1 **Multiple Displacement Amplification Facilitates SMRT Sequencing of Microscopic**

2 **Animals and the Genome of the Gastrotrich *Lepidodermella squamata* (Dujardin, 1841)**

3 Nickellaus G. Roberts[1], ngroberts@crimson.ua.edu

4 Michael J. Gilmore[1], gilmore.michael1999@gmail.com

5 Torsten H. Struck[2], t.h.struck@nhm.uio.no

6 *Kevin M. Kocot[1,3] ,kmkocot @ua.edu

7 [1]Department of Biological Sciences, The University of Alabama, [2]Natural History Museum,

8 University of Oslo, [3]Alabama Museum of Natural History

9 *Corresponding author.

10

11

12

13

14

15

16

17

18

19

20     **Abstract:**

21     **Background:**

22     Obtaining adequate DNA for long-read genome sequencing remains a roadblock to producing

23     contiguous genomes from small-bodied organisms. Multiple displacement amplification (MDA)

24     leverages Phi29 DNA polymerase to produce micrograms of DNA from picograms of input. Few

25     genomes have been generated using this approach, due to concerns over biases in amplification

26     related to GC and repeat content and chimera production. Here, we explored the utility of MDA

27     for generating template DNA for PacBio HiFi sequencing using *Caenorhabditis elegans*

28     (Nematoda) and *Lepidodermella squamata* (Gastrotricha).

29     **Results***:*

30     HiFi sequencing of libraries prepared from MDA DNA produced highly contiguous and

31     complete genomes for both *C. elegans* (102 Mbp assembly; 336 contigs; N50 = 868 Kbp; L50 =

32     39; BUSCO_nematoda: S:92.2%, D:2.7%) and *L. squamata* (122 Mbp assembly; 157 contigs;

33     N50 = 3.9 Mb; L50 = 13; BUSCO_metazoa: S: 78.0%, D: 2.8%). Amplified *C. elegans* reads

34     mapped to the reference genome with a rate of 99.92% and coverage of 99.75% with just one

35     read (of 708,811) inferred to be chimeric. Coverage uniformity was nearly identical for reads

36     from MDA DNA and reads from pooled worm DNA when mapped to the reference genome. The

37     genome of *Lepidodermella squamata*, the first of its phylum, was leveraged to infer the

38     phylogenetic position of Gastrotricha, which has long been debated, as the sister taxon of

39     Platyhelminthes.

40     **Conclusions:**

41    This methodology will help generate contiguous genomes of microscopic taxa whose body size

42    precludes standard long-read sequencing. *L. squamata* is an emerging model in evolutionary

43    developmental biology and this genome will facilitate further work on this species.

44

45    **Background:**

46    In recent years, long-read technologies have become the industry standard for *de novo* genome

47    sequencing. Pacific Biosciences' (PacBio) multiple pass circular consensus sequencing (HiFi)

48    has stood out as a popular choice for *de novo* genome sequencing due to its high accuracy (base-

49    level resolution of >99%) and relatively long reads (~10 Kbp) able to cover complex and

50    repetitive genomic regions (Hon et al. 2020). Using adequate high-molecular-weight (HMW)

51    DNA with as little damage as possible for sequencing library preparation is crucial to obtain

52    reads that are both highly accurate and of adequate length (Pollard et al. 2018). Avoiding pooling

53    multiple individuals serves to increase the assembly quality of genomes by reducing

54    heterozygosity (Birky 1996; Flot et al. 2013; Huang et al. 2017; Simion et al. 2018). It follows

55    that *de novo* genome projects are most likely to succeed when a large amount of unfragmented,

56    HMW DNA with little to no exogenous contamination is available from a single individual.

57    Unfortunately, obtaining such amounts of starting material meeting the above requirements is not

58    always possible. PacBio's Ultra-Low DNA Input Workflow yields data volumes comparable to

59    standard input libraries from as little as 5 ng of HMW template. This approach has successfully

60    been used to produce high-quality genomes from individual small animals such as mosquitos

61    (Kingan et al. 2019) and springtails (Schneider et al. 2021). However, even the Ultra-Low DNA

62    Input Workflow is not applicable to very small animals (e.g., < 1 mm total length). For example,

63    a single hermaphrodite specimen of *Caenorhabditis elegans* contains 959 somatic nuclei (Cohen-

64    Fix and Askjaer 2017) and has a genome size of roughly 100 Mbp (*C. elegans* Sequencing

65    Consortium 1998) meaning that a single specimen may contain as little as 200 picograms of

66    DNA with comparable DNA content in other small-bodied phyla, also ranging in the hundreds of

67    picograms of maximal yield. Thus, obtaining sufficient starting material meeting even the greatly

68    reduced input requirements for the PacBio Ultra-Low DNA Input Workflow is not possible for

69    many organisms. Further, if researchers want to avoid pooling multiple individuals or if the

70    organism of interest is both small and rare, reduced input strategies are not an option thus far.

71    It is of little debate that for most applications leveraging genomes, assemblies that are both

72    contiguous and well-annotated are ideal (reviewed by Laumer 2018). However, due to the

73    challenges associated with sequencing genomes of small-bodied animals, metazoan genome

74    sequencing efforts have been biased towards large-bodied organisms. Of the roughly 30

75    metazoan phyla, nearly all those that are still lacking high quality genome assemblies – or any

76    genomic data at all in some cases – are small-bodied (Worsaae et al. 2023). For example, within

77    Lophotrochozoa (=Spiralia) all phyla that are lacking in genomic data are mostly or exclusively

78    small-bodied (Figure 1, Figure S1). With an average body size of <1 mm**,** these organisms are

79    underrepresented in genomic and phylogenomic studies (Kocot et al. 2017; Laumer et al. 2019).

80    Whole genome amplification (WGA) methodologies are of great use to the fields of diagnostic

81    medicine and single cell research (Evrony et al. 2015; Fu et al. 2015; Li et al. 2018; Zhou et al.

82    2020), but their utility in amplifying DNA for sequencing very small-bodied organisms in the

83    age of long-read sequencing technologies has received little attention (but see Lee et al. 2023).

84    WGA is a popular strategy for generating large amounts of DNA from limiting sample material

85    through techniques such as multiple displacement amplification (MDA) (Dean et al. 2001),

86  degenerate-oligonucleotide-primed-PCR (DOP-PCR) (Telenius et al. 1992), DOP-PCR (Telenius

87  et al. 1992), PEP-PCR (Moghaddaszadeh-Ahrabi et al. 2012), LM-PCR (Carey MF et al. 2009),

88  T-PCR (Tarr AW et al. 2005), and multiple annealing and looping based amplification cycles

89  (MALBAC) (Chapman et al. 2015). However the use of WGA for long-read sequencing library

90  preparation has largely been limited to the detection of structural variants, large deletions and

91  inversions, and amplicon sequencing from degraded material (Evrony et al. 2015; Zhou et al.

92  2020).

93  WGA methodologies typically make use of PCR, isothermal amplification, or both (reviewed by

94  Wang et al. 2022). PCR based protocols make use of specific primers, degenerate oligos, and/or

95  repetitive genomic regions for priming. These strictly PCR-based amplification strategies are

96  exponential and thus uneven or non-specific amplification may occur due to variation in template

97  GC content and amplicon length. PCR-based methodologies have been used in medical and

98  diagnostic approaches for their ability to analyze copy number variation and identify small

99  structural variants or small stretches of SNP variation, but their utility for long-read sequencing

100  is limited due to typical product fragment sizes only being within 100-1000 bp. Long-range PCR

101  based approaches have been demonstrated to be successful for *de novo* genome sequencing of

102  specimens in which material is limited (Laumer 2023; Stevens et al. 2023). Picogram input

103  multimodal sequencing (PiMmS) makes use of oligo-dT beads to separate cDNA and genomic

104  DNA from the same individual, providing transcript evidence from the same individual for

105  downstream gene annotation. PiMmS amplification technique involves DNA purification using a

106  salting out procedure to provide HMW ethanol-precipitated DNA. As PiMmS makes use of

107  PCR, a qPCR step is required to estimate cycle number to reduce artifacts produced during PCR

108  amplification.

109    Isothermal amplification in the case of MDA uses Phi29 DNA polymerase to amplify long

110    fragments of linear or circular DNA (Dean et al. 2001; Dean et al. 2002; Lasken 2007).

111    Amplification using Phi29 polymerase is highly accurate due to high strand-displacement

112    activity, extremely high processivity, and strong proofreading activity (Garmendia et al. 1992).

113    MDA, as opposed to PCR-based WGA approaches, has the benefit of producing highly accurate

114    fragments with an average fragment length greater than 10 Kbp (up to 100 Kbp). Like other

115    WGA approaches, MDA can produce tens of micrograms of DNA from picograms of starting

116    material. MDA typically makes use of short, random exonuclease-resistant primers to allow the

117    reaction to proceed isothermally. Like PCR-based methods, MDA is exponential leading to

118    potential biases in coverage based on GC content (Lasken and Egholm 2003; Borgström et al.

119    2017). In GC-rich regions, there may be reduction of polymerase processivity leading to

120    underrepresentation of that region during amplification. However, differential binding affinities

121    of GC- and AT-rich primers can lead to GC-rich regions being preferentially amplified (Benita et

122    al. 2003; Sahdev et al. 2007). Further, the denaturation strategy used prior to MDA (alkali,

123    thermal or otherwise) may affect access and efficacy of DNA binding in GC rich regions. For

124    example, KOH alkali denaturation has been shown to increase primer access to GC rich template

125    DNA (Pinard et al. 2006). Such amplification bias with respect to GC richness is a function of

126    reaction gain, defined as the ratio of DNA product mass over DNA template mass. Compared to

127    MALBAC and PicoPLEX single cell, increased reaction gain in the case of MDA more

128    negatively influences amplification coverage, uniformity, and the detection of copy number

129    variants (CNVs). However, unlike both MALBAC and PicoPLEX, MDA reactions with high

130    gain do not have a meaningful increase in the rate of single nucleotide errors (de Bourcy et al.

131    2014). Fortunately, lowering amplification time reduces reaction gain, thus reducing the over-

132    amplification of certain genomic regions (Dean et al. 2001), making MDA protocols with more

133    modest amplification time preferable to generate template for de novo sequencing. Lowering

134    reaction gain also can also reduce chimera formation. During MDA, chimeric DNA fragments

135    can form when 3' termini are displaced and anneal to nearby displaced 5' strands, forming

136    chimeric DNA fragments that are eventually made into double stranded DNA after random

137    hexamer annealing and extension (Lasken and Stockwell 2007). Previous studies have suggested

138    a rate of up to 1 chimera/10 Kbp for MDA, which is problematic for the recovery of specific

139    genomic regions, the analysis of copy number variation and *de novo* assembly, although these

140    problems can be reduced with adequate coverage (Rodrigue et al. 2009). Proper preservation of

141    samples and HMW DNA extraction (resulting in fewer 3' termini) can reduce chimera formation

142    (Nelson 2014), especially when combined with lower reaction gain. Moreover, as MDA can be

143    applied directly to lysed cells (Hosono et al. 2003), shearing of DNA during extraction can be

144    avoided altogether when aliquots of cell cultures or entire microscopic organisms are lysed

145    followed by immediate amplification in the same tube. Taken together, MDA's isothermal

146    nature, high accuracy, production of fragments longer than typical PacBio HiFi reads, and ability

147    to amplify DNA directly from freshly lysed cells makes it an attractive choice for template

148    preparation for HiFi sequencing. Under low reaction gain and in low volumes, sequencing of

149    MDA DNA can achieve high data uniformity and coverage (Li et al. 2018). Thus, MDA was

150    chosen for this study.

151    Here, we investigated the utility of MDA for generating adequate template DNA for PacBio

152    long-read sequencing from very small animals (<1 mm). We demonstrate the utility of this

153    approach using the model organism *Caenorhabditis elegans,* which was selected because a high-

154    quality reference genome and HiFi reads generated from unamplified DNA extracted from a pool

155    of worms are already available for comparison. We additionally chose the gastrotrich

156    *Lepidodermella squamata* (Dujardin 1841) as a test of MDA's ability to assist in the *de novo*

157    long-read sequencing of a non-model organism. Gastrotricha (Metschnikoff, 1865), colloquially

158    referred to as "hairy bellied worms" is a phylum of approximately 749 species of microscopic

159    benthic invertebrates found in aquatic habitats worldwide (Margulis and Schwartz 1998). Of

160    these, 346 are found in freshwater habitats, with the rest inhabiting marine and brackish water.

161    Gastrotricha is composed of two orders, Chaetonotida (Remane, 1925), which include both fresh

162    and marine species, and Macrodasyida (Remane, 1925), which is an almost exclusively marine

163    clade. Living mostly in sandy, benthic habitats, key taxonomic characters of Gastrotricha include

164    the variable presence of scales, spines, or hooks on the body wall and the presence of adhesive

165    tubes for attaching to substrate. Gastrotrichs typically range from 70 μm to 1 mm in size,

166    although they may reach up to 3 mm in the case of the genus *Megadasys* (Schmidt, 1974,

167    Fontaneto et al. 2015). Additionally, gastrotrichs are flattened along their dorsal-ventral axis

168    reducing their volume even further in comparison to similar sized animals with a round body

169    plan. The ~250 μm freshwater gastrotrich *L. squamata* has proven to be a useful species for the

170    study of Gastrotricha due to its ease of culture in a freshwater solution containing wheat grains

171    inoculated with fungi, protozoa, and bacteria (Bennett 1979). Gastrotricha is of  great interest to

172    evolutionary biologists as its phylogenetic position has long been debated (Edgecombe et al.

173    2011; Struck et al. 2014; Laumer et al. 2015; Kocot 2016; Bleidorn 2019). Availability of a high-

174    quality reference genome for *L. squamata* would also to the field of evolutionary developmental

175    biology as *L. squamata* is an emerging model species (Hejnol 2015).

176

177    **Results:**

178  Genome Assembly

179  Sequencing of a HiFi library generated from MDA DNA for *C. elegans* resulted in 708,811 HiFi

180  reads (96x coverage of the 102.5 Mbp genome) with an average HiFi read length of 13,885 bp.

181  The *de novo C. elegans* assembly was 102 Mbp contained in 336 contigs with an N50 of 868

182  Kbp and an L50 of 39. The *de novo* assembly was screened for contamination with Blobtools

183  (Laetsch and Blaxter 2017), however all contigs had hits only to Metazoa (336) and thus none

184  were removed. The *de novo* assembly had a BUSCO metazoa_odb10 score of 69.9% (Single

185  copy BUSCOs [S]:67.5%, Duplicate BUSCOs [D]:2.4%) and a BUSCO nematoda_odb10 score

186  of 94.9% (S:92.2%, D:2.7%) compared to 51.3% (S:50.7%, D:0.6%) and 98.5% (S:98.0%,

187  D:0.5%), respectively, for the reference. The reference-guided *C. elegans* assembly was 102.5

188  Mbp contained in 281 contigs with an N50 of 16.9 Mbp and an L50 of 3 (Table 1).

189  Sequencing of *L. squamata* resulted in 917,258 HiFi reads with an average HiFi read length of

190  13,358 bp. The *L. squamata* assembly was 147 Mbp contained in 527 contigs with an N50 of 2.9

191  Mbp and an L50 of 13. After filtering with Blobtools, 370 contigs were removed based on best

192  hits to non-metazoans. The resulting assembly was 122 Mbp contained in 157 contigs with an

193  N50 of 3.9 Mbp and an L50 of 13. BUSCO analysis of the decontaminated assembly using the

194  metazoa_odb10 dataset was 80.8% (Single copy BUSCOs [S]: 78.0%, Duplicate BUSCOs [D]:

195  2.8%; Table 1).

196  **Table 1**. Assembly Statistics for *C. elegans* and *L. squamata*

| Organism: | HiFi Reads: | Genome Size (Mb): | Contigs: | N50 (Mb): | L50 : | BUSCO metazoa_odb 10 | BUSCO nematoda_odb 10 |
|---|---|---|---|---|---|---|---|
| *Lepidodermella squamata [MDA]* | 917258 | 122 | 157 | 3.9 | 13 | 80.8% (S: 78.0%, D: 2.8%) | N/A |
| *Caenorhabditis elegans [MDA]* | 708811 | 102 | 336 | 0.87 | 39 | 69.9% (S: 67.5%, D: 2.4%) | 94.9% (S: 92.2%, D: 2.7%) |
| *Caenorhabditis elegans [Reference Guided MDA]* | N/A | 102.5 | 281 | 16.9 | 3 | 51.3% (S: 50.7%, D:0.6%) | 98.5% (S: 98.0%, D: 0.5%) |

197

## Coverage and Amplification Bias Estimation

199 Due to MDA's potential for amplification bias with respect to GC and repeat content (Lasken

200 and Egholm 2003; Borgström et al. 2017), we investigated the relationship between GC and

201 repeat content and read coverage using the model organism *C. elegans* and its reference genome.

202 To understand the level as to which MDA may have resulted in over amplification of genomic

203 regions with respect to GC content, we compared mean coverage in 100 Kbp blocks of reads

204 from our dataset based on MDA and reads from unamplified DNA extracted from a pool of

205    worms (NCBI SRR22507561) mapped to the reference genome (GCA_000002985.3) (Figure 2,

206    Figure S2). Due to the potential for repetitive regions to have GC content far outside the mean,

207    we also investigated coverage with respect to repeat percentage. The reference *C. elegans* Bristol

208    N2 genome has a GC content of 35% (GCA_000002985.3) and, according to RepeatModeler, it

209    has a repeat content of 18.40%. Mean coverage was normalized for each dataset by subtracting

210    the minimum mean coverage from each datapoint (representing a 100 Kbp region) divided by the

211    range of the dataset [(x - min(x))/range(x)), x = coverage of 100 Kbp region]. Resulting

212    normalized mean coverage, repeat content, and GC content in 100 Kbp blocks across all six

213    chromosomes were plotted to visualize potential over-amplification of genomic regions during

214    MDA (Figure 3, Figure 4A-B, Figure S3, Figure S4). This revealed that while some genomic

215    regions showed over-amplification, amplification across the chromosome was consistent enough

216    to provide adequate coverage. Plotting normalized coverage against GC content revealed

217    virtually no difference between data generated from amplified and unamplified DNA (Figure 2)

218    and the distribution was very uniform with respect to GC content with a few outliers.

219    Surprisingly, the reads generated from unamplified DNA were more skewed towards higher

220    coverage for GC values above the mean.

221    Lorenz curves, defined here as a plot of the cumulative share of reads against cumulative share

222    of genomic positions covered by those reads, ordered from highest to lowest (Zong et al. 2012;

223    de Bourcy et al. 2014) were calculated (Figure 4, D). The Gini index is calculated as the area

224    between the calculated Lorenz curve, from the dataset, and a straight Lorenz curve, representing

225    perfect uniform coverage. In our case, a Gini index of zero would indicate perfect uniform

226    coverage and a Gini index of 1 would indicate maximally non-uniform coverage. The Gini index

227    values for the amplified and unamplified datasets were nearly identical, 0.14 and 0.16,

228    respectively, revealing that the uniformity of coverage between the unamplified and amplified *C.*

229    *elegans* read sets is similar. This suggests little to no effect of MDA amplification bias on

230    coverage uniformity in the case of *C. elegans*. Coverage statistics, including mean coverage

231    above and below mean GC content are summarized in Table 2.

232    **Table 2**. *Coverage Statistics and Gini index values for C. elegans and L. squamata unmapped*

233    *reads*.

| Mapped Reads: | Mean Coverage Across all 100 Kbp Blocks: | Mean Coverage of 100 Kbp regions above mean GC content. | Mean Coverage of 100 Kbp regions below mean GC content. | Gini Index : |
|---|---|---|---|---|
| *C. elegans* [MDA] | 98 [SD +/-62] | 102 [SD +/-77] | 92 [SD +/-19] | 0.14 |
| *C. elegans* [Unamplified] | 136[SD +/- 69] | 136 [SD +/- 77] | 136[SD +/- 19] | 0.16 |
| *L. squamata* [MDA] | 117 [SD+/- 102] | 142 [SD +/- 135] | 118[SD +/- 42] | 0.29 |

234

235    Without a reference genome, no mapping comparisons could be made for *L. squamata*, however

236    reads were mapped to the primary assembly after removal of contamination with Blobtools to

237    assess coverage with respect to GC content and repeat content. The resulting genome assembly

238    of *L. squamata* has a GC content of 43% and a repeat content of 21.98%. A Lorenz curve was

239    plotted and the Gini index was calculated as described above. This resulted in a Gini index of

240    0.29 indicating less uniform coverage than both *C. elegans* reads based on MDA DNA and reads

241    based on unamplified DNA (Figure 4D). *L. squamata* was comparatively not as uniform with

242    respect to GC and repeat content as observed for *C. elegans* showing higher coverage in

243    repetitive regions (Figure 4C), which also happen to have higher than mean GC richness. While

244    few 100 Kbp blocks with average and lower GC content have higher fold coverage (Figure 5,

245    Figure S5), a very substantial part of the coverage at GC content values above 45% exhibits very

246    high coverage at a repeat content higher than 50% (Figure 6, Figure S6). This indicates possible

247    preferential amplification of repetitive regions with high GC content. Coverage statistics,

248    including mean coverage above and below mean GC content are summarized in Table 2.

249    Annotation of *L. squamata* genome

250    BRAKER predicted 28,377 gene models with a BUSCO score of 79.80% (76.4% single-copy

251    and 3.2% duplicated). The final assembly was estimated to be 21.98% repetitive. Of the 122

252    Mbp assembly, ~16 Mbp (13.74%) was composed of unclassified repeats, while the remaining

253    ~26 Mbp (8.24%) was composed of simple repeats.

254    Phylogenetic Analysis

255    To infer the phylogenetic position of Gastrotricha within Lophotrochozoa, we identified

256    orthologs in the genomes of 23 lophotrochozoans (plus five outgroup taxa) broadly spanning the

257    diversity of the group. We assembled a supermatrix of 2,779 genes totaling 676,632 amino acids

258    with 18.68% missing data. Maximum likelihood analysis using the posterior mean site frequency

259    model (PMSF) (Wang et al. 2018) fitting the C60 profile mixture model resulted in a tree with

260    full bootstrap support for all nodes (Figure 7). Gastrotricha was recovered as the sister taxon of

261  Platyhelminthes with maximal support, recapitulating the proposed sister-group relationship

262  recovered previously (Struck et al. 2014) and dubbed Rouphozoa. Rotifera, the only

263  representative of Gnathifera in this study, was recovered as the sister group to all other

264  Lophotrochozoa with Rouphozoa as the next branching lineage and the sister taxon of all

265  remaining taxa (Trochozoa) consistent with Struck et al. (2014). Within Trochozoa, we recover

266  Mollusca as the sister taxon of a clade comprising, Annelida, Nemertea and the Lophophorate

267  phyla (Brachiopoda, Bryozoa, and Phorinida), consistent with several previous studies (Struck et

268  al. 2014; Kocot et al. 2017; Laumer et al. 2019). Due to the tree being based exclusively on

269  genomes, several lineages of the superphylum, specifically most gnathiferan phyla

270  (Chaetognatha, Gnathostomulida and Micrognathozoa), Cycliophora, and Entoprocta could not

271  be included in this analysis.

272

273  **Discussion:**

274  Here, we investigated the utility of MDA for generating adequate template DNA from small-

275  bodied organisms for PacBio HiFi sequencing. Although coverage uniformity with respect to GC

276  and repeat content are known biases of MDA, our results show that MDA of DNA obtained

277  directly from lysed cells can be used to generate template for HiFi sequencing with limited

278  coverage bias and artifacts.

279  Given existing high-quality genomic resources, *C. elegans* was chosen to explore the uniformity

280  of coverage with respect to GC and repeat content by comparing our HiFi data generated from

281  MDA DNA to HiFi data generated from unamplified DNA extracted from a pool of worms.

282  Compared to HiFi reads generated unamplified DNA (Gini Index of 0.16), reads generated from

283     MDA DNA (Gini Index of 0.14) performed similarly in terms of coverage uniformity,

284     demonstrating that starting with samples as small as one-half *C. elegans* does not result in

285     sequence data with significant amplification bias. While our *de novo C. elegans* assembly is not

286     chromosome-level, it is rather contiguous and performed well in terms of BUSCO completeness

287     (Table 1). When compared to PacBio Ultra-Low Input assemblies from single organisms, such as

288     springtails (Schneider et al. 2021), our *C. elegans* assembly exhibits comparable completeness,

289     albeit with lower contiguity. While our assembly performed well in terms of completeness

290     (BUSCO Nematoda score of 94.9%), N50 was notably lower (868 Kbp) and the assembly

291     consisted of a higher number of contigs (336) when compared to both springtail assemblies (N50

292     = 5.6 Mbp, 142 contigs, 211 Mbp size and N50 = 8.7 Mbp 165 Mbp size, 79 contigs for *Desoria*

293     *tigrina* and *Sminthurides aquaticus* respectfully). This suggests that the PacBio Ultra-Low Input

294     Workflow is likely a better choice than our MDA-based strategy for animals that can provide

295     adequate HMW template DNA (>5 ng), but that would not be possible in the case of *C. elegans*

296     and many (if not most) meiofaunal animals where a whole animal contains roughly an order of

297     magnitude less DNA than required.

298     We were also concerned about the potential for formation of chimeras. However, investigation of

299     inverted sequences and tandem regions with Alvis identified only one suspected chimeric contig

300     in our *C. elegans* assembly and none in our *L. squamata* assembly, supporting previous

301     hypotheses about chimera formation being highly correlated with high reaction gain (Lasken and

302     Stockwell 2007; de Bourcy et al. 2014). Also, because we lysed cells and immediately proceeded

303     with MDA in the same tube without a DNA extraction procedure, our template DNA underwent

304     little handling that could cause breakage. This would result in fewer free 3' ends of strands that

305     could anneal to another strand and act as a primer to initiate chimera formation. As chimeric

306    reads may be more problematic in data obtained from more fragmented template DNA, paired-

307    end Illumina data generated from half of the specimen lysate before MDA could be used in

308    future studies to aid in the detection of chimeras (see Lee et al. 2023). Examining the mapping of

309    paired-end short reads to post-amplification HiFi reads could be used to reveal chimeric reads as

310    one of the reads in a pair would be expected to map on one side of the chimeric junction while its

311    mate would be expected to map elsewhere.

312    Our results demonstrate that this approach can be successfully used on small-bodied organisms

313    without a reference genome to produce assemblies that are both rather contiguous and complete.

314    The *L. squamata* assembly was rather contiguous with an assembly consisting of 157 contigs and

315    an N50 of 3.9 Mbp approaching the level of contiguity produced using the PacBio Ultra Low

316    Input Workflow (Schneider et al. 2021). The inferred BUSCO completeness score of 80.8% for

317    *L. squamata* based on the metazoa_odb10 set suggests that this genome is lacking nearly 20% of

318    the nuclear protein-coding genes thought to be single-copy in all metazoans. However, *L.*

319    *squamata* is a long-branched taxon in our phylogenomic analysis and Gastrotricha as a phylum

320    has been shown to be long branched in previous phylogenomic studies (Struck et al. 2014; Kocot

321    et al. 2017; Laumer et al. 2019; Marlétaz et al. 2019). Moreover, there are no gastrotrichs or

322    early-branching flatworms in the metazoa_odb10 reference dataset (Manni et al. 2021). Like the

323    reference genome of *C. elegans* (BUSCO metazoa_odb10 score of 51.3%), *L. squamata* may

324    have performed poorly in this BUSCO analysis due to sequence divergence rather than absence

325    of these genes, although loss of BUSCOs in compact invertebrate genomes has been reported

326    previously (Cunha et al. 2023). Regardless, the score recovered here represents a relatively high

327    level of completeness when compared to other sequenced members of the lophotrochozoan

328    superphylum (Figure S1). Annotation of the genome resulted in 28,377 genes, which BUSCO

329 found to be 79.8% complete, which is comparable to the BUSCO score for the genome

330 assembly, suggesting that the predicted gene model set is of high quality. Taken together, these

331 results demonstrate that MDA combined with highly accurate HiFi sequencing can result in

332 genome assemblies of similar quality to those of other non-model invertebrates (Rayko et al.

333 2020).

334 Despite the high contiguity of the *L. squamata* assembly, more bias with respect to repeat

335 content and GC richness was observed in these data than in those from *C. elegans*. When

336 compared to the *C. elegans* reference genome (18.40% repeats, 35% GC content), *L. squamata*

337 (21.98% repeats, 43% GC content) has a higher repeat content. As a result of increased repeat

338 percentage and higher overall GC content, we observed increased coverage in regions with >45%

339 GC content, particularly in 100 Kbp blocks containing more than 50% repeats and in relatively

340 short contigs. Amplification bias in GC-rich regions is a known artifact of MDA (Dean et al.

341 2002; de Bourcy et al. 2014), and our results provide supporting evidence for this. However,

342 provided ample coverage, we were still able to capture enough breadth of the genome to produce

343 a relatively complete and contiguous assembly.

344 MDA has been explored in the context of Oxford Nanopore (ONT) sequencing to produce

345 relatively contiguous genomes from single nematodes (Lee et al. 2023). For *C. elegans*, the

346 authors generated an assembly consisting of 499 contigs with an N50 of 656.2 Kbp and a

347 BUSCO completeness score of 97.6%, compared to 336 contigs with an N50 of 868.1 Kbp and a

348 BUSCO score of 92.2% in the present study. However, applying this approach to 13 other non-

349 model nematodes, they recovered genome assembles ranging from 136.6 Mbp to 738.8 Mbp in

350 1,557 to 32,219 contigs with N50 values of 26.3 Kbp to 441.0 Kbp. Drops in coverage were

351 observed around highly repetitive regions but coverage with respect to GC content was not

352    examined in detail. Notably, T7 endonuclease digestion was needed in the case of ONT

353    sequencing, as MDA's secondary branching products can clog sequencing pores, reducing yield.

354    Due to the nature of library preparation following shearing and size selection, endonuclease

355    digestion is not needed for HiFi sequencing, a benefit of the approach presented here.

356    Given the successful application of this method to the organisms studied here, an important

357    direction for future work would be to apply this strategy to other small-bodied phyla that are

358    currently lacking published genomes (Chaetognatha, Gnathostomulida, Micrognathozoa,

359    Entoprocta and Cycliophora). Increased genomic sampling of small-bodied phyla will lead to an

360    elevated understanding about the evolution of small-bodied organisms and will assist in

361    uncovering potential signatures and commonalities underlying the genomic architecture of small-

362    bodied animals. The genomic consequences of miniaturization are poorly understood given the

363    dearth of genomes from microscopic animals. Work to date has revealed that while some small-

364    bodied lineages exhibit drastic simplification of genome architecture and reduction in content

365    with respect to their large-bodied relatives (Seo et al. 2001; Mikhailov et al. 2016), have taken a

366    more conservative route in regards to genome compaction (Martín-Durán et al. 2021; reviewed

367    by Worsaae et al. 2023). Increased genomic sampling broadly spanning the metazoan tree stands

368    to provide insight into genomic signatures of miniaturization. While our results demonstrate the

369    utility of an MDA-based genome sequencing strategy for small-bodied animals with small

370    genomes (i.e., 102 Mbp and 122 Mbp), the efficacy of this approach for organisms that are

371    small-bodied, but with larger genomes still needs rigorous examination.

372    Combination of this approach with other low input sequencing techniques, such as PiMmS

373    (Laumer, 2022), may be a valuable strategy as these techniques are expected to differ in their

374    biases with respect to amplification uniformity. In contrast to long-range PCR amplification used

375    in PiMmS and the PacBio Ultra-Low Input Workflow, the isothermal amplification used in

376    MDA eliminates PCR duplicates at the cost of reduced uniformity across GC rich regions,

377    particularly those high in repeats (Figure 3, Figure 4A, B). Here, DNA was amplified directly

378    from lysed cells, bypassing HMW DNA extraction that may result in DNA loss from

379    exceptionally small specimens and fragmentation, which we view as a beneficial characteristic of

380    this approach. Careful consideration of the potential amount of DNA in a sample will aid in

381    decision making regarding amplification directly from lysed cells as performed here versus

382    HMW DNA extraction. Library diversity and complexity are both important considerations

383    following any amplification protocol, as low complexity libraries, resulting from degraded or too

384    little template DNA, can lead to poor quality assemblies. Finally, MDA's yield of large amounts

385    of DNA from limiting starting material allows multiple sequencing libraries to be produced from

386    a single round of amplification, making it possible to increase coverage across under-amplified

387    genomic regions, such as those with lower GC content.

388

389    **Conclusion:**

390    The amount of starting DNA has long been a limiting factor in *de novo* genome projects and,

391    despite recent advances, many organisms still contain too little DNA for genome sequencing

392    using available long-read sequencing library preparation protocols. In this study, the efficacy of

393    whole genome sequencing following multiple displacement amplification with phi29 DNA

394    polymerase was assessed. We demonstrate that HiFi sequencing data produced from MDA DNA

395    can successfully lead to a fairly contiguous and complete assembly for both model organisms

396    such as *C. elegans* (Nematoda) and for non-model organisms such as *L. squamata* (Gastrotricha)

397 with limited coverage bias and an extremely low incidence of chimeras. This methodology could

398 be further expanded upon or combined with other scaffolding techniques or amplification

399 strategies to increase the contiguity of genomes from amplified material. This strategy has the

400 potential to greatly expand the availability of whole genome sequencing data from a wide variety

401 of DNA-limited sources including but not limited to other microscopic metazoans.

402

403 **Materials and Methods**

404 Laboratory Methods

405 Live cultures of *C. elegans* (Bristol N2) and *L. squamata* were purchased from Carolina

406 Biological Supply Co. and kept alive in the laboratory under manufacturer recommended

407 conditions. Half of a single individual *C. elegans* hermaphrodite (~500 μm) cut with a sterile

408 razor blade and a single individual *L. squamata* (~190 μm) were added to 4 μl of Qiagen REPLI-

409 g Advanced Single Cell Storage Buffer using a sterilized Irwin Loop as input for amplification

410 using the Qiagen REPLI-g Advanced Single Cell kit following the manufacturer's protocol (i.e.,

411 lysis with SDS and denaturation at room temperature, incubation at 30°C for 2 hours, 60°C for

412 10 minutes, and 4°C until purification). Resulting amplified DNA, 14.92 μg from *L. squamata*

413 and 8.74 μg from *C. elegans* were sent to Hudson Alpha Discovery for sequencing on a shared

414 single PacBio Sequel II flow cell, resulting in ½ a flow cell each.

415 Assembly

416 HiFi reads were extracted using the pbbioconda extracthifi package and each set of reads was

417 assembled with Hifiasm v.0.15.2. A reference guided assembly, using the genome of *C. elegans*

418    Bristol N2 was generated for *C. elegans* using RagTag v.2.1.0 (-scaffold) (Alonge et al. 2022)

419    with default settings. Genome statistics were generated using QUAST v.5.0 (Gurevich et al.

420    2013), and completeness was measured with BUSCO 4.0.2 (Simão et al. 2015). [Table 1,2]

421    Completeness of all assemblies was analyzed with the BUSCO metazoa_odb10 BUSCO set, and

422    additionally, both the *de novo* and reference guided *C. elegans* assemblies were assessed with the

423    nematoda_odb10 BUSCO set. HiFi reads were mapped back to their respective assemblies using

424    Minimap2 v.2.22 (Li 2018). The *C. elegans* and *L. squamata* assemblies were searched against

425    the UniProt reference proteomes database with Diamond 2.0.14 (Buchfink et al. 2015) using the

426    following settings: diamond blastx --db ~/databases/uniprot_taxonmap.dmnd --query

427    assembly.bp.p_ctg.fasta --outfmt 6 qseqid staxids bitscore qseqid sseqid pident length mismatch

428    gapopen qstart qend sstart send evalue bitscore --sensitive --max-target-seqs 1 --evalue 1e-25 --

429    threads 32. The output of these tools as well as BUSCO were used as input for Blobtools.

430    Further, the *de novo C. elegans* assembly was aligned to the publicly available reference *C.*

431    *elegans* N2 genome (GCA_000002985.3) to assess completeness and screen for the presence of

432    chimeric reads as identified by Alvis v.1.0 (Martin and Leggett 2021), a program that identifies

433    chimeras in the assembly based on mapping statistics from mapping raw reads to the assembled

434    genome. All *de novo* assemblies (*C. elegans* and *L. squamata*) were screened for contamination

435    using Blobtools2 v2.6.1 by removing all contigs that had top BLAST hits to taxa other than

436    Metazoa or "no-hit" when searched against the NCBI non-redundant protein database.

437    Annotation

438    The repeat content of the contaminant filtered *L. squamata* genome was assessed using

439    RepeatModeler (Flynn et al. 2020) and soft masked using RepeatMasker (Smit et al. 2015). Raw

440    transcript reads from the publicly available *L. squamata* transcriptome (SRR1273732) were

441    downloaded, quality trimmed with TrimGalore v.0.6.10 from NCBI and mapped against the soft

442    masked genome using STAR (Dobin et al. 2013). The soft masked genome and mapped

443    transcripts were given to the BRAKER (Hoff et al. 2019) pipeline as input to produce gene

444    annotations with the following command: braker.pl --useexisting --cores 12 --softmasking --

445    UTR=on --crf --makehub --gff3 --species=Lepidodermella_squamata --genome

446    Lepidodermella_sp.asm.bp.p_ctg_filtered.fasta.masked --bam RNAseq.bam. Gene annotations

447    were measured for completeness using the BUSCO metazoa_odb10 database.

448    Read Coverage Estimation

449    Assemblies of all three data sources (*L. squamata* and amplified and non-amplified *C. elegans*)

450    were used as subjects to map HiFi reads as a query to estimate genome coverage and coverage

451    bias. Following read mapping with Minimap2 v.2.22, Bedtools v.2.30.0 (Quinlan 2014) was used

452    to separate the genomes into independent 100 Kbp windows in which both overall read coverage,

453    mean coverage, repeat content and GC content were calculated. Resulting graphs were

454    constructed in the R computing environment (R Core Team 2022). Lorenz graphs were

455    constructed utilizing the package *gglorenz*. Additionally, we assessed coverage uniformity of our

456    *C. elegans* HiFi data mapped against the reference genome, as well as PacBio HiFi reads

457    generated from unamplified DNA extracted for a pool of worms of the same strain

458    (SRR22507561). Read depth, breadth, and mapping rate for *C. elegans* to the reference assembly

459    were calculated using SAMtools.

460    Orthology Inference and Supermatrix Construction

461    To investigate the phylogenetic position of Gastrotricha, orthology determination and

462    supermatrix assembly were performed following a modification of the pipeline of Kocot et al.

463    (2017b) implementing new versions of the programs used if applicable. Lophotrochozoan

464    genomes (including *L. squamata*) broadly spanning the diversity of the group were selected

465    along with five outgroup taxa resulting in 28 OTUs. Homologous sequences were identified

466    using OrthoFinder v.2.4.0 (Emms and Kelly 2019) and sequences shorter than 100 amino acids

467    were removed, keeping the longest sequence. Homogroups were retained if they included greater

468    than 75% of the total of all taxa, with those with less being removed. Fasta files were aligned

469    with MAFFT v.7.310 (Katoh 2002), suspected mistranslated regions removed with HmmCleaner

470    (Di Franco et al. 2019), and regions high in gaps or ambiguity were trimmed with BMGE

471    v.1.12.2 (Criscuolo and Gribaldo 2010). Maximum likelihood gene trees were generated using

472    IQ-Tree 2 (Minh et al. 2020) with the best-fitting model of amino acid substitution for each gene.

473    Subsequent fasta files and trees were used as input for PhyloPyPruner v.0.9.5

474    (https://pypi.org/project/phylopypruner) to retrieve a set of strict [1:1] orthologs with the settings

475    phylopypruner --min-taxa 15 --min-support 0.9 --mask pdist --trim-lb 3 --trim-divergent 0.75 --

476    min-pdist 0.01 --prune LS. The –min-taxa flag guaranteed removal of pruned alignments that did

477    not include at least 15 taxa from the original dataset.

478    Maximum Likelihood Estimation

479    The final supermatrix from PhyloPyPruner of 676,632 amino acid positions was used for

480    phylogeny reconstruction in IQ-Tree2. Using maximum likelihood with the simplest profile

481    mixture model available, LG+C20+F, a guide tree for PMSF (Wang et al. 2018) was produced.

482    The PMSF analysis fitted the C60 profile mixture model, LG+C60+F, to the data and

483    reconstructed the phylogeny assessing nodal support with 1000 rapid bootstraps.

484

485     **Availability of data and materials**

486     Genomic data generated for this study have been deposited in the NCBI SRA and Genome

487     databases under the BioProject accession number PRJNA1063365. Scripts used for assembly,

488     annotation, read coverage investigation, and visualization are available at

489     https://github.com/ngroberts/MDA_Hifi_Roberts_2023.

490

491     **Competing interests**

492     The authors declare that they have no competing interests.

493

494     **Funding**

495     KMK, NGR, funding from NSF (NSF DEB-1846174) and the University of Alabama the

496     College Academy of Research, Scholarship, and Creative Activity (CARSCA) grant. THS,

497     funding from the Research Council of Norway (FRIMEDBIO Project Number 300587).

498

499     **Authors' contributions**

500     Study design, NGR, KMK. Genome assembly, sequencing, and annotation, NGR, KMK, THS.

501     Statistical analyses and data visualization, NGR, MJG, THS. Manuscript preparation, NGR,

502     KMK.

503

504 **References:**

505 Alonge M, Lebeigle L, Kirsche M, Jenike K, Ou S, Aganezov S, Wang X, Lippman ZB, Schatz
506      MC, Soyk S. 2022. Automated assembly scaffolding using RagTag elevates a new
507      tomato system for high-throughput genome editing. *Genome Biology* 23:258.

508 Thalen, F. phylopypruner · PyPI. Available from: https://pypi.org/project/phylopypruner/

509 Benita Y, Oosting RS, Lok MC, Wise MJ, Humphery-Smith I. 2003. Regionalized GC content of
510      template DNA as a predictor of PCR success. *Nucleic Acids Research* 31:e99.

511 Bennett LW. 1979. Experimental Analysis of the Trophic Ecology of *Lepidodermella*
512      *squammata* (Gastrotricha: Chaetonotida) in Mixed Culture. *Transactions of the American*
513      *Microscopical Society* 98:254–260.

514 Birky CW Jr. 1996. Heterozygosity, Heteromorphy, and Phylogenetic Trees in Asexual
515      Eukaryotes. *Genetics* 144:427–437.

516 Bleidorn C. 2019. Recent progress in reconstructing lophotrochozoan (spiralian) phylogeny. *Org*
517      *Divers Evol* 19:557–566.

518 Borgström E, Paterlini M, Mold JE, Frisen J, Lundeberg J. 2017. Comparison of whole genome
519      amplification techniques for human single cell exome sequencing.Li Y, editor. *PLoS*
520      *ONE* 12:e0171566.

521 de Bourcy CFA, De Vlaminck I, Kanbar JN, Wang J, Gawad C, Quake SR. 2014. A Quantitative
522      Comparison of Single-Cell Whole Genome Amplification Methods.Wang K, editor.
523      *PLoS ONE* 9:e105585.

524 Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using DIAMOND.
525      *Nat Methods* 12:59–60.

526 C. elegans Sequencing Consortium. 1998. Genome sequence of the nematode *C. elegans*: a
527      platform for investigating biology. *Science* 282:2012–2018.

528 Chapman AR, He Z, Lu S, Yong J, Tan L, Tang F, Xie XS. 2015. Single Cell Transcriptome
529      Amplification with MALBAC.Lee JW, editor. *PLoS ONE* 10:e0120889.

530 Cohen-Fix O, Askjaer P. 2017. Cell Biology of the *Caenorhabditis elegans* Nucleus. *Genetics*
531      205:25–59.

532 Criscuolo A, Gribaldo S. 2010. BMGE (Block Mapping and Gathering with Entropy): a new
533      software for selection of phylogenetic informative regions from multiple sequence
534      alignments. *BMC Evol Biol* 10:210.

535  Cunha TJ, de Medeiros BAS, Lord A, Sørensen MV, Giribet G. 2023. Rampant loss of universal
536      metazoan genes revealed by a chromosome-level genome assembly of the parasitic
537      Nematomorpha. *Current Biology* 33:3514-3521.e4.

538  Dean FB, Hosono S, Fang L, Wu X, Faruqi AF, Bray-Ward P, Sun Z, Zong Q, Du Y, Du J, et al.
539      2002. Comprehensive human genome amplification using multiple displacement
540      amplification. *Proc. Natl. Acad. Sci. U.S.A.* 99:5261–5266.

541  Dean FB, Nelson JR, Giesler TL, Lasken RS. 2001. Rapid Amplification of Plasmid and Phage
542      DNA Using Phi29 DNA Polymerase and Multiply-Primed Rolling Circle Amplification.
543      *Genome Res.* 11:1095–1099.

544  Di Franco A, Poujol R, Baurain D, Philippe H. 2019. Evaluating the usefulness of alignment
545      filtering methods to reduce the impact of errors on evolutionary inferences. *BMC Evol*
546      *Biol* 19:21.

547  Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras
548      TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29:15–21.

549  Dujardin F. 1841. Histoire naturelle des zoophytes: Infusoires, comprenant la physiologie et la
550      classification de ces animaux, et la manière de les étudier à l'aide du microscope:
551      Ouvrage accompagné de planches. Roret

552  Edgecombe GD, Giribet G, Dunn CW, Hejnol A, Kristensen RM, Neves RC, Rouse GW,
553      Worsaae K, Sørensen MV. 2011. Higher-level metazoan relationships: recent progress
554      and remaining questions. *Org Divers Evol* 11:151–172.

555  Emms DM, Kelly S. 2019b. OrthoFinder: phylogenetic orthology inference for comparative
556      genomics. *Genome Biol* 20:238.

557  Evrony GD, Lee E, Mehta BK, Benjamini Y, Johnson RM, Cai X, Yang L, Haseley P, Lehmann
558      HS, Park PJ, et al. 2015. Cell Lineage Analysis in Human Brain Using Endogenous
559      Retroelements. *Neuron* 85:49–59.

560  Flot J-F, Hespeels B, Li X, Noel B, Arkhipova I, Danchin EGJ, Hejnol A, Henrissat B, Koszul R,
561      Aury J-M, et al. 2013. Genomic evidence for ameiotic evolution in the bdelloid rotifer
562      Adineta vaga. *Nature* 500:453–457.

563  Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF. 2020.
564      RepeatModeler2 for automated genomic discovery of transposable element families.
565      *Proceedings of the National Academy of Sciences* 117:9451–9457.

566  Fontaneto D, Kieneke A, Kristensen R, Maas A, Riemann O, Rothe B, Sørensen M, Sterrer W,
567      Taraschewski H, Warner L. 2015. Gastrotricha and Gnathifera.

568  Fu Y, Li C, Lu S, Zhou W, Tang F, Xie XS, Huang Y. 2015. Uniform and accurate single-cell
569      sequencing based on emulsion whole-genome amplification. *Proc Natl Acad Sci USA*
570      112:11923–11928.

Garmendia C, Bernad A, Esteban JA, Blanco L, Salas M. 1992. The bacteriophage phi 29 DNA polymerase, a proofreading enzyme. *Journal of Biological Chemistry* 267:2594–2599.

Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* 29:1072–1075.

Hejnol A. 2015. Gastrotricha. In: Wanninger A, editor. Evolutionary Developmental Biology of Invertebrates 2: Lophotrochozoa (Spiralia). Vienna: Springer. p. 13–19. Available from: https://doi.org/10.1007/978-3-7091-1871-9_2

Hoff K, Lomsadze A, Borodovsky M, Stanke M. 2019. Whole-Genome Annotation with BRAKER. *Methods Mol Biol* 1962:65–95.

Hon T, Mars K, Young G, Tsai Y-C, Karalius JW, Landolin JM, Maurer N, Kudrna D, Hardigan MA, Steiner CC, et al. 2020. Highly accurate long-read HiFi sequencing data for five complex genomes. *Sci Data* 7:399.

Hosono S, Faruqi AF, Dean FB, Du Y, Sun Z, Wu X, Du J, Kingsmore SF, Egholm M, Lasken RS. 2003. Unbiased Whole-Genome Amplification Directly From Clinical Samples. *Genome Res.* 13:954–964.

Huang S, Kang M, Xu A. 2017. HaploMerger2: rebuilding both haploid sub-assemblies from high-heterozygosity diploid genome assembly.Berger B, editor. *Bioinformatics* 33:2577–2579.

Katoh K. 2002. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Research* 30:3059–3066.

Kocot KM. 2016. On 20 years of Lophotrochozoa. *Org Divers Evol* 16:329–343.

Kocot KM, Struck TH, Merkel J, Waits DS, Todt C, Brannock PM, Weese DA, Cannon JT, Moroz LL, Lieb B, et al. 2017. Phylogenomics of Lophotrochozoa with Consideration of Systematic Error. *Systematic Biology* 66:256–282.

Laetsch DR, Blaxter ML. 2017. BlobTools: Interrogation of genome assemblies. *F1000Res* 6:1287.

Lasken RS. 2007. Single-cell genomic sequencing using Multiple Displacement Amplification. *Current Opinion in Microbiology* 10:510–516.

Lasken RS, Egholm M. 2003. Whole genome amplification: abundant supplies of DNA from precious samples or clinical specimens. *Trends in Biotechnology* 21:531–535.

Lasken RS, Stockwell TB. 2007. Mechanism of chimera formation during the Multiple Displacement Amplification reaction. *BMC Biotechnol* 7:19.

Laumer CE. 2018. Inferring Ancient Relationships with Genomic Data: A Commentary on Current Practices. *Integrative and Comparative Biology* 58:623–639.

605   Laumer, C. Picogram input multimodal sequencing (PiMmS) v1.
606          https://doi.org/10.17504/protocols.io.rm7vzywy5lx1/v1 (2022).

607   Laumer CE, Bekkouche N, Kerbl A, Goetz F, Neves RC, Sørensen MV, Kristensen RM, Hejnol
608          A, Dunn CW, Giribet G, et al. 2015. Spiralian Phylogeny Informs the Evolution of
609          Microscopic Lineages. *Current Biology* 25:2000–2006.

610   Laumer CE, Fernández R, Lemer S, Combosch D, Kocot KM, Riesgo A, Andrade SCS, Sterrer
611          W, Sørensen MV, Giribet G. 2019. Revisiting metazoan phylogeny with genomic
612          sampling of all phyla. *Proc. R. Soc. B.* 286:20190831.

613   Lee Y-C, Ke H-M, Liu Y-C, Lee H-H, Wang M-C, Tseng Y-C, Kikuchi T, Tsai IJ. 2023. Single-
614          worm long-read sequencing reveals genome diversity in free-living nematodes. *Nucleic*
615          *Acids Research* 51:8035–8047.

616   Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences.Birol I, editor.
617          *Bioinformatics* 34:3094–3100.

618   Li J, Lu N, Tao Y, Duan M, Qiao Y, Xu Y, Ge Q, Bi C, Fu J, Tu J, et al. 2018. Accurate and
619          sensitive single-cell-level detection of copy number variations by micro-channel multiple
620          displacement amplification (μcMDA). *Nanoscale* 10:17933–17941.

621   Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM. 2021. BUSCO Update: Novel and
622          Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for
623          Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. *Molecular Biology and*
624          *Evolution* 38:4647–4654.

625   Margulis L, Schwartz KV. 1998. Five Kingdoms: an illustrated guide to the Phyla of life on
626          earth. 3rd edition. Available from:
627          https://www.marinespecies.org/aphia.php?p=sourcedetails&id=3

628   Marlétaz F, Peijnenburg KTCA, Goto T, Satoh N, Rokhsar DS. 2019. A New Spiralian
629          Phylogeny Places the Enigmatic Arrow Worms among Gnathiferans. *Current Biology*
630          29:312-318.e3.

631   Martin S, Leggett RM. 2021. Alvis: a tool for contig and read ALignment VISualisation and
632          chimera detection. *BMC Bioinformatics* 22:124.

633   Martín-Durán JM, Vellutini BC, Marlétaz F, Cetrangolo V, Cvetesic N, Thiel D, Henriet S,
634          Grau-Bové X, Carrillo-Baltodano AM, Gu W, et al. 2021. Conservative route to genome
635          compaction in a miniature annelid. *Nat Ecol Evol* 5:231–242.

636   Mikhailov KV, Slyusarev GS, Nikitin MA, Logacheva MD, Penin AA, Aleoshin VV, Panchin
637          YV. 2016. The Genome of *Intoshia linei* Affirms Orthonectids as Highly Simplified
638          Spiralians. *Curr Biol* 26:1768–1774.

Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, Lanfear R. 2020. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Molecular Biology and Evolution* 37:1530–1534.

Nelson JR. 2014. Random-Primed, Phi29 DNA Polymerase-Based Whole Genome Amplification. *Current Protocols in Molecular Biology* [Internet] 105. Available from: https://onlinelibrary.wiley.com/doi/10.1002/0471142727.mb1513s105

Pinard R, de Winter A, Sarkis GJ, Gerstein MB, Tartaro KR, Plant RN, Egholm M, Rothberg JM, Leamon JH. 2006. Assessment of whole genome amplification-induced bias through high-throughput, massively parallel whole genome sequencing. *BMC Genomics* 7:216.

Pollard MO, Gurdasani D, Mentzer AJ, Porter T, Sandhu MS. 2018. Long reads: their purpose and place. *Human Molecular Genetics* 27:R234–R241.

Quinlan AR. 2014. BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Current Protocols in Bioinformatics* 47:11.12.1-11.12.34.

R Core Team. 2022. R: A Language and Environment for Statistical Computing. Vienna, Austria: R Foundation for Statistical Computing Available from: https://www.R-project.org/

Rayko M, Komissarov A, Kwan JC, Lim-Fong G, Rhodes AC, Kliver S, Kuchur P, O'Brien SJ, Lopez JV. 2020. Draft genome of *Bugula neritina*, a colonial animal packing powerful symbionts and potential medicines. *Sci Data* 7:356.

Rodrigue S, Malmstrom RR, Berlin AM, Birren BW, Henn MR, Chisholm SW. 2009. Whole Genome Amplification and De novo Assembly of Single Bacterial Cells. *PLOS ONE* 4:e6864.

Sahdev S, Saini S, Tiwari P, Saxena S, Singh Saini K. 2007. Amplification of GC-rich genes by following a combination strategy of primer design, enhancers and modified PCR cycle conditions. *Molecular and Cellular Probes* 21:303–307.
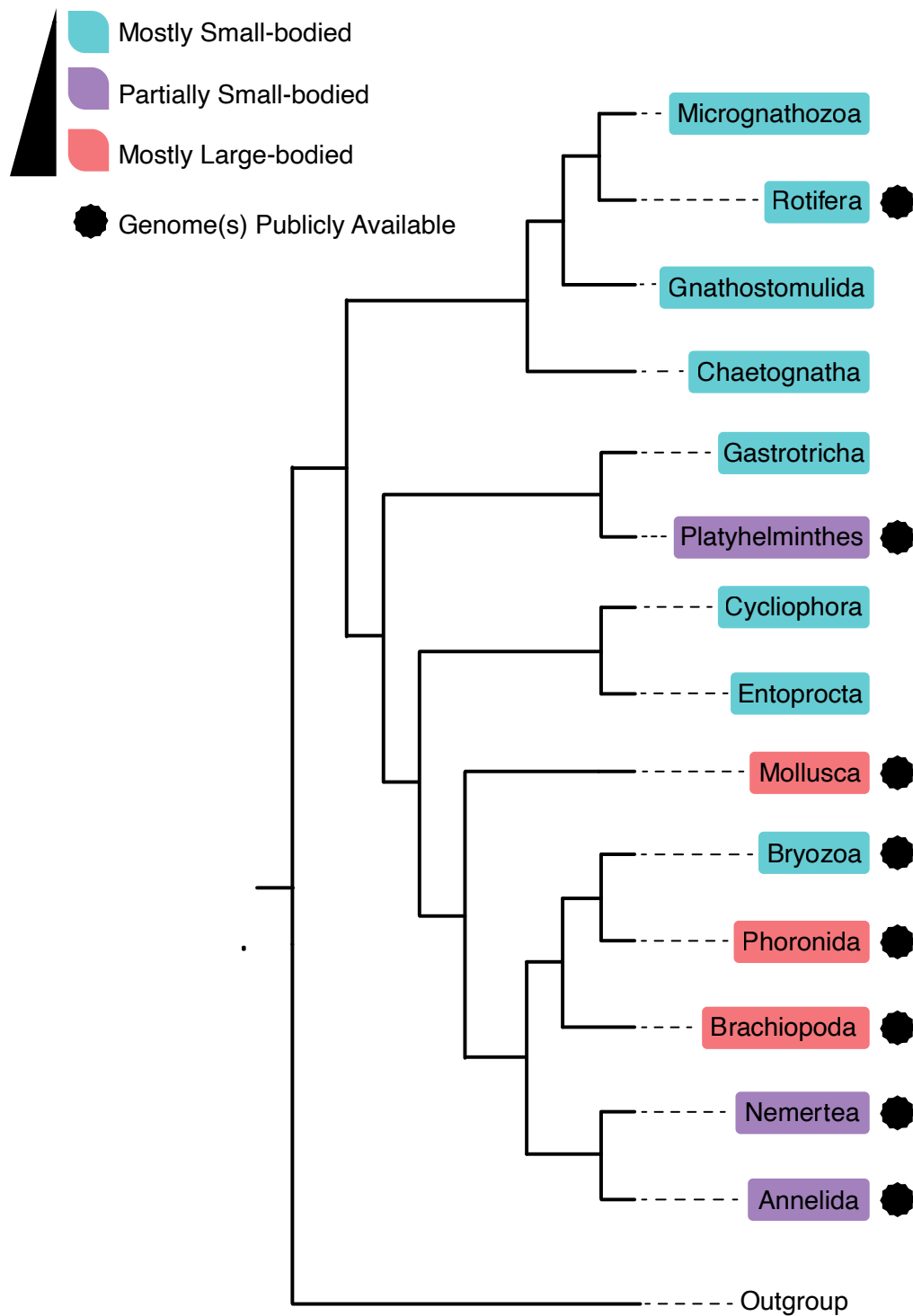
Schneider C, Woehle C, Greve C, D'Haese CA, Wolf M, Hiller M, Janke A, Bálint M, Huettel B. 2021. Two high-quality *de novo* genomes from single ethanol-preserved specimens of tiny metazoans (Collembola). *GigaScience* 10:giab035.

Seo H-C, Kube M, Edvardsen RB, Jensen MF, Beck A, Spriet E, Gorsky G, Thompson EM, Lehrach H, Reinhardt R, et al. 2001. Miniature Genome in the Marine Chordate *Oikopleura dioica*. *Science* 294:2506–2506.

Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31:3210–3212.

673   Simion P, Belkhir K, François C, Veyssier J, Rink JC, Manuel M, Philippe H, Telford MJ. 2018.
674         A software tool 'CroCo' detects pervasive cross-species contamination in next generation
675         sequencing data. *BMC Biol* 16:28.

676   Smit A, Hubley R, Green P. 2015. RepeatMasker Open-4.0. 2013–2015.

677   Stevens L, Martínez-Ugalde I, King E, Wagah M, Absolon D, Bancroft R, Gonzalez de la Rosa
678         P, Hall JL, Kieninger M, Kloch A, et al. 2023. Ancient diversity in host-parasite
679         interaction genes in a model parasitic nematode. *Nat Commun* 14:7776.

680   Struck TH, Wey-Fabrizius AR, Golombek A, Hering L, Weigert A, Bleidorn C, Klebow S,
681         Iakovenko N, Hausdorf B, Petersen M, et al. 2014. Platyzoan Paraphyly Based on
682         Phylogenomic Data Supports a Noncoelomate Ancestry of Spiralia. *Molecular Biology*
683         *and Evolution* 31:1833–1849.

684   Telenius H, Carter NP, Bebb CE, Nordenskjo¨ld M, Ponder BAJ, Tunnacliffe A. 1992.
685         Degenerate oligonucleotide-primed PCR: General amplification of target DNA by a
686         single degenerate primer. *Genomics* 13:718–725.

687   Wang H-C, Minh BQ, Susko E, Roger AJ. 2018. Modeling Site Heterogeneity with Posterior
688         Mean Site Frequency Profiles Accelerates Accurate Phylogenomic Estimation.
689         *Systematic Biology* 67:216–235.

690   Worsaae K, Vinther J, Sørensen MV. 2023. Evolution of Bilateria from a Meiofauna
691         Perspective—Miniaturization in the Focus. In: New Horizons in Meiobenthos Research.
692         Springer, Cham. p. 1–31. Available from: https://link.springer.com/chapter/10.1007/978-
693         3-031-21622-0_1

694   Zhou W, Emery SB, Flasch DA, Wang Y, Kwan KY, Kidd JM, Moran JV, Mills RE. 2020.
695         Identification and characterization of occult human-specific LINE-1 insertions using
696         long-read sequencing technology. *Nucleic Acids Research* 48:1146–1163.

697   Zong C, Lu S, Chapman AR, Xie XS. 2012. Genome-Wide Detection of Single Nucleotide and
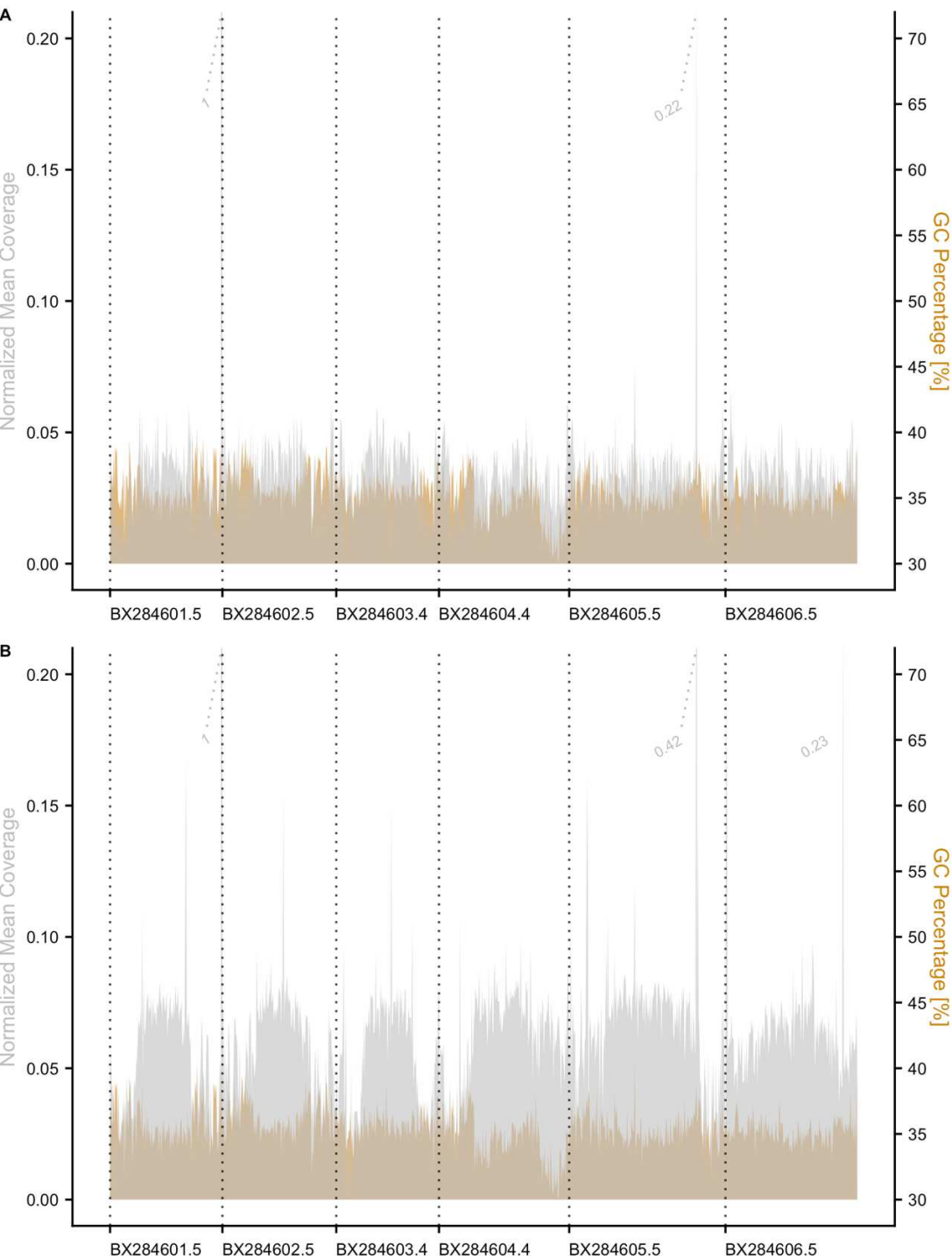698         Copy Number Variations of a Single Human Cell. *Science* 338:1622–1626.

699

700

701

702

703

704

**Figure 1.** Phyla within Lophotrochozoa that are exclusively, partially, or mostly small-bodied.

Those with one or more publicly available genomes are indicated. Tree modified from Laumer et

al. (2019).

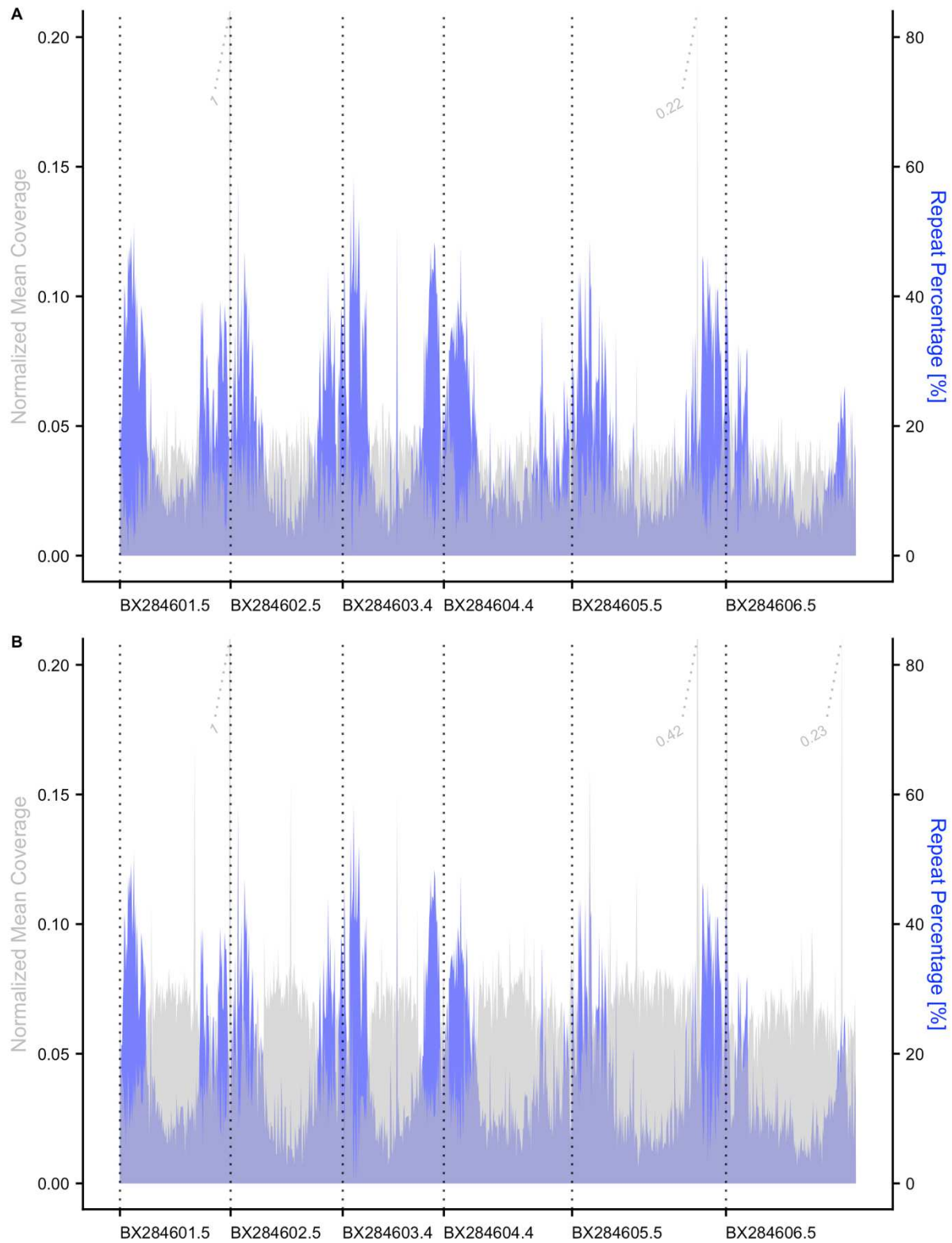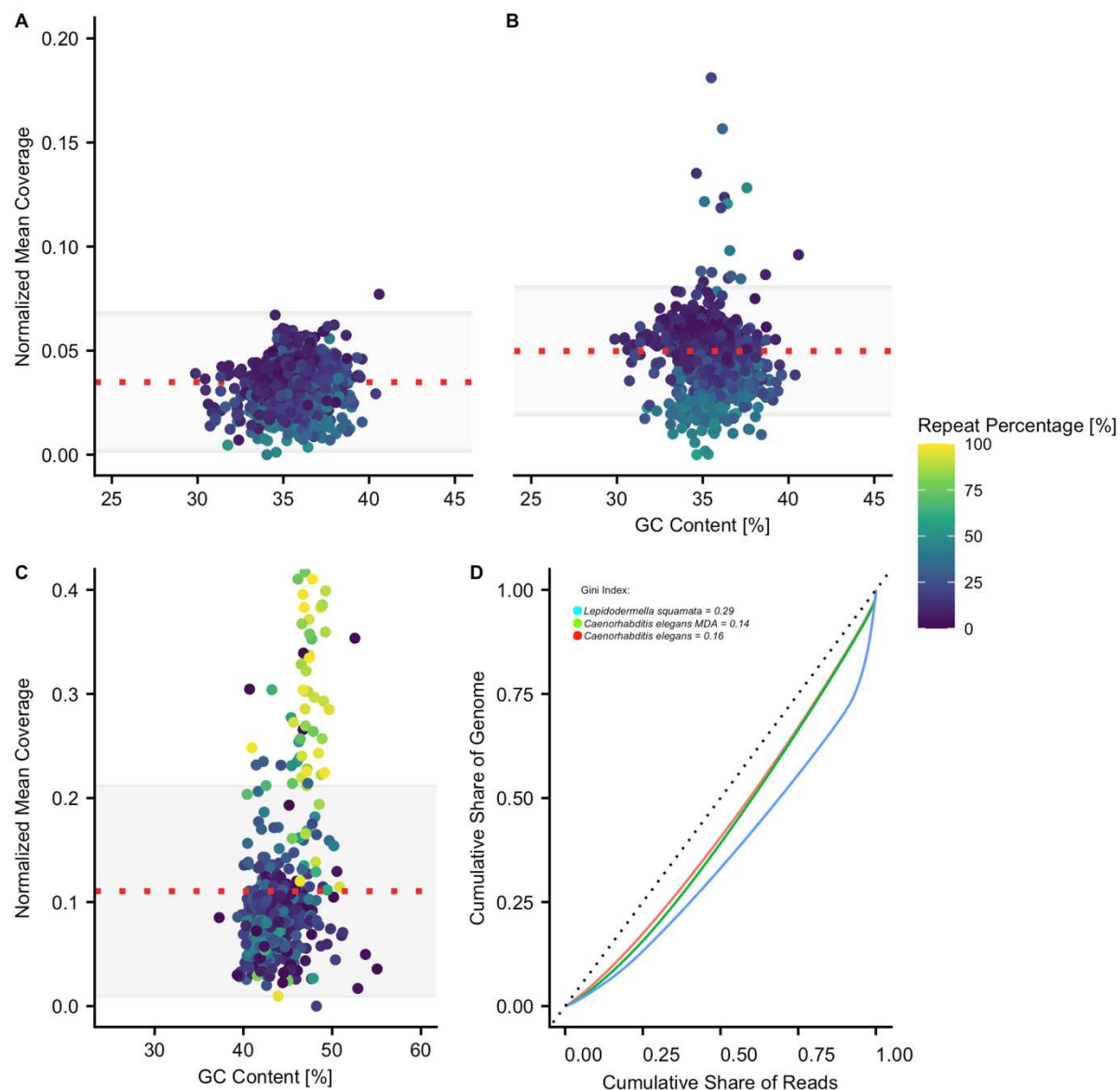*Caenorhabditis elegans*



708

709

710

711     **Figure 2.** Mean coverage and GC content across 100 Kbp blocks of the *C. elegans* reference

712     genome. A. Coverage of PacBio HiFi reads from MDA DNA aligned to the reference *C. elegans*

713     genome [Gray] alongside GC percentage [Orange]. B. Coverage of unamplified PacBio HiFi

714     reads aligned to the reference *C. elegans* genome [Gray] alongside GC percentage [Orange].
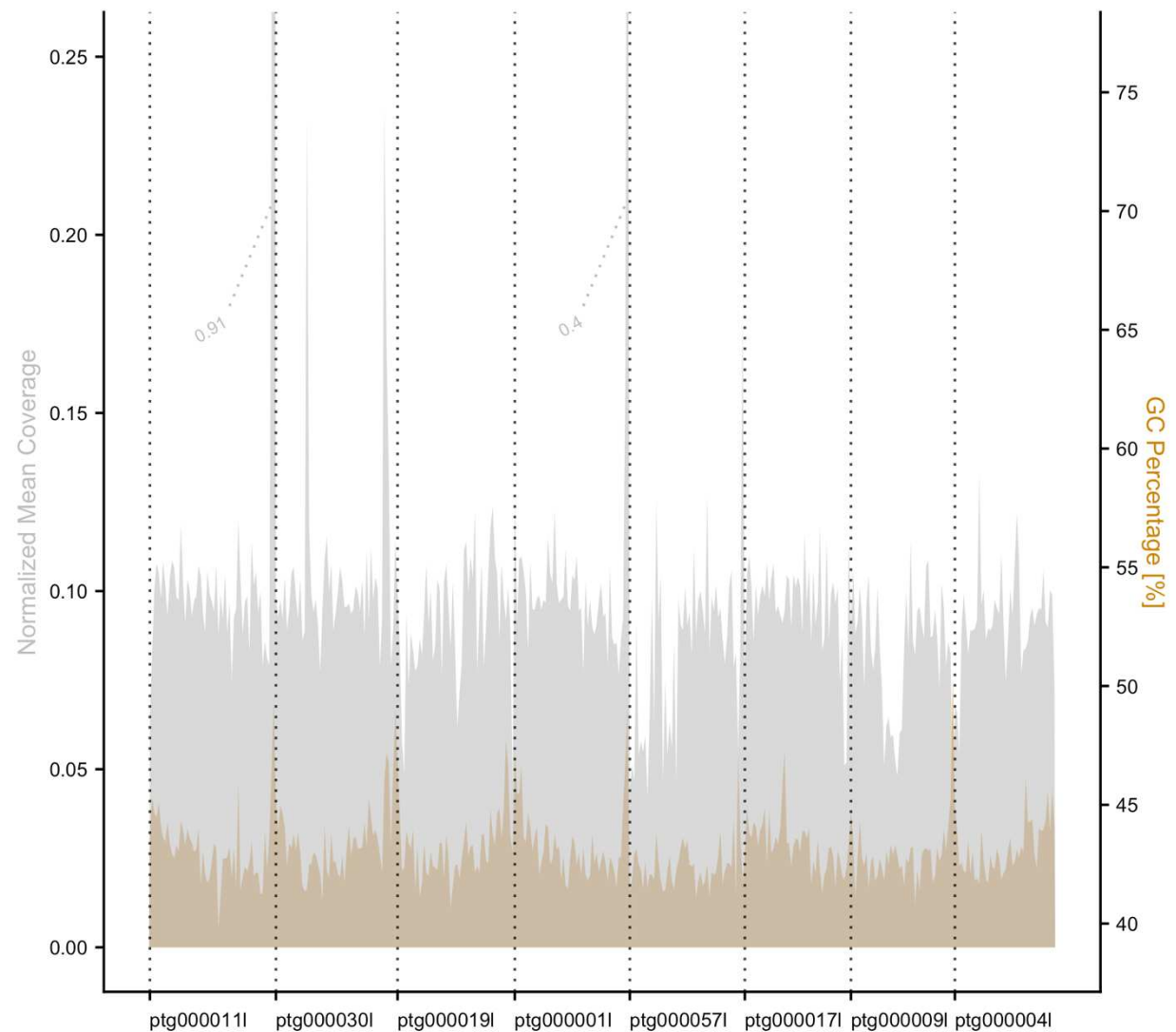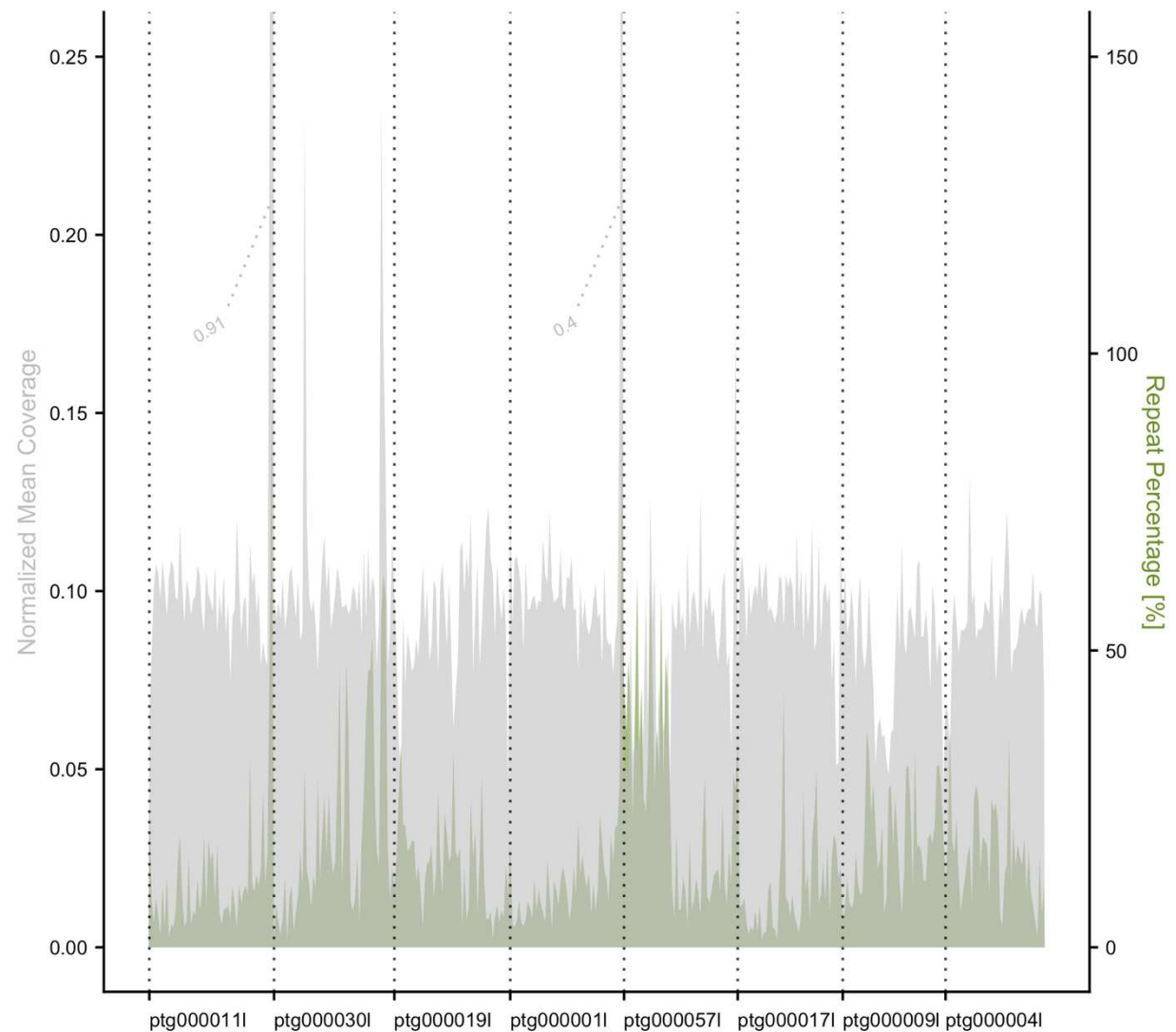
*Caenorhabditis elegans*



715

716    **Figure 3.** Mean coverage and repeat content across 100 Kbp blocks of the *C. elegans* reference

717    genome. A. Coverage of PacBio HiFi reads from MDA DNA aligned to the reference *C. elegans*

718    genome [Gray] alongside repeat content [Blue]. B. Coverage PacBio HiFi reads from

719    unamplified DNA from a pool of worms aligned to the reference genome [Gray] alongside repeat

720    content [Blue].
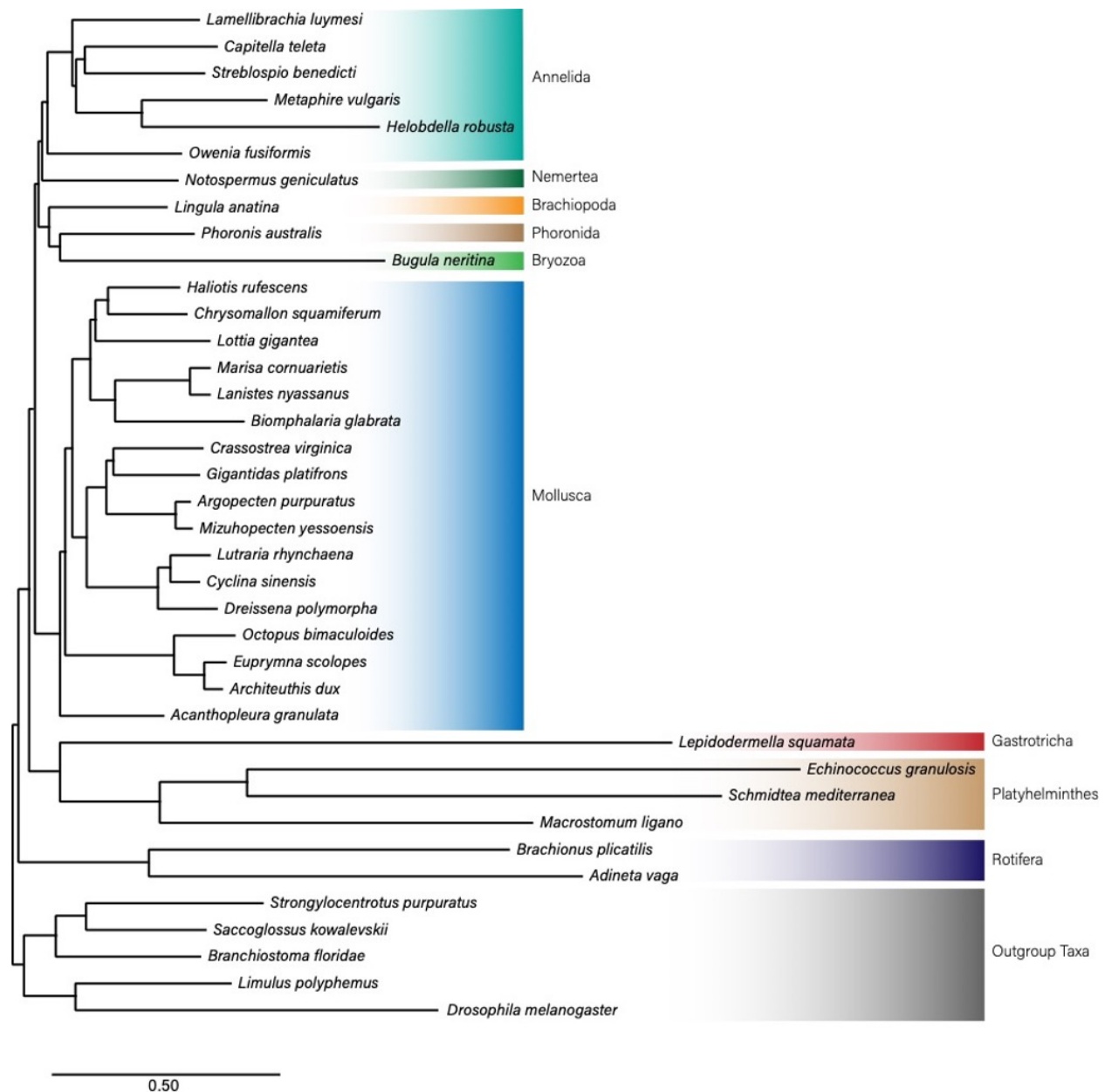


721

722

723    **Figure 4.** Normalized coverage, GC content, and repeat percentage of each genome assembly in

724    100 Kbp blocks. A. Dotplot showing coverage of MDA amplified reads with respect to GC

725    content for 100 Kbp blocks of the *C. elegans* reference genome. B. Dotplot showing coverage of

726    non-amplified reads from a pool of worms with respect to GC content for 100 Kbp blocks of the

727    *C. elegans* reference genome. C. Dotplot showing coverage with respect to GC content for

728    contigs in the *L. squamata* genome assembly based on MDA DNA. A-C. Repeat content of each

729    contig is indicated according to the key at the top right of the figure. D. Lorenz curves and Gini

730    values for each assembly indicating coverage uniformity. The diagonal line represents perfect

731    read coverage uniformity.

**Figure 5.** Coverage and GC content across 100 Kbp blocks of the eight largest contigs of the *L. squamata* genome. Coverage of PacBio HiFi reads from MDA DNA aligned to the 8 largest contigs of the *de novo L. squamata* genome assembly [Gray] alongside GC percentage [Orange].
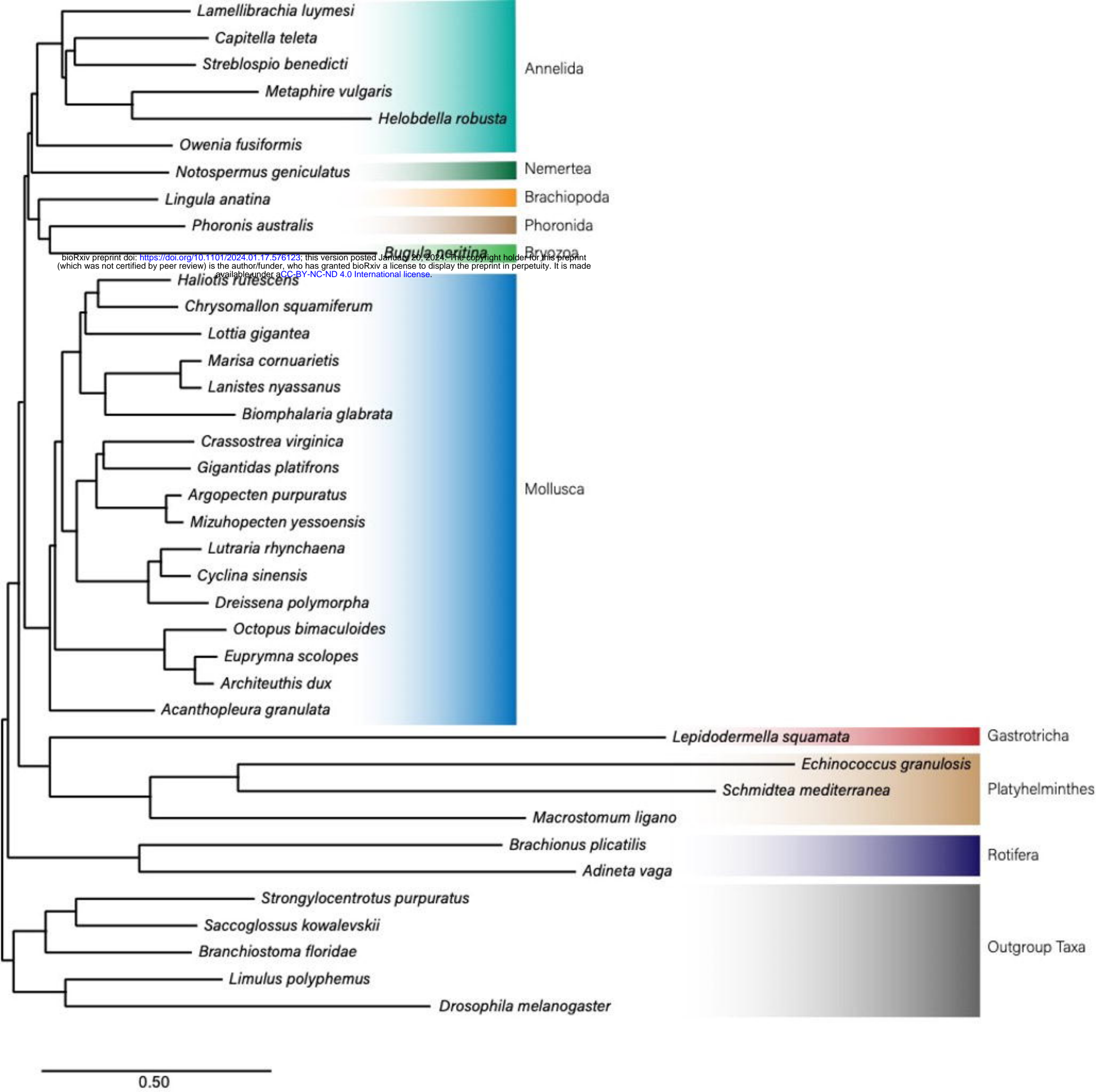
**Figure 6.** Coverage and repeat across 100 Kbp blocks of the eight largest contigs of the *L. squamata* genome. Coverage of PacBio HiFi reads from MDA DNA aligned to the 8 largest contigs of the *de novo L. squamata* genome assembly [Gray] alongside repeat percentage [Green]
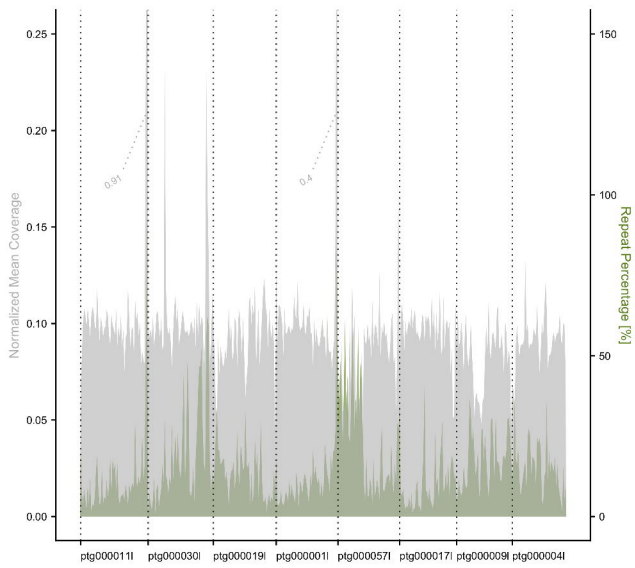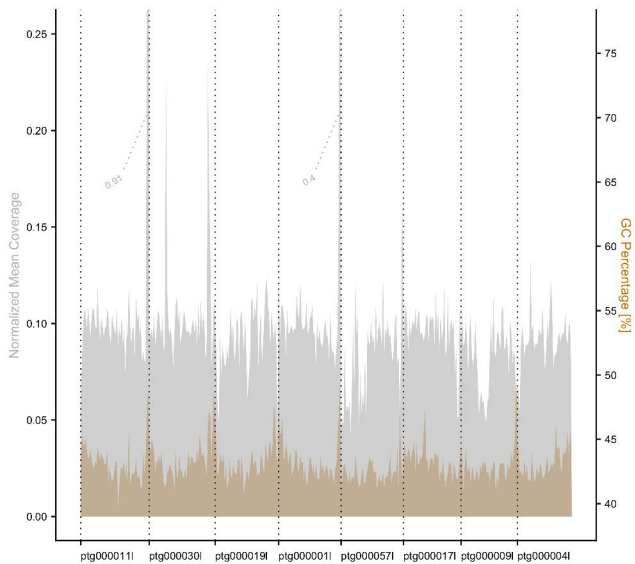
747

**Figure 7.** Evolutionary relationships of lophotrochozoan taxa with genomic data sequenced to

date placing Gastrotricha as the sister taxon of Platyhelminthes (Rouphozoa). IQ-TREE 2

analysis of 676,632 amino acid positions with 18.68% missing data using the PMSF

(LG+C60+F) with LG+C20+F as the guide tree. All nodes have 100% bootstrap support.

*Lamellibrachia luymesi*
*Capitella teleta*
*Streblospio benedicti*
*Metaphire vulgaris*
*Helobdella robusta*
*Owenia fusiformis*                    Annelida
*Notospermus geniculatus*              Nemertea
*Lingula anatina*                      Brachiopoda
*Phoronis australis*                   Phoronida
**Bugula neritina**                    Bryozoa

*Haliotis rufescens*
*Chrysomallon squamiferum*
*Lottia gigantea*
*Marisa cornuarietis*
*Lanistes nyassanus*
*Biomphalaria glabrata*
*Crassostrea virginica*
*Gigantidas platifrons*
*Argopecten purpuratus*                Mollusca
*Mizuhopecten yessoensis*
*Lutraria rhynchaena*
*Cyclina sinensis*
*Dreissena polymorpha*
*Octopus bimaculoides*
*Euprymna scolopes*
*Architeuthis dux*
*Acanthopleura granulata*

*Lepidodermella squamata*              Gastrotricha
*Echinococcus granulosis*
*Schmidtea mediterranea*               Platyhelminthes
*Macrostomum ligano*
*Brachionus plicatilis*
*Adineta vaga*                         Rotifera
*Strongylocentrotus purpuratus*
*Saccoglossus kowalevskii*
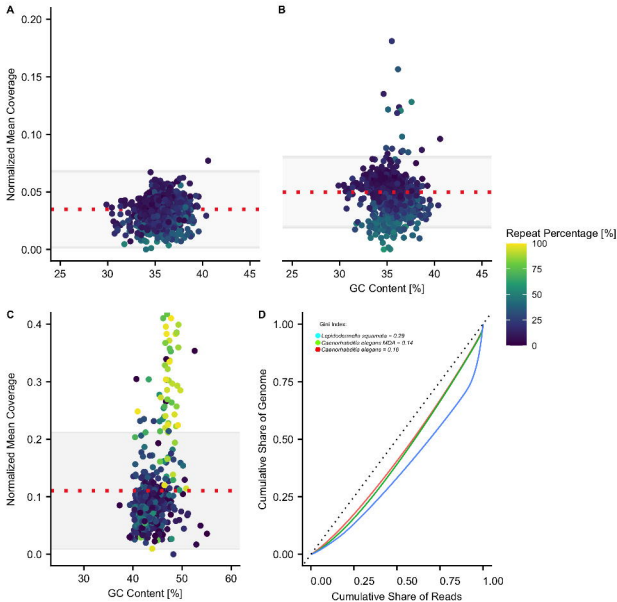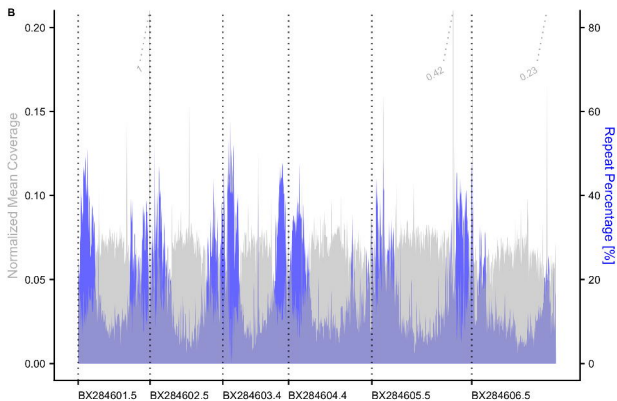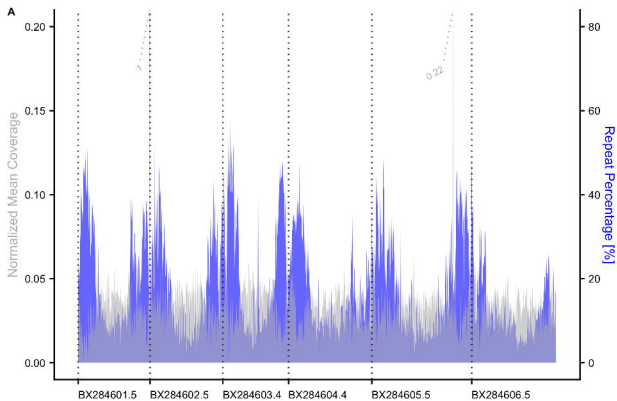*Branchiostoma floridae*               Outgroup Taxa
*Limulus polyphemus*
*Drosophila melanogaster*

0.50

*Caenorhabditis elegans*

*Caenorhabditis elegans*