. "Full-Genome Evolutionary Histories of Selfing, Splitting, and Selection in Caenorhabditis." *Genome Research* 25 (5): 667–78. https://doi.org/10.1101/gr.187237.114.

Thompson, Owen A, L Basten Snoek, Harm Nijveen, Mark G Sterken, Rita J M Volkers, Rachel Brenchley, Arjen van't Hof, et al. 2015. "Remarkably Divergent Regions Punctuate the Genome Assembly of the Caenorhabditis Elegans Hawaiian Strain CB4856." *Genetics* 200 (3): 975–89. https://doi.org/10.1534/genetics.115.175950.

Van't Hof, Arjen E, Pascal Campagne, Daniel J Rigden, Carl J Yung, Jessica Lingley, Michael A Quail, Neil Hall, Alistair C Darby, and Ilik J Saccheri. 2016. "The Industrial Melanism Mutation in British Peppered Moths Is a Transposable Element." *Nature* 534 (7605): 102–5. https://doi.org/10.1038/nature17951.

Vergara, Ismael A, Allan K Mah, Jim C Huang, Maja Tarailo-Graovac, Robert C Johnsen, David L Baillie, and Nansheng Chen. 2009. "Polymorphic Segmental Duplication in the Nematode Caenorhabditis Elegans." *BMC Genomics* 10 (1): 329. https://doi.org/10.1186/1471-2164-10-329.

Wicks, Stephen R, Raymond T Yeh, Warren R Gish, Robert H Waterston, and Ronald H A Plasterk. 2001. "Rapid Gene Mapping in Caenorhabditis Elegans Using a High Density Polymorphism Map." *Nature Genetics* 28 (2): 160–64. https://doi.org/10.1038/88878.

Yoshimura, Jun, Kazuki Ichikawa, Massa J Shoura, Karen L Artiles, Idan Gabdank, Lamia Wahba, Cheryl L Smith, et al. 2019. "Recompleting the Caenorhabditis Elegans Genome." *Genome Research* 29 (6): 1009–22. https://doi.org/10.1101/gr.244830.118.

Zhao, Yuehui, Lijiang Long, Wen Xu, Richard F Campbell, Edward E Large, Joshua S Greene, and Patrick T McGrath. 2018. "Changes to Social Feeding Behaviors Are Not Sufficient for Fitness Gains of the Caenorhabditis Elegans N2 Reference Strain." *eLife* 7 (October). https://doi.org/10.7554/eLife.38675.