1  **TITLE:** Single-molecule sequencing of long DNA molecules allows high contiguity *de*

2  *novo* genome assembly for the fungus fly, *Sciara coprophila*

3

4  **AUTHORS:** John M. Urban[1,2], Michael S. Foulk[1,3], Jacob E. Bliss[1], C. Michelle Coleman[4], Nanyan

5  Lu[4], Reza Mazloom[4], Susan J. Brown[4], Allan C. Spradling[2] and Susan A. Gerbi[1]

6

7  [1]Brown University Division of Biology and Medicine, Department of Molecular Biology, Cell

8  Biology and Biochemistry, Providence, Rhode Island 02912, USA

9  [2]Carnegie Institution for Science, Department of Embryology, 3520 San Martin Drive, Baltimore,

10  Maryland 21218, USA

11  [3]Present Address: Mercyhurst University, Department of Biology, Erie, PA 16546, USA

12  [4]Kansas State University Division of Biology, KSU Bioinformatics Center, Ackert Hall, Manhattan,

13  Kansas 66502, USA

14

15  CORRESPONDING AUTHORS:

16      Susan A. Gerbi

17      Brown University Division of Biology and Medicine

18      Sidney Frank Hall for Life Sciences

19      185 Meeting Street

20      Providence, RI 02912 USA

21      TEL: 401-863-2359/ FAX: 401-863-1201

22      E-mail: Susan_Gerbi@Brown.edu

23

24      John M. Urban

25      Carnegie Institution for Science

26      Department of Embryology

27      3520 San Martin Drive

28      Baltimore, Maryland 21218 USA

29      TEL: 410-246-3001

30      E-mail: jurban@carnegiescience.edu

31

32      **RUNNING TITLE:** The fungus fly genome

33

34      **KEYWORDS:** genome assembly, single molecule sequencing, long reads, optical maps,

35      nanopore sequencing, DNA modifications, non-model organism, emerging model system, insect

36      genomes, fungus fly *Sciara* (*Bradysia*) coprophila.

37

38      **MANUSCRIPT TYPE:** RESEARCH

39

40      **DEDICATION:** Dedicated to Ellen M. Rasch (January 31, 1927 – July 31, 2016), a leader in

41      Feulgen-DNA cytophotometry who quantified the genome size of *Sciara coprophila*. Active in the

42      American Society for Cell Biology (Council) and the Histochemical Society (Council, Secretary

43      and Treasurer), she was a Fellow of the American Association for Advancement of Science and

44      the Royal Microscopic Society.

**ABSTRACT**

45

46      The lower Dipteran fungus fly, *Sciara coprophila*, has many unique biological features. For

47    example, *Sciara* undergoes paternal chromosome elimination and maternal X chromosome

48    nondisjunction during spermatogenesis, paternal X elimination during embryogenesis,

49    intrachromosomal DNA amplification of DNA puff loci during larval development, and germline-

50    limited chromosome elimination from all somatic cells. Paternal chromosome elimination in *Sciara*

51    was the first observation of imprinting, though the mechanism remains a mystery. Here, we

52    present the first draft genome sequence for *Sciara coprophila* to take a large step forward in aiding

53    these studies. We approached assembling the *Sciara* genome using multiple sequencing

54    technologies: PacBio, Oxford Nanopore MinION, and Illumina. To find an optimal assembly using

55    these datasets, we generated 44 Illumina assemblies using 7 short-read assemblers and 50 long-

56    read assemblies of PacBio and MinION sequence data using 6 long-read assemblers. We ranked

57    assemblies using a battery of reference-free metrics, and scaffolded a subset of the highest-

58    ranking assemblies using BioNano Genomics optical maps. RNA-seq datasets from multiple life

59    stages and both sexes facilitated genome annotation. Moreover, we anchored nearly half of the

60    *Sciara* genome sequence into chromosomes. Finally, we used the signal level of both the PacBio

61    and Oxford Nanopore data to explore the presence or absence of DNA modifications in the *Sciara*

62    genome since DNA modifications may play a role in imprinting in *Sciara*, as they do in mammals.

63    These data serve as the foundation for future research by the growing community studying the

64    unique features of this emerging model system.

## **INTRODUCTION**

The fungus gnat, *Sciara coprophila* (also known as *Bradysia coprophila*), is a Dipteran fly that is both an old and emerging model system rich with opportunities for studying fundamental biology, especially chromosomal biology due to its dynamic genome. In contrast to the rule that the amount of nuclear DNA is constant in all cells of an organism (Boivin et al. 1948), the nuclear DNA in *Sciara* cells exhibits copy number differences at the levels of loci, chromosomes, and the genome. Genomic copy numbers vary from canonical haploid and diploid tissues to the endocycling larval salivary glands that result in cells with over 8000 copies of each chromosome held closely together to form giant polytene chromosomes (Rasch 1970b). Locus-specific copy number regulation occurs at the "DNA puff" loci in polytene chromosomes where site-specific re-replication results in intrachromosomal DNA amplification (Rasch 1970a; Gerbi et al. 2002). Whole chromosome copy number gains and losses are seen in spermatogenesis, fertilization, in somatic cells of early embryos, and in the germ-line during development (Gerbi 1986).

The chromosome cycle of *Sciara* gives rise to numerous research opportunities not found in *Drosophila*, the standard Dipteran model organism. In *Sciara*, there are "L" chromosomes limited to the germ-line of both sexes (Gerbi 1986). Whereas oogenesis has orthodox chromosome movements, they are unusual in spermatogenesis leading to sperm cells that are haploid for each autosome, diploid for the X, and variable for the L with 0-4 copies. X diploidy in sperm is due to developmentally programmed X chromosome nondisjunction in male meiosis (Gerbi 1986). Fertilization ultimately produces zygotes and early embryos that are temporarily triploid for the X chromosome, and variable for the L. The fates of the X and L chromosomes in early embryonic nuclei are subsequently determined by whether a cell is somatic or germline, and by whether it is male or female. All L chromosomes are eliminated from somatic cell nuclei in early embryos. As part of the sex determination pathway, X diploidy is restored in female somatic cells (XX) by the elimination of one X, but the elimination of two X chromosomes in male somatic cells

4

91    (XO) leads to X haploidy (Gerbi 1986). Diploidy for the X and L is restored in the germline through

92    elimination events later in development (Gerbi 1986).

93

94        The X chromosomes eliminated during early embryo development are always paternally

95    derived. Moreover, all paternally derived chromosomes, except L, are eliminated in the first

96    meiotic division of spermatogenesis in the only known case of a naturally occurring monopolar

97    spindle (Gerbi 1986). The ability to differentiate between the maternal and paternal chromosomes

98    gave rise to the term "imprinting" (Crouse 1960) and was the first description of this phenomenon

99    in any system. L chromosomes apparently escape imprinting in *Sciara* as maternal and paternal

100   copies are both eliminated from all nuclei destined to become somatic cells (Crouse et al. 1971),

101   and they are not eliminated with the paternal cohort during male meiosis. The mechanism for

102   imprinting in *Sciara* remains unknown. It is of interest to learn if DNA modifications occur in the

103   *Sciara* genome, since imprinting in mammals utilizes DNA methylation (Li et al. 1993).

104

105       This black fungus gnat and its unusual chromosomal features are part of one of the most

106   interesting yet little-studied groups of Dipteran flies, the suborder Nematocera. The group of

107   Nematocerans contains agricultural pests as well as disease vectors, such as mosquitoes

108   (Matthews et al 2018). Nematocera diverged from higher Dipteran flies, the suborder Brachycera

109   that includes the fruit fly *Drosophila melanogaster*, ~200 million years ago (Wiegmann et al. 2011).

110   *Bradysia (Sciara) coprophila* is classified as part of the infraorder Bibionomorpha in the Sciaroidea

111   super family, which also comprises the family Cecidomyiidae (gall midges) and the Hessian fly in

112   particular, a notorious wheat pest (Stuart et al 2012). Sciarid flies also include the Mycetophilidae,

113   a fungus gnat family where members have been shown to withstand freezing and thawing (Sformo

114   et al 2009). Indeed, we also have unpublished observations that *Sciara coprophila* embryos and

115   larvae can be stored in the cold from a few months to over a year in a diapause-like state before

116   returning to room temperature and resuming development. Despite flies making up at least 10%

117  of all metazoan diversity, there are only 157 Dipteran genomes described

118  (i5k.github.io/arthropod_genomes_at_ncbi), most of which are highly fragmented assemblies,

119  and the majority of which are from the higher Dipteran order and limited to only two suborders

120  therein (Muscomorpha, Stratiomyomorpha). Thus, there is a real need for high quality genomes

121  across the Dipteran tree, and particularly for the lower Dipteran suborder that includes *Sciara*.

122

123       The complete *Sciara* genome comprises three autosomes (chromosomes II, III and IV),

124  an X chromosome, and the germ-line limited L chromosome (Figure 1; Gerbi 1986). L

125  chromosomes are eliminated from nuclei destined to become somatic cells in the 5th or 6th

126  nuclear division, ~3 hours after egg deposition (Gerbi 1986). *Sciara* lacks a Y chromosome, and

127  sex is determined by whether or not the mother carries a variant of the X, called X', that has a

128  long paracentric inversion. Females that are XX have only sons, whereas X'X females have only

129  daughters. The XX or X'X genotype of adult females can be determined by phenotypic wing

130  markers (Figure 1). The *Sciara* genome has ~38% GC content (Gerbi 1971) and is ~280 Mb in

131  somatic cells and ~363 Mb in germ cells that contain L chromosomes (Rasch 2006)

132  (Supplemental Table S1A-D).

133

134       There are many ways to assemble a genome, but no universal recipe of sequencing

135  technologies, pre-assembly practices (e.g. quality filtering, error correction), assembly algorithms,

136  parameter tuning, and post-assembly steps exists that guarantees the best assembly for a given

137  genome. Therefore, to maximize contiguity and quality, we sequenced the *Sciara* genome with

138  multiple technologies, including 100 bp Illumina paired-end reads, long reads from Pacific

139  Biosciences (PacBio) (Eid et al. 2009) and the Oxford Nanopore MinION (Ip et al. 2015), and

140  generated optical maps from the BioNano Genomics Irys platform (Lam et al. 2012). We produced

141  assemblies using combinations of these technologies with multiple algorithms and ranked each

142  using a battery of reference-free metrics. Highly contiguous assemblies that were most complete
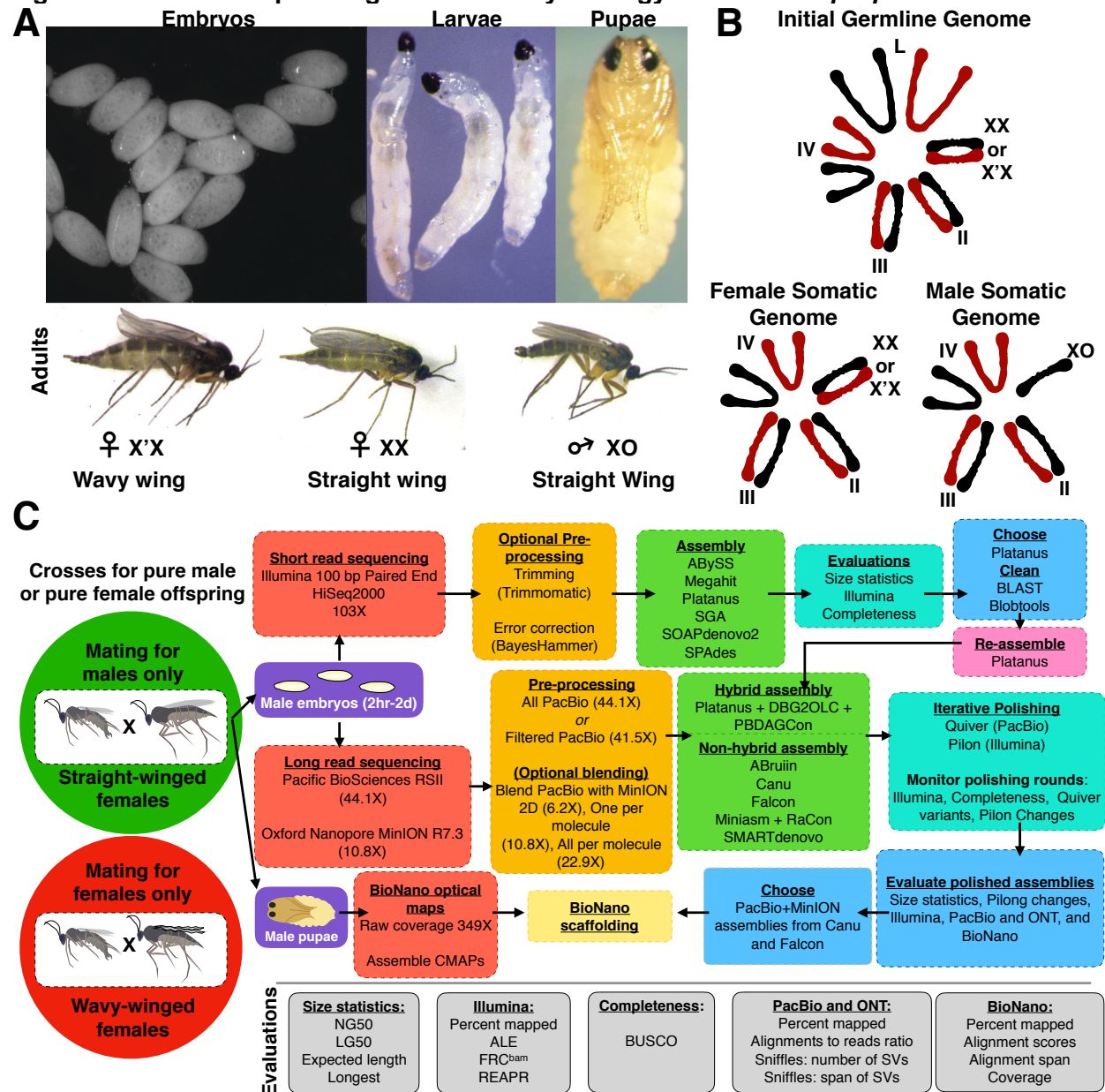
143    in expected gene content and which were judged to be most consistent with our Illumina, PacBio,

144    MinION, BioNano, and RNA-seq datasets were identified. These evaluations allowed us to

145    monitor steps (e.g. polishing), to choose a few assemblies for BioNano scaffolding, and to make

146    a final selection for the *Sciara* draft genome.

147

148        We report here the first draft genome assembly for *Sciara coprophila*, and its

149    accompanying gene and repeat annotations. The *Sciara* genome sequence will be a valuable

150    resource for future comparative genomics analyses, as one of the highest-quality Nematoceran

151    genome sequences available, as the only sequenced member of the Sciaridae family, and due to

152    its phylogenetic position at the gateway between lower and higher Dipterans. More than half of

153    the *Sciara* genome is contained on contigs $\geq$1.9 Mb and scaffolds $\geq$6.8 Mb. This exceeds the

154    contiguity of ~90% of all Dipteran genome assemblies

155    (i5k.github.io/arthropod_genomes_at_ncbi). More specifically, the contig sizes in this release of

156    the *Sciara* genome are longer than 42 of the 43 Nematoceran genome assemblies described,

157    only outshined by the assembly for the mosquito, *Aedes aegypti* (Matthews et al 2018). The

158    megabase-scale contigs and scaffolds will aid in efforts to improve the contiguity of more

159    fragmented assemblies of related species by synteny. The genome annotation contains >97% of

160    expected gene content. Up to 49% of the *Sciara* genome sequence was anchored into specific

161    loci of chromosomes X, II, III, and IV; and 100% was classified as either X or autosomal, allowing

162    an analysis of dosage compensation of the single male X. A *Rickettsia* genome was co-

163    assembled with the *Sciara* genome, suggesting it may be an endosymbiont. The signal data from

164    both PacBio and MinION both suggest the presence of DNA modifications in the *Sciara* genome.

165    Finally, candidate L sequences were briefly explored. Sequencing, assembly, and annotation of

166    the *Sciara* genome reported here serves as the foundation for future studies of the many unique

167    features of this emerging model organism.

168 **Figure 1: Genome sequencing and assembly strategy for *Sciara coprophila***



169
170 **Figure 1: Genome sequencing and assembly strategy for *Sciara coprophila***

171
172 **(A)** Images of different lifecycle stages of *Sciara coprophila*: embryos, larvae, pupae, and adults.
173 The adult figures show a male and the two different types of females that can be distinguished
174 based on the Wavy wing phenotype that marks the X' chromosome. **(B)** Examples of different
175 chromosome compositions in *Sciara* cells. We focused on the male somatic genome. Red
176 chromosomes are paternal, black are maternal. **(C)** The genome assembly and evaluation
177 workflow up until BioNano scaffolding. The workflow begins by highlighting that crosses can be
178 conducted to generate only male (green) or only female (red) offspring using the Wavy wing
179 phenotypic marker. We used only matings for males to obtain genomic DNA for sequencing,
180 illustrated by the arrows from the green circle that point to subsequent steps in the pipeline. Both
181 male (green) and female (red) offspring were used for transcriptomes.
182

183      **RESULTS**

184    **Data collection**

185        Using wing phenotypic markers, XX *Sciara* adult females were crossed with XO males to

186    produce only male progeny (Figure 1). DNA isolated from purely male embryos was used for

187    sequencing (Illumina, PacBio, MinION), thereby avoiding assembly complications from the

188    heteromorphic X' chromosome found in female-producing females (Figure 1B), as well as

189    minimizing possible complications from later life stages due to polytenization and contamination

190    from the gut microbiome. Moreover, although early embryos were included to potentially capture

191    sequences from chromosome L, the somatic genome is over-represented in these samples and

192    we do not expect L sequences from the germline genome to be well-represented. Separate

193    preparations of male embryo genomic DNA were made for 100 bp paired-end Illumina, and long-

194    read PacBio and Oxford Nanopore MinION sequencing resulting in 103X, 50-55X, and 10-11X

195    coverage, respectively (Table 1). We used male pupae to collect nearly 350X coverage from a

196    third single molecule technology: optical maps from the BioNano Genomics (BNG) Irys (Lam et

197    al. 2012) (Table 1). Finally, to facilitate gene annotation, we acquired sex- and stage-specific 100

198    bp paired-end RNA-seq datasets from whole embryos, larvae, pupae, and adults using the

199    appropriate crosses for only males (XX x XO) or only females (X'X x XO) (Supplemental Table

200    S2).

201 **Table 1: Genome sequencing datasets for *Sciara coprophila***

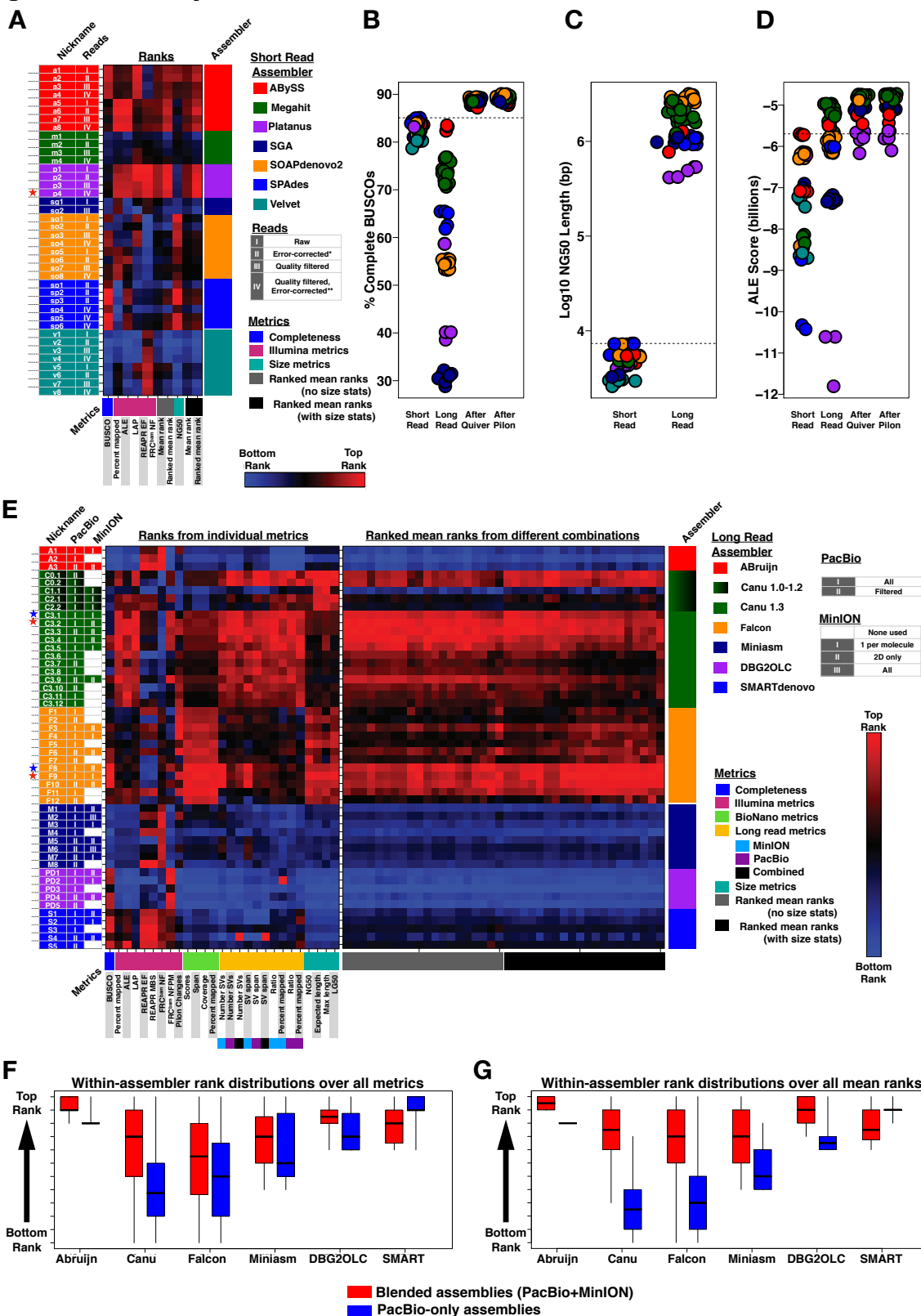| | Illumina HiSeq 2000 | PacBio RSII | Oxford Nanopore MinION MkI | BioNano Genomics Irys |
|---|---|---|---|---|
| **Source** | Male Embryos | Male Embryos | Male Embryos* | Male pupae |
| **Library** | Paired-End* | SMRTBell | MAP002-006 (2D) | IrysPrep |
| **Details** | - | P5-C3 | Pores R7.3-R7.3 70bps 6mer | BssSI |
| **Read Length N50 (kb)** | 0.1 | 9.681 | 9.934 | 132.613 |
| **Mean Read Length (kb)** | 0.1 | 6.607 | 5.883 | 62.531 |
| **Count** | 301,513,554 | 1,949,427 | 532,714 | 1,628,681 |
| **Span (Gb)** | 30.15 | 12.88 | 3.15 | 101.84 |
| **Coverage >0 kb** | 103.26 | 44.11 | 10.77 | 348.78 |
| **>20 kb** | 0 | 1.28 | 2.91 | 330.22 |
| **>30 kb** | 0 | 0.01 | 1.72 | 323.31 |
| **>50 kb** | 0 | 0 | 0.71 | 303.02 |
| **>100 kb** | 0 | 0 | 0.28 | 226.1 |
| **>150 kb** | 0 | 0 | 0.2 | 148.5 |

202    *A minority of the MinION data came from male adults (see Methods).

10

203 **<u>Short-read assemblies</u>**

204       Using the Illumina dataset, both with and without quality filtering and/or error-correction

205 steps, we generated 44 assemblies using 7 popular short-read genome assemblers (Figures 1C

206 and 2A). The assemblies ranged from ~226-348 Mb in size (Supplemental Table S3), with a mean

207 assembly size of ~280 Mb, exactly the expected somatic genome size. Evaluating these

208 assemblies with several reference-free evaluation tools (see Methods) allowed us to determine

209 the highest quality assemblies (Figures 1C, 2A-D, Supplemental Figure S1). Rankings from these

210 metrics were generally correlated with each other (Figure 2A, Supplemental Figure S2A).

211 Platanus and ABySS assemblies most consistently returned the best rankings across metrics with

212 Platanus assemblies having higher mean ranks overall (Figure 2A and Supplemental Figure S1).

213 All Illumina assemblies did moderately well in terms of gene content, most containing between

214 80-85% of the expected Arthropod BUSCOs (Figure 2B). Nonetheless, all of the Illumina

215 assemblies were highly fragmented, containing up to hundreds of thousands of contigs mostly

216 less than 1 kb in length. The NG50 values ranged from 2.5-7.3 kb (Figure 2C, Supplemental Table

217 S3). Although some scaffolds in the assemblies reached up to the Mb range, they were all

218 bacterial, a common observation for assemblies from whole animals (Supplemental Figure S3).

219 Of recognizable bacterial sequence, at the genus level, ~90% was characterized as *Delftia* and

220 ~5% as *Rickettsia*. Amongst all Illumina assemblies, the longest scaffolds of apparent insect origin

221 were 50-60 kb. Filtering for bacterial contamination and re-assembling with the filtered data did

222 not improve the contiguity (Supplemental Table S3). Although short-read-only assemblies were

223 not pursued further, the highest quality Platanus assembly was used in hybrid assemblies with

224 long reads, and the high accuracy Illumina short reads were useful for polishing long-read

225 assemblies.

226 **Figure 2: Assembly evaluations**

228 **Figure 2: Assembly evaluations. (A)** Rank matrix for 40 Illumina assemblies. Each column
229 corresponds to a metric. Each row corresponds to an assembly. The columns and rows are
230 organized by metric class and assembler, respectively. Multiple assemblies were generated for
231 each assembler differing by the input reads, parameters used, or both. Assembly nicknames allow
232 finding the assemblies in supplementary tables and methods. Assembly ranks for each metric
233 span from lowest (blue) to highest (red) in each column. Assemblies (rows) that do well across
234 the metrics tend to be mostly shades of red. The red star marks the Platanus assembly that
235 performed best overall and was used as the input for hybrid assemblies. **(B-D)** Use the short-read
236 assembly color scheme from (A) and the long-read color scheme from (E) to visualize **(B)** percent
237 of complete BUSCOs found, **(C)** Log10 NG50 lengths, and **(D)** ALE scores for short-read and
238 long-read assemblies. B and D show the long-read scores before and after polishing steps. The
239 dotted lines in (C) represent the maximum NG50 from short-read assemblies. **(E)** Rank matrix for
240 50 long-read assemblies organized as described in A. Red and blue stars mark assemblies
241 brought into BioNano scaffolding. Red stars represent the scaffolded assemblies that were
242 chosen after BioNano scaffolding. **(F-G)** Box and whisker plots of within-assembler rank
243 distributions comparing blended (red) to PacBio-only (blue) inputs to each assembler. The
244 boxplots are not comparable between assemblers. The boxes show the $25^{th}$-$75^{th}$ percentile, the
245 black line is the median, and the whiskers span the range (min to max). Assemblies from a given
246 assembler were ranked either using (F) all individual metrics from E or (G) all ranked mean ranks
247 from different combinations of metric ranks from E. The ranks were then partitioned into those
248 from blended versus PacBio-only assemblies. In both cases (F-G), blended assemblies from all
249 assemblers except SMARTdenovo had significantly higher ranks by Wilcoxon Rank Sum Test
250 than PacBio-only assemblies from the same assembler.

251 **<u>Long-read datasets and assemblies</u>**

252       A route to obtaining more contiguous assemblies is incorporating data from single

253 molecule, long-read technologies, such as Single Molecule Real Time (SMRT) sequencing from

254 Pacific Biosciences (PacBio) and nanopore sequencing with the MinION from Oxford Nanopore

255 Technologies (ONT). These technologies are more error-prone than Illumina, but the errors are

256 approximately randomly distributed allowing high quality consensus sequences with enough

257 coverage (Eid et al. 2009; Ip et al. 2015; Loman et al. 2015). Both long-read technologies

258 produced read lengths that exceeded the scaffold lengths in the Illumina short-read assemblies,

259 particularly MinION reads obtained using our modified protocols (Supplemental Figure S4; Urban

260 et al 2015). Thus, even before attempting to assemble the long reads, we had a richer source of

261 long-distance information than the short-read assemblies provided.

262

263       The majority of long-read coverage (50-55X total) was from PacBio (44.1X; Table 1;

264 Supplemental Figure S4), and we were able to produce high quality assemblies using PacBio

265 reads alone. However, despite having four times lower coverage, the MinION data (10.77X) had

266 in excess of two times more coverage from molecules greater than 20 kb and over a hundred

267 times more coverage from molecules exceeding 30 kb than the PacBio data (Table 1). Over 10%

268 of the MinION data was from molecules that surpassed the longest PacBio read length of 36 kb,

269 approximately a third of which came from high quality 2D reads (Table 1, Supplemental Figure

270 S4). Validation of the MinION reads on assemblies generated from the PacBio data alone showed

271 many high quality 1D and 2D reads (Supplemental Figure S5). These included hundreds of 2D

272 reads exceeding 50 kb and several >100 kb that aligned across their full lengths with percent

273 identities up to 94.6%. One notable 131 kb MinION 2D read aligned with 91.1% accuracy to the

274 PacBio data. This gave us an opportunity to test whether even a small amount of ultra-long

275 MinION reads could improve upon the PacBio assemblies. Therefore, we also generated

276 assemblies from a blend of both single-molecule technologies, referred to here as "blended

277     assemblies" to differentiate them from "hybrid assemblies" that refers to combining short-read and

278     long-read technologies (Figure 1C).

279

280          In total, we generated 50 assemblies using long reads (Figure 1C), including hybrid

281     assemblies that started from Illumina contigs. We evaluated the long-read assemblies with the

282     same metrics used to rank the short-read assemblies (Figure 2 B-D, Supplemental Figure S1).

283     Before polishing, ABruijn and Canu assemblies rose highest in most rankings (Figure 2 B-D,

284     Supplemental Figure S1), perhaps because these assemblers had the best consensus sequence

285     modules. Even before polishing, most long-read assemblies outperformed short-read assemblies

286     for percent error-free bases (REAPR) and had comparable or better scores in other metrics (e.g.

287     LAP, ALE, FRC). However, most underperformed in terms of gene content with fewer than 80%

288     BUSCOs detected (Figure 2 B-D, Supplemental Figure S1).

289

290     **Long-read assembly polishing and monitoring**:

291          To ensure that the assembly evaluations primarily reflected the structural integrity of each

292     assembly rather than differences in consensus quality, we employed extensive post-assembly

293     polishing using Quiver (Chin et al 2013) and Pilon (Walker et al 2014) (Figure 1C). We monitored

294     the outputs from each round of polishing using the metrics discussed above as well as the number

295     of variants detected and changes made by the polishing algorithms (Figure 1C). The assemblies

296     started out with up to millions of Quiver variants and converged to just a few thousand, and

297     evaluations improved across Quiver rounds, with the biggest impact occurring in the first round

298     (Supplemental Figure S6). After Quiver polishing, Canu assemblies continued to take many of

299     the highest ranks whereas ABruijn assemblies lost their lead (Figure 2 B-D, Supplemental Figure

300     S1). Quiver polishing also closed the gaps between the highest and lowest scoring assemblies in

301     each metric. For example, whereas the percent of BUSCOs detected ranged from 30-83% prior

302     to Quiver polishing, ~90% were detected in all assemblies after (Figure 2B). Moreover, all polished

15

303 long-read assemblies outperformed the best scoring short-read assemblies in each metric, with

304 the exception of the hybrid assemblies that still underperformed on the ALE metric (Figure 2D).

305 The Illumina-based metrics favored non-hybrid long-read assemblies over both the short-read

306 and the hybrid assemblies that were constructed from the same Illumina data. This speaks to the

307 structural and consensus quality of the contig sequences derived from long reads alone (Figure

308 2D, Supplemental Figure S1; "After Quiver"). Nevertheless, Illumina-polishing with Pilon improved

309 the consensus further, fixing 19.2-25.8 thousand base and small indel errors (~60-90 errors/Mb)

310 in the first round, and 0.9-2.4 thousand (~3-8 errors/Mb) in the second. The small number of

311 corrections introduced in the final round indicates long stretches (hundreds of kb) of high-quality

312 consensus sequences between any remaining errors in the final assemblies. Accordingly, Pilon

313 tended to improve evaluations modestly over what Quiver had already accomplished (Figure 2B,

314 2D, Supplemental Figure S1; "After Pilon"). For example, it resulted in detecting up to an additional

315 1.05% of BUSCOs (0.63% on average).

316

317 **Selecting assemblies for BioNano scaffolding**:

318 After polishing, the number of variants or genes detected and other metrics that reflect

319 consensus sequence quality converged to similar scores across assemblies. This allowed us to

320 focus on the size and long-range integrity of contigs when making selections for scaffolding with

321 optical maps. We used an expanded battery of reference-free metrics to guide our choice of which

322 assemblies to scaffold (Figures 1C and 2E). The additional metrics were based on long reads and

323 optical maps (see Methods). There was general agreement on assembly rankings among metrics

324 from the four orthogonal technologies (Supplemental Figure S2B).

325

326 Long-read assembly sizes ranged from 281.5-306.6 Mb (Supplemental Table S4), close

327 to the expected *Sciara* male somatic genome size of 280 Mb. All long-read assemblers produced

328 assemblies that were orders of magnitude more contiguous than short-read assemblies. NG50s

16

329    were typically in the Mb range and all exceeded 100 kb (Figure 2C, Supplemental Figure S1F,

330    Supplemental Table S4). For all size metrics, assemblies from Canu and Falcon ranked highest

331    (Figure 2C, 2E), with the largest NG50s of 3.08 Mb and 3.17 Mb, respectively (Figure 2C "Long

332    Read", Supplemental Table S4). Canu and Falcon assemblies also had the lowest LG50s

333    containing 50% of the expected genome size on just 21 and 23 contigs, respectively

334    (Supplemental Figure S1F, Supplemental Table S4). The highest normalized expected contig

335    sizes (Salzberg et al. 2012) for assemblies from Canu and Falcon exceeded 5 Mb and the longest

336    contigs from each exceeded 20 Mb (Supplemental Figure S1F, Supplemental Table S4).

337

338    Longer contigs can simply be a consequence of more aggressively joining reads at the

339    cost of more misjoins. Therefore, we interrogated whether Canu and Falcon assemblies, which

340    had the longest contigs, suffered from higher error rates. However, in direct opposition, Canu and

341    Falcon assemblies were consistent rank leaders in our battery of evaluations (Figure 2E). Canu

342    assemblies led most Illumina-based and long-read metrics. Falcon assemblies led BioNano

343    metrics and gene content (Figure 2E; Supplemental Figure S1), although differences in gene

344    content were negligible (Figure 2B). Canu and Falcon assembles had fewer putative mis-

345    assemblies than others as proxied, for example, by long-read detection of structural variants

346    (Supplemental Figure S7J). They also had apparently higher long-range integrity according to

347    BioNano map alignments, which spanned a range of 237-252 Mb in Falcon and Canu assemblies,

348    but only 181-230 Mb in others (Supplemental Figure S7H, S7J, S7L). In sum, Canu and Falcon

349    assemblies had longer contigs and ranked higher than other assemblies in most metrics (Figure

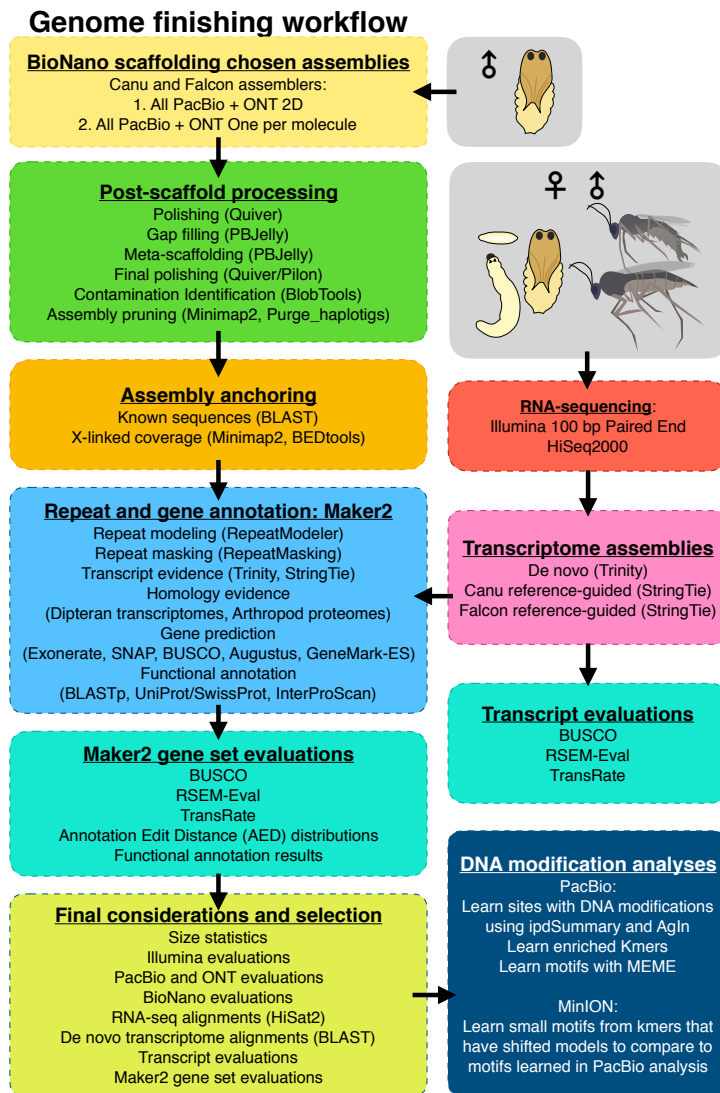350    2E), the latter arguing against more misjoins.

351

352    To select a final subset of assemblies for BioNano hybrid scaffolding, we sorted the

353    assemblies by taking mean ranks across 40 combinations of the 27 metrics (Figure 2E). In

354    general, blended assemblies tended to rank higher than their PacBio-only counterparts for 5 of

355    the 6 long-read assemblers, although this often reflected modest improvements in the actual

356    scores. The largest variation amongst scores tended to reflect the assembler used (Figure 2 F-

357    G, Supplemental Figure S7). Blended assemblies from Canu and Falcon were the clear rank

358    leaders again in this final analysis (Figure 2 E-G), and two assemblies from each were chosen for

359    BioNano hybrid scaffolding. The chosen assemblies were constructed from all 44X PacBio data

360    and either only 2D MinION reads (6.2X) or 1D and 2D reads (10.8X).

361

362  **FIGURE 3: Post-assembly work flow:**



363

364  **FIGURE 3: Post-assembly work flow:**

365  Workflow starting after selecting assemblies for BioNano scaffolding. Chosen assemblies were
366  scaffolded, polished, gap-filled, filtered for contamination, anchored into chromosomes by
367  sequences with known chromosomal addresses, and anchored to the X or autosomes by haploid
368  or diploid coverage. Repeats were characterized and RNA-seq was used to facilitate
369  transcriptome assembly and gene annotation. The single-molecule datasets were re-used to
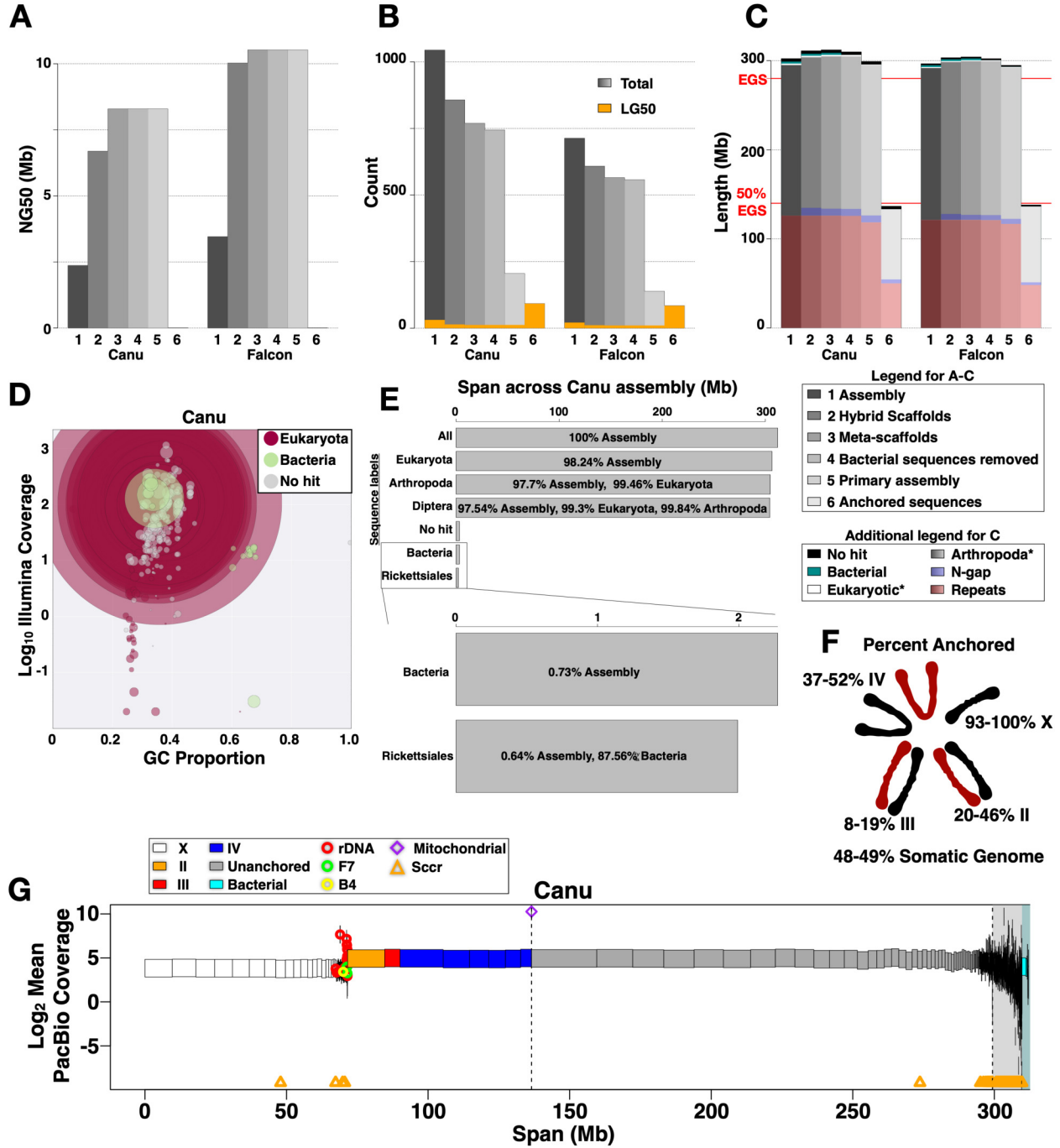370  investigate DNA modifications.

**Scaffolding with optical maps**

We obtained BioNano Irys optical map data from male pupae (Figure 3, Table 1). The raw molecule N50 was 214.1 kb for molecules >150 kb. The genomic consensus maps (CMAPs) produced from them had a map N50 of 712 kb and a cumulative length of 325.5 Mb. Thus, the genome length estimated from the CMAPs was between the expected sizes of the somatic and germline genomes. The CMAPs spanned 266-278 Mb of the Canu and Falcon contigs. The CMAPs and sequence assemblies were used to produce the hybrid scaffold maps. Both the CMAPs and sequence contigs had similar spans across the hybrid scaffold maps of approximately 275-280 Mb. We found that the hybrid scaffolds derived from both Canu assemblies and from both Falcon assemblies were nearly identical as determined by evaluation metrics and whole genome alignments (Supplemental Figures S8-S9, Supplemental Table S5). Therefore, we chose the single scaffold set from each pair that was evaluated to be slightly better, hereafter referred to as "Canu" and "Falcon". Hybrid scaffolding approximately tripled the contiguity of the assemblies (Figure 4A, Supplementary Tables S6, S7). Throughout the following text, Canu assembly statistics will be described with corresponding Falcon statistics in parentheses. The total numbers of sequences in the Canu (Falcon) assembly decreased from 1044 to 857 (713 to 608) while increasing the NG50 of 2.3 Mb to 6.7 Mb (3.5 Mb to 10 Mb). The assembly size also increased from 302 Mb to 311 Mb (296 Mb to 303 Mb) (Figure 4A-C). The Canu (Falcon) scaffolds had 187 (105) gaps summing to 8.7 Mb (6.7 Mb) with a maximum gap size of 677 kb (965 kb) and median of 20.8 kb (30.5 kb) (Supplemental Table S8).

We next iteratively filled and polished the gaps using PBJelly (English et al. 2012) and Quiver, respectively. In the Canu (Falcon) scaffolds, 31 (14) gaps were completely closed and over 972 kb (1.06 Mb) of gap sequence was filled in (Figure 4C, Supplemental Table S8). In the final round of gap filling, we allowed PBJelly to "meta-scaffold" the hybrid scaffolds using connections from long-read alignments. This decreased the total number of sequences in Canu

20

397    (Falcon) from 857 to 769 (608 to 565) while increasing the NG50 of 6.7 Mb to 8.3 Mb (10.0 Mb to

398    10.5 Mb) and the assembly size from 311 Mb to 312 Mb (303 Mb to 304 Mb) (Figure 4A,

399    Supplemental Table S6, S7). We used both Quiver and Pilon to correct errors in the gap-filled

400    meta-scaffolds. In the final round, Pilon made only 18-27 changes to the consensus sequences,

401    translating to 1 change per 16.9 Mb and 11 Mb of non-gap sequence for Canu and Falcon,

402    respectively.

## 403 **Figure 4: Assembly scaffolding and anchoring**

**Figure 4: Assembly scaffolding and anchoring**

**(A)** NG50 of the assembly at different stages 1-6 as defined in "Legend for A-C" within the figure.

**(B)** Number of sequences in the assembly at different stages 1-6 as in A. The orange portions are LG50 counts, or the number of the longest sequences in each set needed to reach 50% of the estimated genome size (EGS = 280 Mb) for the somatic genome. **(C)** The total length of the assembly at different stages 1-6 as in A. The "Additional legend for C" defines colored portions of the bars. *The length of the Eukayotic and Arthropod labeled sequences includes everything up through that color. **(D)** Log10 Illumina coverage versus GC content over the Canu assembly (similar results for Falcon), colored by taxonomy information, and with circle sizes proportional to the contig sizes they represent. **(E)** The proportion of the assembly taxonomically labeled as Eukaryotic, Arthropoda, Diptera, Bacteria, and Rickettsiales. **(F)** The percentage of the expected genome size and chromosome sizes that has been anchored. Ranges represent range in Canu and Falcon assemblies. Colors as in Figure 1. **(G)** The Canu assembly with scaffolds drawn as rectangles corresponding to their lengths, colored according to the chromosome they were anchored into (or unanchored), and at their mean coverage from PacBio reads, the dataset used to determine X-linked sequences by haploid level coverage. The white background highlights sequences in the primary assembly whereas the grey and cyan backgrounds are set behind associated and bacterial sequences, respectively. All sequences to the left of the first vertical dashed line are anchored.

425 **Assembly cleaning**

426    BlobTools (Laetsch and Blaxter 2017) was used to identify contaminating contigs in the

427 final scaffolds (Figure 4C-E, Supplemental Figure S10, S11). *Sciara* male embryo coverage from

428 Illumina, PacBio, and the MinION all gave similar results (Supplemental Figure S10). The vast

429 majority of the final Canu and Falcon scaffolds (>97.7% of the total sequence length) was

430 identified as Arthopoda, >99% of which was also Dipteran (Figure 4C, 4E, Supplemental Figure

431 S11). Canu and Falcon had 25 and 8 bacterial contigs respectively, with total lengths of 2.0-2.3

432 Mb (<1% of the total sequence length) and bacterial contig N50s of 1.0-1.3 Mb (Figure 4C, 4D,

433 4E, 4G; Supplemental Figure S11, Supplemental Table S9). There were no BioNano optical map

434 alignments over the bacterial contigs, and accordingly there were no bacterial contigs attached to

435 or found in any of the final Arthropod-associated scaffolds. Removing bacterial contigs only

436 marginally affected contig size statistics of the *Sciara* assemblies (Figure 4G; Supplemental

437 Tables S6, S7).

438

439    No bacterial contigs were labeled as Delftia in the long-read assemblies despite it being

440 the major bacterial representation in short-read assemblies. The majority of the bacterial

441 sequence (87-96%) in the Canu and Falcon scaffolds was labeled as Rickettsiales (Figure 4D-E,

442 Supplemental Figure S11), nearly all of which was associated with *Rickettsia prowazekii* (88.5-

443 90.1%) and *Rickettsia peacockii* (9.9-10.8%). Given that the published genome sizes for these

444 *Rickettsia* species range from ~1.1-1.3 Mb (Andersson et al. 1998; Felsheim et al. 2009), it is

445 possible that a complete *Rickettsia* genome sequence was co-assembled. The genus *Rickettsia*

446 includes obligate intracellular bacteria that may be the closest extant relatives to the ancestor of

447 the mitochondrial endosymbiont (Andersson et al. 1998). *Rickettsia* is closely related to

448 *Wolbachia* that is found in many strains of *Drosophila melanogaster* (Clark et al. 2005). The

449 *Rickettsia* or *Rickettsia*-like species in our *Sciara* datasets may be an important part of *Sciara*

450 biology. Interestingly, in the Illumina, PacBio, and MinION datasets, the *Rickettsia* genome has

24

451    nearly the same coverage as the *Sciara* genome (Figure 4D, 4G, Supplemental Figure S10). This

452    indicates that there is approximately one *Rickettsia* genome per haploid *Sciara* genome or two

453    *Rickettsia* for each diploid *Sciara* cell in male embryos on average. The current evidence can only

454    suggest that this correspondence is coincidental. Despite high Rickettsia coverage in embryos,

455    there were no *Rickettsia* optical maps from pupae. This may reflect the DNA plug isolation

456    procedure used and/or a far lower abundance of *Rickettsia* in pupal cells or its restriction to a

457    small subset of cells.

458

459        After removing bacterial sequences, each assembly was partitioned into "primary" and

460    "associated" sequences. Primary sequences represent one haplotype of the genome whereas

461    associated sequences consist of short redundant contigs called haplotigs that represent other

462    haplotypes of specific loci (Figure 4G). The Canu (Falcon) assembly contained 744 (557)

463    sequences of which 205 (138) were primary and 539 (419) were associated, giving a primary

464    assembly size of ~299 Mb (~295 Mb) with ~13 Mb (9.4 Mb) of associated sequences (Figure 4A-

465    C, Supplemental Tables S6, S7). The difference of ~4 Mb between the Canu and Falcon primary

466    assembly sizes is in part owed to Canu having ~2.2 Mb more gap length than Falcon. Given that

467    the associated sequences are generally short (~23 kb on average), computing size statistics on

468    the primary assembly has relatively large effects on the mean and median contig sizes

469    (Supplemental Tables S6, S7). For example, the mean scaffold size in Canu (Falcon) increased

470    from ~416 kb to 1.5 Mb (542 kb to 2.1 Mb).

471 **Table 2: Anchoring into chromosomes using previously known sequences**

| Sequence | Location | Canu contig size | Falcon contig size | Reference |
|---|---|---|---|---|
| DNA puff II/9A | Chr II locus 9A | 13.1 Mb | 28.5 Mb | DiBartolomeis and Gerbi 1989; Bienz-Tadmor et al. 1991; Urnov et al 2002; Foulk et al. 2006 |
| RNA Puff III/9B | Chr III locus 9B | 5.4 Mb | 12.5 Mb | Wu et al 1993; Foulk et al. 2006 |
| Ecdysone receptor | Chr IV locus 12A | 3.8 Mb | 9.6 Mb* | Foulk et al 2013 |
| Ultraspiracle | Chr IV locus 10A | 9.3 Mb | 5.5 Mb | Foulk et al 2013 |
| Hsp70 | Chr IV locus 4A or 12C | 5.4 Mb | 13 Mb | Mok et al. 2000 |
| Hsp70 | Chr IV locus 4A or 12C | 6.8 Mb | 2.6 Mb | Mok et al. 2000 |
| ScoHet1 | Chr IV locus 5A | 15.2 Mb | (9.6 Mb)* | Greciano et al. 2009 |
| ScoHet2 | Chr IV locus 12C-13A | 5.9 Mb | 4 Mb | Greciano et al. 2009 |
| rDNA | End of Chr X | 5 primary contigs and 11 associated contigs (Σ 1.3 Mb) | 2 primary contigs and 41 associated contigs (Σ 1.7 Mb) | Pardue et al. 1970; Gerbi 1971; Crouse et al. 1977; Kerrebrock et al. 1989 |
| Microclone B4 | End of Chr X | 69.8 kb | 59 kb | Escribá et al. 2011 |
| Microclone F7 | Near centromere of Chr X**, non-centromeric Chr IV, L chromosomes | 3 associated contigs (Σ 66.8 kb) | 1 primary and 1 associated contig (Σ 161.6 kb) | Escribá et al. 2011 |
| Microclone G2 (Sccr) | Centromeres of all chromosomes | 20 primary and 85 associated contigs (Σ 1.3 Mb) | 6 primary and 42 associated contigs (Σ 604 kb) | Escribá et al. 2011 |

472 *Ecdysone receptor (EcR) and ScoHet1 identified the same 9.6 Mb contig in Falcon. The locus
473 inconsistency may represent a misassembly in Falcon or misannotation from Greciano et al
474 (2009). Nevertheless, both EcR and ScoHet1 results agree it is from chr IV.
475 ** Coverage analyses confirm contigs with F7 as chromosome X sequence.
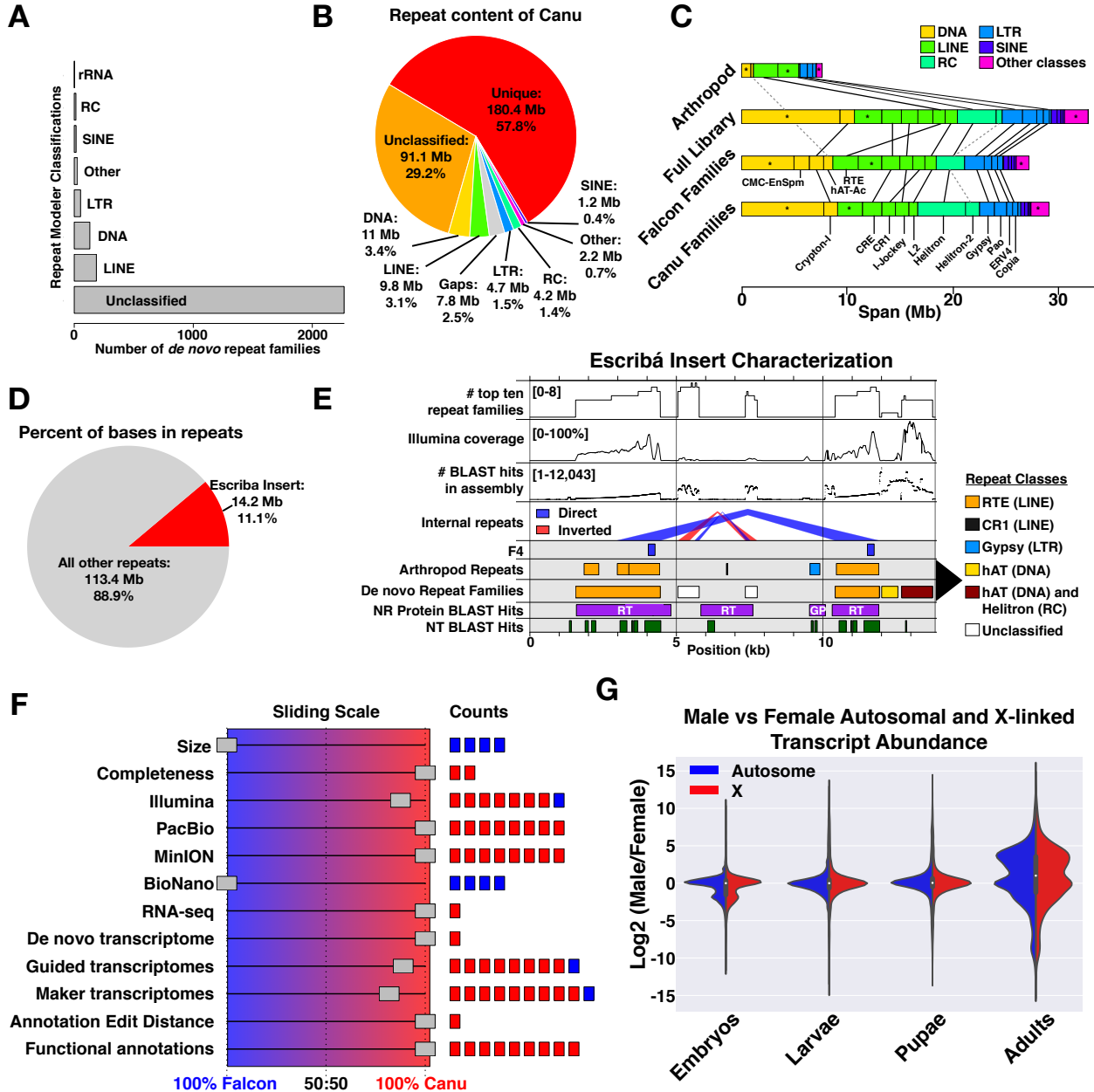
26

**Assembly anchoring**

We used previously known sequences with associated *in situ* hybridization results from polytene chromosomes to anchor some of the scaffolds into chromosomes (Table 2). The results span all 3 autosomes and the X chromosome. We anchored 7-8 primary autosome-linked contigs from each assembly that sum to 64.9-75.6 Mb, or ~23-27% of the expected somatic genome size and 28-33% of the expected autosomal sequence length. Given the number of regions determined for each chromosome from polytene banding patterns (Gabrusewycz-Garcia 1964), we expect chromosomes II, III, and IV to be approximately 62-66 Mb, 66-71 Mb, and 88-94 Mb, respectively (Supplementary Table S1E). Therefore, across both assemblies we expect to have anchored 20-46% of II, 8-19% of III, and 37-52% of IV. Since it is possible to transfer anchoring information from one assembly to the other, the overall anchoring percentages for both assemblies are essentially the higher end of each range above. We also anchored between 1-2 Mb of X-linked contigs using repetitive sequences specific to the X (Table 2, e.g. rDNA). In addition to chromosome-specific sequences, "Sccr" (*Sciara* centromere consensus sequence) that hybridized to the centromeres of all *Sciara* chromosomes (Escribá et al. 2011) mapped to 48-105 contigs, the majority of which were not primary sequences (Table 2).

The majority of genomic DNA sequenced from male embryos came from somatic cells that are haploid for the X and diploid for all autosomes. Therefore, X-linked contigs could be defined as primary contigs with haploid level coverage across 80% or more of their lengths. The Canu (Falcon) assembly contained 69 (36) contigs called as X that summed to 71 Mb (62 Mb) with the longest X-linked contig reaching 9.68 Mb (12 Mb) and an X-linked contig N50 of 5.95 Mb (7.3 Mb). In both assemblies, contigs containing the X-hybridized repetitive sequences (Table 2: rDNA, B4) were called as X as expected (Figure 4G, Supplemental Figure S11C). Upon closer inspection, contigs with rDNA not called as X had haploid level coverage regions interrupted by regions of unusually high coverage over collapsed rDNA repeats, and were therefore consistent

502   with being X-linked sequences as well. We also found X-linked contigs that contained the F7

503   repeat sequence known to be on X, IV, and L (Escribá et al. 2011) (Figure 4G, Supplemental

504   Figure S11C). The X chromosome is estimated to be ~50 Mb based on DNA-Feulgen

505   cytophotometry or ~62 Mb based on the number of polytene bands (Rasch 2006; Gabrusewycz-

506   Garcia 1964; Supplementary Table S1 A-E). Therefore, >93% of the X chromosome sequence

507   could be anchored. In total, at least 136.6-138.0 Mb of *Sciara* sequence, or ~49% of the expected

508   somatic genome size, was anchored into specific chromosomes with 100% of the assembly

509   characterized as either X or autosomal from the coverage analysis.

**Figure 5: Repeats and genes in the chosen assembly**



**Figure 5: Repeats and genes in the chosen assembly**

**(A)** Number of *de novo* repeat families trained on Canu with each classification. **(B)** Pie chart partition of the Canu assembly into the major repeat categories. Note that the "DNA" label used by RepeatMasker refers to DNA transposons. **(C)** The major sub-classes of repeats in each repeat class in the Canu assembly, showing the results when using different repeat libraries for

29

517     masking. **(D)** Pie chart representing the number of bases masked by the Escribá insert (Escribá

518     et al. 2011) alone compared to all masked bases. **(E)** Characterization of the Escribá insert,

519     highlighting major repeats in the *Sciara* genome. Black arrowhead on right side pointing to repeat

520     classes legend corresponds only to the two repeat family rows. RT = reverse transcriptase. GP =

521     gag-pol. **(F)** The ranking results of the final two assemblies demonstrating how many metrics in

522     each category for which each assembler scored better. **(G)** Split violin plots showing the log2 of

523     the male to female ratio of transcript abundance for the X (red) and autosomes (blue) across

524     multiple life stages.

525 **Repeats in the *Sciara* genome**

526     To learn more about the repeat content of the *Sciara* genome and to facilitate repeat masking,

527     *de novo* repeat families were created from both the Canu and Falcon assemblies using

528     RepeatModeler (Smit and Hubley 2008). There were close to 2700 repeat families in each library

529     of which 15-19 were classified as SINEs, 160-186 as LINEs, 48-53 as LTR, 130-131 as DNA

530     elements, and 43-50 as other repeat classes (Figure 5A, Supplemental Figure S11D,

531     Supplemental Tables S10, S11). The majority of repeats in each library were unclassified. For

532     repeat masking, the *de novo* repeat libraries were combined with the few previously known repeat

533     sequences (see Methods) as well as repeats from across Arthropods. Using this comprehensive

534     repeat library, RepeatMasker (Smit et al. 2013) classified ~121-126 MB or 39-41% of the Canu

535     and Falcon assemblies as repeats (Figure 5B, Supplemental Figure S11E, Supplemental Tables

536     S12, S13). Assuming that scaffold gaps also correspond to repeats leaves ~180 Mb of unique

537     sequence (~58%) in the Canu assembly (Figure 5B). The majority (76.6-76.9%) of repeats were

538     unclassified and spanned 93.3-96.7 Mb (Figure 5B) whereas SINE, LINE, LTR, RC, and DNA

539     elements each constitute 0.4-3.4% of the assemblies (Figure 5B). DNA elements made up the

540     largest class in terms of total span and Crypton-I was the largest sub-class therein (Figure 5C).

541     However, Helitron elements from the RC class were the largest sub-class in the assembly overall

542     (Figure 5C). Simple repeats made up ~1% of the assemblies (Supplemental Table S12). Similar

543     results were obtained when repeat masking with only the *de novo* repeat libraries (Figure 5C).

544     However, using only known arthropod repeats found fewer, had a higher composition of LINE

545     elements, and found the RTE sub-class therein to be the most abundant sub-class in the

546     assembly (Figure 5C).

547

548     Previously, Escribá et al. (2011) published a 13.8 kb lambda phage insert sequence that

549     contains two copies of the non-LTR retrotransposon named ScRTE. A corresponding probe (F4)

550     predominantly labeled pericentromeric regions of all *Sciara* chromosomes by FISH (Escribá et al.

31

551    2011). We found that the 13.8 kb Escribá insert contains some of the most abundant sequences

552    in the *Sciara* genome, although there was only one full-length copy of the lambda insert in each

553    assembly presumably from the locus that was originally cloned (Supplemental Figure S12).

554    Otherwise, pieces of the insert were scattered across the assembly corresponding to nearly

555    60,000 alignments spanning ~14.2 Mb, or ~11% of bases labeled as repeats (Figure 5D-E). Of

556    the top ten most abundant *de novo* repeat families, eight map to the Escribá insert across most

557    of their lengths at sites that are consistently over-represented in DNA sequencing coverage and

558    BLAST hits from other regions of the genome, and correspond to direct repeats of the ScRTE

559    element, unclassified inverted and direct repeats, as well as hAT and Helitron elements (Figure

560    5E).

561

562    **<u>Transcriptome assembly and gene annotation</u>**

563        We annotated protein-coding genes in the Canu and Falcon genome assemblies with

564    Maker2 (Holt and Yandell 2011) guided by transcriptome assemblies from poly-A enriched RNA-

565    seq datasets from *Sciara* male and female embryos, larvae, pupae, and adults (Figure 3). With

566    the gene sets available from each assembly, we performed a final set of reference-free

567    evaluations to choose a final assembly: Canu or Falcon (Figure 3). The Falcon assembly had

568    slightly longer contig size statistics and a corresponding lead in metrics from optical map

569    alignments (Figure 5F, Supplemental Figure S13). However, the Canu assembly outperformed

570    Falcon in completeness metrics, RNA-seq and *de novo* transcriptome alignments, as well as

571    metrics from Illumina, PacBio, and MinION datasets (Figure 5F, Supplemental Figure S13).

572    Moreover, both the Canu-guided transcriptome assembly and the transcripts in the final Canu

573    annotation received higher evaluation scores than their Falcon counterparts (Figure 5F,

574    Supplemental Tables S14, S15). Finally, the Canu annotation had lower annotation edit distances,

575    more genes with GO terms, Pfam domains, and/or BLAST hits in the UniProt-SwissProt database,

576    more BUSCOs, as well as more hits with proteomes from *Drosophila melanogaster* and

577  *Anopheles gambiae* (Figure 5F, Supplemental Figure S14, Supplemental Table S16). We

578  conclude that the Canu assembly had higher consistency with the genome sequencing datasets

579  and yielded the superior gene set. We therefore chose the Canu assembly as the first draft

580  genome for *Sciara coprophila.*

581

582  The final annotation of the Canu assembly had 23,117 protein-coding gene models with

583  28,870 associated transcripts (Supplemental Table S15A). There are more genes than expected

584  from other fly genomes, which may be a result of gene splitting in the annotation. To increase the

585  quality of the *Sciara* gene set, the annotation was deposited at the i5k-workspace for community-

586  enabled manual curation (https://i5k.nal.usda.gov/). Nevertheless, in its current form, the

587  annotation contains nearly all expected Dipteran genes: 94.2% complete Dipteran BUSCOs were

588  found in the final gene set, 97% when including fragmented BUSCOs (Supplemental Figure S14E,

589  Supplemental Table S15A). The majority of genes in the annotation (87.5%) had only a single

590  associated transcript isoform (Supplemental Figure S14B). The median gene and transcript

591  lengths are ~2.6 kb and ~1.3 kb, respectively (Supplemental Table S15A). Genes had a median

592  of 4 exons, ranging from just one (10.8% of genes) to over 100 exons. There are 10,801 genes

593  with both 5' and 3' UTRs annotated and 13,335 with one or the other. Exons, introns, 5' and 3'

594  UTRs had median lengths of 182 bp, 80 bp, 165 bp and 184 bp, respectively. Of all genes, we

595  were able to attach functional information to ~65%. Specifically, 8671 (37.5%) have Ontology

596  Terms, 13745 (59.5%) have UniProt/SwissProt hits, 13789 (59.6%) have Pfam descriptions (El-

597  Gebali et al. 2019), 8252 (35.7%) have all three, and 14961 (64.7%) have one or more

598  (Supplemental Figure S14F, Supplemental Table S16). Genes spanned over 54% of the Canu

599  assembly, mostly attributable to introns, and ~20% of the assembly was both unique (not

600  repetitive) and intergenic (Supplemental Figure S14H).

601

602     In the standard Dipteran model, *Drosophila melanogaster*, where males are XY and

603     females are XX, male flies exhibit dosage compensation of X-linked genes. We used the *Sciara*

604     gene annotation and anchoring information to explore dosage compensation in *Sciara.* The

605     majority of cells in *Sciara* embryos, larvae, pupae, and adults are somatic where X ploidy differs

606     between males and females. Males are haploid and females are diploid for the X, respectively.

607     Both sexes are diploid for autosomes. We defined genes as X-linked if they were on contigs

608     anchored into the X chromosome by the coverage analysis described above. We then determined

609     if there was dosage compensation for X-linked genes, or if they consistently had 2-fold lower

610     transcript abundances in male samples. Across each stage of development sequenced, the

611     distributions of log2 fold changes between male and female transcript abundance were the same

612     for autosomal and X-linked genes (Figure 5G, Supplemental Figure S15). There were many

613     examples of both autosomal and X-linked genes that were differentially expressed between males

614     and females, but there was no difference between males and females for the majority of genes in

615     both classes. Therefore, the evidence strongly supports the existence of dosage compensation

616     of most X-linked genes in *S. coprophila*.

617

618     **DNA modification signatures in single-molecule data**

619     The mechanism for imprinting in *Sciara* remains unknown. Since imprinting in mammals

620     utilizes DNA methylation (Li et al. 1993), it was if interest to determine whether DNA modifications

621     are present in *Sciara*. The gene annotation contains the proteins involved in cytosine and adenine

622     methylation pathways (reviewed in Armstrong et al 2019; Rausch et al 2020) that are expected to

623     be found in Dipterans, including putative homologs for DNMT2, TET-family proteins, DAMT-

624     1/METTL4, N6AMT1, ALKBH1, jumu, and several proteins with methyl-CpG binding domains

625     (Supplemental Table S17A-C). There was also evidence of DNA modifications in the *Sciara*

626     genome found using anomalies in the raw signal of both single-molecule, long-read sequencing

627     datasets (Flusberg et al. 2010; Clark et al. 2012; Simpson et al. 2017). The high-coverage PacBio

628    dataset was used to call site-specific modifications in the assembly for 5mC, 4mC, and 6mA. The

629    low-coverage MinION dataset was used to find kmers with signal distributions that were shifted

630    compared to expected models, which could result from DNA modifications. These kmers were

631    used to find sub-motifs for comparison to motifs obtained in the PacBio analysis.

632

633         Using the PacBio SMRT kinetics data, we estimated that ~0.13-0.24% of adenine sites in

634    the *Sciara* male embryo genome were potentially modified with up to ~0.04-0.06% of adenine

635    sites exhibiting the 6-methyl-Adenine (6mA) signature (Figure 6A, Supplemental Table S18A),

636    which is similar to 6mA densities seen for humans (~0.05%; Xiao et al. 2018), some fungi

637    (~0.05%; Mondo et al. 2017), *Drosophila embryos* (0.07%; Zhang et al. 2015), and pig (0.05%;

638    Liu et al. 2016). The tens of thousands of modified adenines were distributed ubiquitously

639    throughout the assembly, including in genes and repeats as well as on both autosomal and X-

640    linked sequences (Supplemental Figure S16A-C). Over 50% of the reads aligning to the majority

641    of 6mA sites were estimated to contain 6mA (Figure 6B), suggesting that while the mark may be

642    rare in the genome it is common at those sites. Although adenine modifications were found in

643    many dimer and trimer contexts, AG and GAG were most enriched (Figure 6C, Supplemental

644    Figure S16D, Supplemental Tables S19, S20). GAG sites were modified up to 7-8 times more

645    frequently than the rate for A alone, with 0.9-1.7% GAG sites flagged as modified and 0.3-0.5%

646    flagged as 6mA specifically (Supplemental Table S18B). The frequencies of bases surrounding

647    the 6mA position in enriched 7mers showed a prominent 4 bp GAGG motif (Figure 6D), which did

648    not differ between X and autosomal sequences (Supplemental Figure S17). Other motifs

649    associated with 6mA in the *Sciara* genome included CAG within them (Supplemental Figure S18).

650    AG, GAG, GAGG, and CAG motifs were also previously found associated with 6mA sites in

651    human, rice, and *C. elegans* genomes (Greer et al. 2015; Xiao et al. 2018; Zhou et al. 2018). We

652    found that 6mers defined by the sequence logo from enriched 7mers showed shifted MinION

653    signal distributions whereas other control kmers fully agreed with the expected model (e.g.

654     CGAGGT; Figure 6E-F, Supplemental Figure S19). From the set of all kmers with shifted MinION

655     signals, we found similar motifs to those found in the analysis of 6mA sites identified in the PacBio

656     analysis (Figure 6G, Supplemental Figure S18).

657

658         We also used the PacBio SMRT kinetics data to look at cytosine methylation, which has

659     been previously shown to mark heterochromatic regions in *Sciara* chromosomes by

660     immunofluorescence (Eastman et al. 1980; Wei et al. 1981; Greciano et al. 2009). Up to 0.6-1.1%

661     of cytosines were modified with up to 0.11-0.24% and 0.26-0.43% showing 4-methylcytosine

662     (4mC) and 5-methylcytosine (5mC) signatures, respectively (Figure 6A, Supplemental Table

663     S17C). Modified cytosines were present throughout autosomal and X-linked sequences

664     (Supplemental Figure S16A-C). The frequency of methylation at the majority of 4mC and 5mC

665     sites was estimated to be over 80% (Figure 6B), despite being rare in the genome overall.

666     Modified cytosines were found in all dimer and trimers, but the most enriched were CG and GCG

667     (Figure 6C, Supplemental Figure S16D). This is reflected in the sequence logos constructed from

668     enriched 7mers centered on the modified C position (Figure 6D, Supplemental Tables S19, S20),

669     and was the same for autosomes and the X (Supplemental Figure S17). Up to ~1.3-2.5% of CpG

670     dinucleotides were estimated to be modified with 0.26-0.57% and 0.55-0.96% specifically

671     classified as 4mCpG and 5mCpG, respectively (Supplemental Table S18D). A more sensitive

672     algorithm (Suzuki et al. 2016) estimated as high as 6.4% of CpG dinucleotide sites in the genome

673     as targets for methylation in male *Sciara* embryos (Supplemental Table S18E). GCG sites were

674     modified up to ~4-5 times more frequently than the rate for C alone and 2 times more than CG,

675     with 2.5-4.9% of GCG sites flagged as modified and 0.5-1.2% and 0.9-1.5% of GCG sites flagged

676     as 4mC and 5mC, respectively (Supplemental Table S18F). Interestingly, GCG trimers are

677     depleted in the genome sequence whereas GTG trimers are enriched (Supplemental Figure S20).

678     This suggests that GCG may be a methylation target in the germline where 5mC deamination and

679     conversion to thymine can deplete GCG trimers over evolutionary time. We found that 6mers

680   defined by the sequence logos from enriched 7mers displayed shifted MinION signal distributions

681   (e.g. for TTCGGT and GGCGGA) whereas control kmers did not (Figure 6E-F, Supplemental

682   Figure S19), and that many motifs similar to those in the PacBio analysis specific to 4mC and

683   5mC were found when looking for motifs in kmers with shifted MinION signal distributions (e.g.

684   GCG; Figure 6G, Supplemental Figure S18).

685

686        The distribution of distances between adjacent DNA modifications, for both methylated C

687   and A, showed an enrichment of shorter distances than expected by chance (Figure 6H). There

688   were spikes of enrichment with a periodicity of 10 bp out to distances of at least 200 bp when

689   looking at both strand-agnostic and strand-specific spacings (Figure 6H, Supplemental Figures

690   S21-22). This periodicity is highly suggestive of turns of the DNA helix. Periodic spacing of 10 bp

691   between methylation sites and target motifs has been observed enriched over nucleosome

692   positions in *Arabidopsis* and mammals (Jia et al. 2007; Chodavarapu et al. 2010; Collings and

693   Anderson 2017). Moreover, 6mA was shown to be phased between nucleosomes in

694   *Chlamydomonas* and *Tetrahymena* (Fu et al. 2015; Wang et al. 2017; Luo et al. 2018). Indeed,

695   ~175 bp is one of the most enriched distances separating two modifications in our *Sciara* male

696   embryo data (Figure 6H, Supplemental Figures S21-22), a length reminiscent of nucleosomal

697   spacing in general and the exact length found for nucleosome intervals in *Drosophila* (Mavrich et

698   al. 2008).

699

700        We searched for relationships between DNA modifications and genomic features. The

701   trends were the same for all modification types (6mA, 4mC, 5mC). With respect to annotated

702   protein-coding genes, DNA modifications were random or slightly depleted, though there were

703   slight depletions in exons and promoters and slight enrichments in introns (Supplemental Figure

704   S16B, Supplemental Table S21A-B). These trends were the same when using gene locations

705   defined by the StringTie transcriptome assembly (Supplemental Table S21C) and were generally

706    true even when splitting the genes into categories of not expressed, lowly expressed, and highly

707    expressed using male embryo RNA-seq data (Supplemental Table S21D). Repeat regions in the

708    genome had more modifications than expected, and conversely the non-repeat regions had fewer

709    than expected (Supplemental Figure S16B, Supplemental Table S21E-F). In the *de novo* repeat

710    library, there were repeat families, including simple repeats, with 2-100 fold more modifications

711    than expected and many families with no modifications indicating that specific classes of repeats

712    are targeted for DNA methylation.

713

714    **Candidate germline-limited L sequences:**

715    L chromosome sequences are likely to be absent or of very low abundance in our datasets

716    from late male embryos. Nevertheless, an effort was made to identify candidate L-sequences.

717    Similar to identifying X-linked contigs in the assembly based on haploid-level coverage, L

718    candidates were gathered based on very low coverage, which may include junk or redundant

719    contigs. There were 25 contigs summing to ~230 kb that had at most 3X PacBio coverage across

720    their lengths in contrast to the genomic average of ~34X. These sequences were comprised of

721    ~60% repeats compared with ~40% genome-wide. The most abundant repeats were unclassified,

722    but Gypsy, Pao, and Zator transposons were highly represented. There were 15 mRNA isoforms

723    annotated in these low-coverage contigs corresponding to 13 genes (Supplemental Table S22).

724    Seven genes had putative functional information. Six had best hits to UniProt-SwissProt proteins

725    corresponding mostly Drosophila proteins, including Facilitated trehalose transporter Tret1, Ras-

726    related protein Rab-3, Ubiquitin carboxyl-terminal hydrolase 36, RNA-directed DNA polymerase

727    (jockey reverse transcriptase), Vitellogenin-2, Rho GTPase-activating protein. One was a protein

728    of unknown function, but was classified in the ribosomal L22e protein family from the Pfam domain

729    analysis.

730

38

731       An additional attempt to identify L candidate sequences was made by rescuing

732    unassembled reads. Since sequences of relatively low abundance may not have been

733    assembled, short and long reads that do not map to the Canu *Sciara* genome assembly were

734    used to generate new assemblies using Platanus and Canu, respectively. Using the unmapped

735    long reads, Canu returned 250 contigs summing to 2.55 Mb in length and 2247 unassembled

736    reads summing to 7.45 Mb for a total of 10 Mb. The majority of the sequences were identified as

737    bacterial: 89% of the total contig length and 69% of the total unassembled length for a total of 7.4

738    Mb of bacterial sequence (Supplemental Figure S23). Similarly, when assembling the unmapped

739    short read pairs, Platanus returned 330 scaffolds summing to 6.8 Mb, and ~96% was identified

740    as bacterial (Supplemental Figure S23).

741

742       We focused on the 698 long-read sequences (1.7 Mb) and the 159 short-read scaffolds

743    (141.5 kb) that were either identified as Arthropod or had no hits and 20-45% GC content

744    (Supplemental Figure S23). Only 12% of the total length of these sequences aligned to the

745    assembly, and just 14.5% and 28.3% of the targeted long-read and short-read sequences were

746    identified as repeats. The most abundant repeats present were unclassified. The repeats were

747    also enriched for simple repeats as well as transposons, such as Helitron (RC), Pao (LTR), RTE

748    (LINE), Jockey (LINE), and Mariner (DNA). The centromere-associated Sccr repeat (Escribá et al

749    2011) was on 14 contigs. There were also contigs with rDNA and rDNA transposons R1 or R2. A

750    small fraction of short reads contained the peri-centromeric tandem repeat B4 (Escribá et al

751    2011). Neither B4 nor rDNA has been observed on L chromosomes by in situ hybridization

752    (Escribá et al 2011), suggesting that at least some of these sequences are not from L

753    chromosomes. About 8% of the combined long- and short-read sequence length was covered by

754    hits from 116 proteins, 16% of which were transposon-related and another 65% of which had

755    other functional information. The most convincing alignments matched proteins (Supplemental

756    Figure S23), such as (i) an Integrator complex subunit that is involved in snRNA transcription and

757    processing, (ii) two zinc-finger transcription factors, (iii) a TATA-box binding protein, (iv) proteins

758    involved with chromosome cohesion, recombination, and segregation like Wings apart-like protein

759    (WAPL), Structural maintenance of chromosomes protein 6 (SMC6), and MOB kinase activator-

760    like 1 (mats), and (v) proteins involved in the nervous system. The majority of these were on
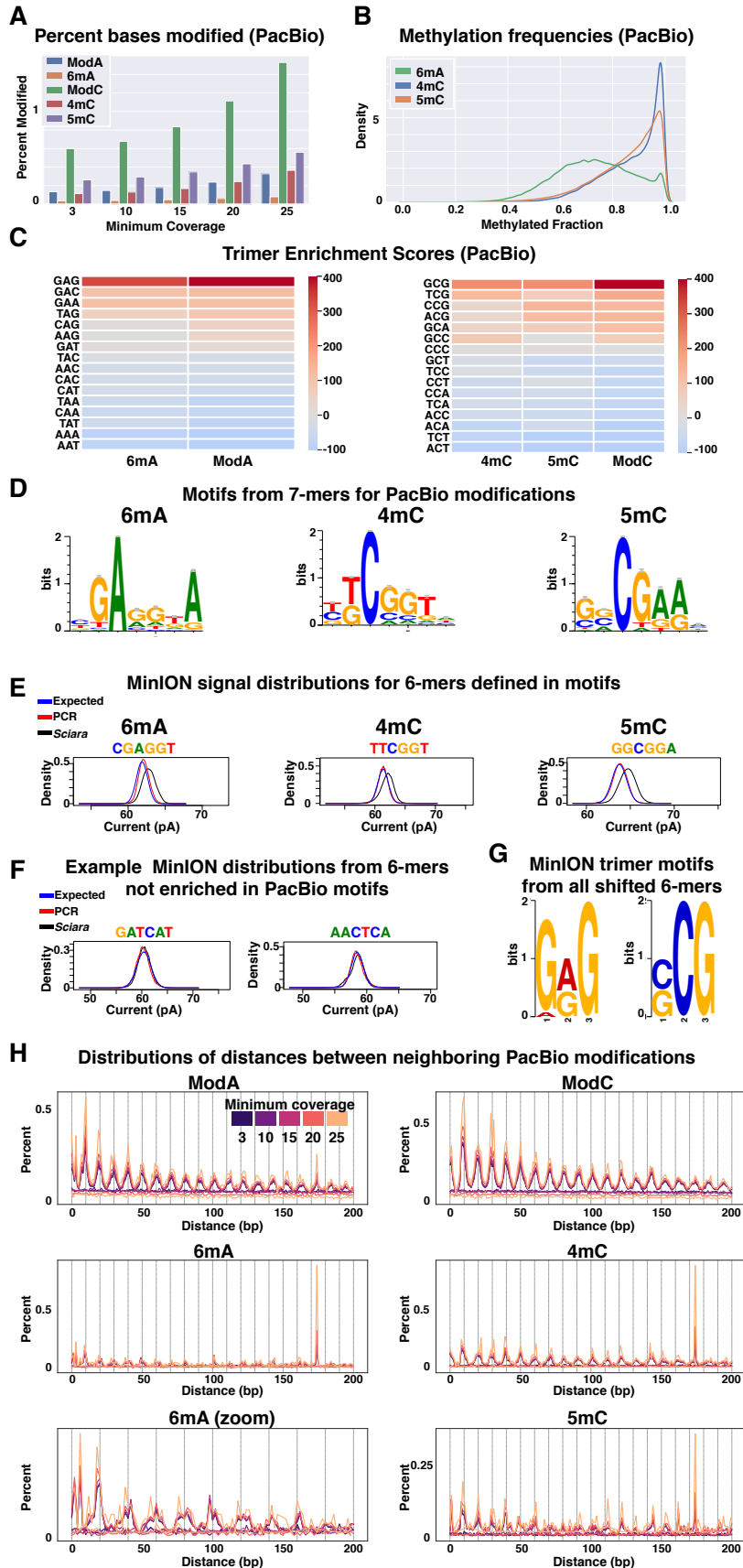
761    contigs with no matches to the genome assembly.

**Figure 6: DNA modifications in male embryo genomic DNA of *Sciara coprophila***

**(A)** Percent of adenines or cytosines assigned to a modification class given a minimum coverage level in the PacBio analysis. ModA and ModC are the sets of all adenines or cytosines, respectively, flagged as modified whereas 6mA, 4mC, and 5mC are the subsets of adenines or cytosines therein with those specific classifications.

**(B)** Methylation frequencies from the PacBio analysis at sites classified as having the given methylation type.

**(C)** Chi-square standardized residuals (enrichment scores) indicating how many standard deviations away each observation is from expectation for trimers with middle adenines or middle cytosines from the PacBio analysis.

**(D)** Position weighted motifs from the sets of 7-mers (where the modified base occurs at position 3) enriched for 6mA, 4mC, or 5mC.

**(E)** The distribution of ionic current means from the MinION data for 6-mers defined by the PacBio motifs in (D). The blue line shows the expected distribution given the MinION model for each kmer. The red line shows the distribution learned from whole *E. coli* genome PCR data (Simpson et al. 2017) using only canonical nucleotides. The black line shows the distribution learned from native genomic DNA from *Sciara*. Distributions are from template reads.

**(F)** As in (E), but showing examples of 6-mers not defined by motifs learned in the PacBio analysis.

**(G)** Two of the top three trimer motifs learned from the set of all 6-mers that had shifted MinION signal distributions with respect to the expected models.

**(H)** Distributions of distances between

762

41

763 **DISCUSSION**

764

765 **Derivation of the *Sciara* genome and anchoring**

766      We report here the assembled sequence of the male somatic genome of the lower Dipteran

767 fly, *Sciara coprophila*, as well as its gene annotation from transcriptomes covering both sexes and

768 all life stages. To find the assembly approach that worked best for the *Sciara* genome, we used

769 a battery of reference-free metrics to evaluate assemblies generated from different technologies,

770 algorithms, inputs, and parameters. *Sciara* genome sequences assembled with Canu and Falcon

771 from a blend of PacBio and MinION data, and polished with PacBio and Illumina data, performed

772 best and were selected for scaffolding using optical maps from the BioNano Genomics Irys

773 platform. Ultimately, the Canu scaffolds were the final selection for the first draft sequence owing

774 to their higher quality gene annotation. This release of the *Sciara* male somatic genome assembly

775 contains 299 Mb of sequence on 205 primary contigs with 50% of the expected genome size on

776 only 12 scaffolds that range from 8.2-23 Mb long. Annotating the *Sciara* protein-coding genes with

777 guidance from RNA-seq data gave a gene set that contained 97% of Dipteran BUSCOs,

778 suggesting it is essentially complete. We have anchored a significant amount of the *Sciara*

779 genome sequence on the three autosomes using previous *in situ* hybridization data, accounting

780 for 20-46% of chromosome II, 8-19% of chromosome III and 37-52% of chromosome IV. As *Sciara*

781 male somatic cells have only one X chromosome in contrast to two of each autosome, we were

782 able to use coverage levels in addition to *in situ* hybridization to anchor most or all of the X

783 chromosome sequences. In total, ~137-138 Mb of sequence, or ~49% of the expected genome

784 size, was anchored into chromosomes. Future research with targeted approaches to study the L

785 chromosome and variations associated with the X' chromosome will be of interest beyond the

786 current male somatic genome assembly presented here.

787

42

788      In its current state, the *Sciara* genome assembly is already more contiguous than up to 95%

789      of all Arthropod genomes described (http://i5k.github.io/arthropod_genomes_at_ncbi). Its

790      contiguity statistics exceed 42 of the 43 currently available lower Dipteran genome assemblies,

791      over 75% of which have sub-100 kb N50s. The low contiguity of most available Dipteran genome

792      sequences and the lack of anchoring to chromosomes limits their utility. However, the *Sciara*

793      genome assembly presented here may be useful for scaffolding currently available and future

794      *Nematoceran* genomes by synteny. The long contigs in the *Sciara* genome assembly reflect the

795      success of using long read technologies and optical maps, both of which span repeats. The long-

796      read datasets and the resulting assembly will be important and extremely useful for analyzing

797      regions of repetitive DNA, like rDNA, centromeres, telomeres, and transposable elements.

798

799      **<u>Comparative phylogenomics</u>**

800      Comparative genomics provides an understanding into the rates and patterns of evolution

801      of genes as well as populations and species (Wiegmann and Richards 2018). The phylogenetic

802      position of *Sciara (Bradysia) coprophila* makes its genome and transcriptome sequences valuable

803      for future comparative genomics studies. *Sciara* is a lower Dipteran fly (Nematocera) whereas

804      *Drosophila* is a higher Dipteran fly (Brachycera) and they diverged from one another ~200 MYA

805      (Wiegmann et al 2011). The *Sciara* genome size of 280 Mb (362 Mb with the L chromosomes) is

806      larger than the 175 Mb size of the *Drosophila melanogaster* genome (Elllis et al 2014), but similar

807      to the 264 Mb genome of the Nematoceran *Anopholes gambia* (Sharakhova et al 2007). Dipteran

808      phylogenetics has been much studied (Hennig 1973; McAlpine and Wood 1989) but some

809      unresolved questions remain. Previously, morphological criteria suggested that the Brachycera

810      (containing *Drosophila*) and the Nematocera (containing *Sciara*) diverged from a common

811      ancestor. However, more recent molecular data supports a model where the Nematoceran

812      infraorder Bibionomorpha ultimately gave rise to the Brachycera (Wiegmann et al 2011). The

813      *Sciara* genome and transcriptome sequences reported here will be valuable resources to further

814    describe Dipteran phylogenetic relationships, and will further our understanding of the evolution

815    and molecular structure of genes and pathways in Dipterans including Drosophila.

816

817    **<u>Evolution of sex determination</u>**

818        The evolution of sex determination is a topic of much current interest. The most common

819    occurrence is male heterogamety where males are XY and females are XX. In contrast, in female

820    heterogamety, females have heteromorphic sex chromosomes (e.g., ZW), and males are

821    homomorphic (e.g., ZZ). Female heterogamety is rare in insects (Blackmon et al 2017), but is

822    exhibited by *Sciara* where males have a single X in their soma and females have two (Gerbi

823    1986). Female *Sciara* can be either XX or X'X where the X' chromosome carries a long paracentric

824    inversion that inhibits crossing over with the X. Thus, the heterogametic *Sciara* female determines

825    the sex of her offspring. In *Sciara coprophila*, XX mothers have only sons and X'X mothers have

826    only daughters. Presumably, the ooplasm is conditioned by the *Sciara* mother to determine the

827    sex of the offspring via X chromosome elimination. In agreement with this hypothesis, sex is

828    determined by a temperature-sensitive maternal effect that controls X-chromosome elimination in

829    *Sciara ocellaris* (Nigro et al. 2007). As for the single X in male soma, *Sciara* males are haploid

830    only for the X but diploid for the autosomes, unlike haplodiploid males that are haploid for their

831    entire genome. This is accomplished by X chromosome elimination in the early *Sciara* embryo

832    and was noted by White (1949) to occur in the Nematoceran families of Sciaridae and

833    Cecidomyidae (including the Hessian fly *Mayetiola destructor*). Comparisons of the

834    genomes/transcriptomes of *Sciara* and *M. destructor* might help to elucidate the molecular

835    regulation of X chromosome elimination.

836

837        Cytoplasmic sex determination, as suggested above for *Sciara*, occurs if sex is under the

838    control of cytoplasmic elements, such as endosymbionts. *Wolbachia* and *Rickettsia* are related

839    groups of intracellular alpha proteobacteria that can distort the sex ratio of their arthropod hosts

44

840    (Lawson et al, 2001, Serbus et al 2008). They are transmitted through the egg cytoplasm and

841    alter reproduction in their arthropod hosts in various ways, including cytoplasmic incompatibility,

842    feminization of genetic males, and male killing (Werren and Windsor 2000; Serbus et al 2008).

843    Both can induce parthogenesis (Blackmon et al 2017). The latter is of interest since (i)

844    parthenogenetic *Sciara* embryos have been observed, but their development arrests in

845    embryogenesis (de Saint Phalle and Sullivan 1998), and (ii) although we did not find *Wolbachia*

846    sequences in *Sciara* genomic DNA, we essentially co-assembled an entire *Rickettsia* genome.

847    Moreover, our genomic copy number analyses suggest there are two *Rickettsia* cells per *Sciara*

848    cell on average in 1-2 day old male embryos. Further evidence is needed to ascertain if *Rickettsia*

849    plays a role in *Sciara* sex determination.

850

851    **<u>Paternal chromosome imprinting</u>**

852        The first example of a chromosome or a chromosomal locus "remembering" its maternal

853    or paternal origin was noted in *Sciara* and the term "imprinting" was coined (Crouse 1960).

854    Specifically, in *Sciara* male meiosis I, the paternally derived chromosomes move away from the

855    single pole of the naturally occurring monopolar spindle and are discarded in a bud of cytoplasm.

856    This is an example of paternal genome elimination (PGE) that can give rise to haplodiploidy in

857    other systems (Blackmon et al 2017). Thus, only the maternal genome is passed down through

858    sperm in *Sciara*. Although sperm in *Sciara* is haploid for autosomes, it is diploid for the X

859    chromosomes due to non-disjunction of the X in *Sciara* male meiosis II (Gerbi 1986). After

860    fertilization of the haploid egg, diploidy is re-established for the autosomes, but the X chromosome

861    is temporarily triploid. Either one or both copies of the paternally-derived X are eliminated from

862    female or male embryos, respectively, during the 7th-9th embryonic cleavage division, representing

863    another example of imprinting in *Sciara* (de Saint Phalle and Sullivan 1996). Nevertheless, the

864    mechanism for imprinting in *Sciara* remains elusive. It is of interest to learn if DNA modifications

865    occur in *Sciara* since different imprints in mammalian genomes are laid down in eggs and sperm

866    through a DNA methylation mechanism, leading to differential gene expression at imprinted loci

867    in the offspring (Li et al 1993).

868

869        DNA methylation typically occurs at CpG sites where it is established *de novo* by DNA

870    methyltransferase 3 (DNMT3) and is maintained by DNMT1 (Goll and Bestor 2005, Kato et al

871    2007). In contrast to vertebrates, DNA methylation in invertebrates is relatively sparse (Bird 1980).

872    DNMT1 is found in all orders of insects except Diptera, which also lack DNMT3 (Bewick et al

873    2017). In agreement, our gene annotations suggest that *Sciara* also lacks DNMT1 and DNMT3.

874    Some bisulfite sequencing studies revealed that CpG DNA methylation is found in all insect

875    Orders except Dipteran flies (Bewick et al 2017) and failed to find specific patterns for methylated

876    C in *Drosophila* embryos (Zemach et al 2010; Raddatz et al 2013). Other studies have asserted

877    that *Drosophila melanogaster* has DNA methyltransferase activity and CpC methylation (Panikar

878    et al 2015), has low levels of 5-methylcytosine (5mC) (Capuano et al 2014, Takayama et al 2014,

879    Deshmukh et al. 2018), and has more cytosine methylation in stage 5 *Drosophila* embryos than

880    oocytes (Takayama et al 2014). Moreover, 6-methyladenine (6mA) has been recently reported to

881    be in the genomic DNA of *Drosophila* and other eukaryotes (Fu et al. 2015; Greer et al. 2015;

882    Zhang et al. 2015; Liu et al. 2016; Mondo et al. 2017; Wang et al. 2018; Xiao et al. 2018). Typically,

883    the level of 6mA is quite low, such as 0.001% in *Drosophila* but rises to 0.07% in early embryos

884    (Zhang et al 2015). DAMT-1 appears to be the methyltransferase for 6mA in insect cells and

885    DMAD has 6mA demethylating activity in *Drosophila* (Luo et al 2015, Zhang et al 2015). Our gene

886    annotations suggest that *Sciara* has both DAMT-1 and DMAD.

887

888        Before it can determined whether or not imprinting in *Sciara* involves DNA modifications,

889    it needs to be determined if the *Sciara* genome harbors DNA modifications at all. Previous

890    immunofluorescence studies have suggested the presence of 5-methylcytosine in *Sciara*

891    chromosomes (Eastman 1980, Greciano 2009). Similarly, our sequencing data support the

892    presence of base modifications in the *Sciara* genome. Overall, up to 0.6-1.1% of cytosines may

893    be modified in the *Sciara* genome, especially at GCG sites, with specifically 0.1-0.2% and 0.3-

894    0.4% identified as 4mC and 5mC, respectively. In addition, 0.13-0.24% of adenine sites in the

895    *Sciara* male embryo genome were potentially modified with up to ~0.04-0.06% of adenine sites

896    containing 6mA, especially GAG sites. Moreover, the PacBio analysis suggests that these DNA

897    modifications are phased with 10 bp and 175 bp periodicities, suggesting physical interactions

898    between the 10 bp turns of the DNA helix and methylation machinery as well as relationships with

899    nucleosome spacing, both of which have been seen previously using orthogonal methods (Jia et

900    al. 2007; Chodavarapu et al. 2011; Fu et al. 2015; Collings and Anderson 2017; Wang et al. 2017;

901    Luo et al. 2018). Lastly, the distribution of modifications we observed with respect to genes and

902    repeats are concordant with previous observations (i) of methyl-C in *Drosophila* (Takayama et al.

903    2014) and (ii) that heterochromatic regions of the *Sciara* genome, where most repeats reside, are

904    enriched for 5mC (Eastman et al. 1980; Greciano et al. 2009). Overall, the evidence from single-

905    molecule sequencing lends support to the presence of methylated cytosines and adenines in the

906    autosomes and X chromosome in the male embryo genome of *Sciara*. However, the analyses

907    suggest that the levels of DNA modifications are low. Their abundance in females and other

908    developmental stages and tissues as well as their biological significance remains to be

909    determined in future investigations. Nevertheless, given the evidence from the current study and

910    previous work, base modifications may be a promising avenue for the study of imprinting in *Sciara*.

911

912    **Summary**

913    The *Sciara* genome sequence provides a foundation for future studies to delve into the

914    many unique biological properties of *Sciara* (reviewed by Gerbi 1986) that include **(i)** chromosome

915    imprinting; **(ii)** sex determination by the mother; **(iii)** a monopolar spindle in male meiosis I; **(iv)**

916    non-disjunction of the X chromosome in male meiosis II; **(v)** chromosome elimination in early

917     embryogenesis; **(vi)** germ line limited L chromosomes; **(vii)** DNA amplification in late larval

918     salivary gland polytene chromosomes; **(viii)** high resistance to radiation.

919 **METHODS**

920

921 **Tissue collection, DNA extraction, and DNA sequencing:**

922 *Sciara coprophila* (renamed *Bradysia coprophila)* was used for these studies. *Sciara*

923 (stock: Holo2) matings were performed to produce only male offspring from which embryos aged

924 2 hours – 2 days (genome sequencing datasets) or pupae (BioNano Irys genome mapping

925 datasets) were collected. For a minority of MinION sequencing data, adult males were used.

926 Genomic DNA (gDNA) was isolated using DNAzol (ThermoFisher) as per the manufacturer's

927 instructions with some modifications. gDNA was cleaned with AMPure beads (Beckman Coulter).

928 Purity was checked with NanoDrop (ThermoFisher) and concentration was checked with Qubit

929 (ThermoFisher).

930

931 For Illumina HiSeq 2000 sequencing, male embryo gDNA was sonicated to a size range

932 of 100-600 bp, prepared using the NEBNext kit (New England Biolabs) following the

933 manufacturer's directions, run on a 2% NuSieve agarose (Lonza) gel, size-selected near the 500

934 bp marker, gel purified (Qiagen), and sequenced to obtain 100 bp paired-end reads.

935

936 For Pacific Biosciences RSII Single Molecule Real Time sequencing datasets (P5-C3

937 chemistry), male embryo gDNA was given to the Technology Development Group at the Institute

938 of Genomics and Multiscale Biology at the Icahn School of Medicine at Mount Sinai for library

939 construction and sequencing. Two DNA libraries were prepared and sequenced across 24

940 SMRTcells as described further in the Supplemental Methods.

941

942 MinION data was collected using multiple early iterations of the technology (original

943 MinION and MkI), kits (SQK-MAP002, MAP004, MAP005, MAP006), and pores (R7.3 and R7.3

944 70 bps 6mer). We prepared 15 libraries from male *Sciara* embryo gDNA (making up >97% of the

49

945    data) and 2 from male adult gDNA. The manufacturer's instructions were followed with

946    modifications to increase read lengths (Urban et al. 2015 and Suppl. Methods). Libraries were

947    loaded onto the MinION, sequenced, and basecalled with Metrichor. Reads were extracted from

948    Fast5 files and analyzed using our own custom set of tools (Fast5Tools:

949    github.com/JohnUrban/fast5tools).

950

951        For BioNano Genomics (BNG) Irys optical maps, male pupae were flash frozen ground in

952    liquid nitrogen and high molecular weight gDNA was isolated (Suppl. Methods), nicked with BssSI

953    (CACGAG, New England BioLabs), labeled, and repaired according to the IrysPrep protocol

954    (BioNano Genomics).

955

956    **Genome assemblies**

957        After optional trimming/filtering with Trimmomatic (Bolger et al 2014) and/or error-

958    correction with BayesHammer (Nikolenko et al 2013), short-read assemblies were generated

959    using ABySS (Simpson et al. 2009), Megahit (Li et al. 2015), Platanus (Kajitani et al. 2014), SGA

960    (Simpson and Durbin 2010), SOAP (Luo et al. 2012), SPAdes (Bankevich et al. 2012), Velvet

961    (Zerbino and Birney 2008). Hybrid assemblies were generated using DBG2OLC (Ye et al. 2016)

962    and PBDagCon (http://bit.ly/pbdagcon) starting with Platanus contigs and long reads. Non-hybrid

963    long-read assemblies were generated with Canu (Koren et al. 2017), Falcon (Chin et al. 2016),

964    Miniasm (Li 2016) with RaCon (Vaser et al. 2017), ABruijn (Lin et al. 2016), and SMARTdenovo

965    (https://github.com/ruanjue/smartdenovo). For all assemblers, we varied filtering, error correction,

966    inputs, and parameters as detailed further in the Suppl. Methods. Long-read assemblies were

967    polished with Quiver (Chin et al. 2013) and Pilon (Walker et al. 2014). BlobTools (Laetsch and

968    Blaxter 2017) was used to identify contaminating contigs.

969

970    **Assembly evaluations**

50

971    Assembly evaluations included subsets of the following: contig size statistics, percent of

972    Illumina reads that mapped using Bowtie2 (Langmead and Salzberg 2012), probabilistic scores

973    from LAP (Ghodsi et al. 2013) and ALE (Clark et al. 2013), number of features from FRC[bam] (Vezzi

974    et al. 2012), percent error-free bases and/or the mean base score from REAPR (Hunt et al. 2013),

975    completeness of gene content with BUSCO (Simão et al. 2015), the percent of long reads that

976    aligned with BWA (Li and Durbin 2009), the average number of split alignments per long read,

977    structural variations using Sniffles (Sedlazeck et al. 2018), the percent of raw BioNano map

978    alignments using Maligner (Mendelowitz et al. 2015), the resulting optical map alignment M-

979    scores, the number of bases covered by optical maps (span), and the total coverage from aligned

980    optical maps. Evaluations were automated and parallelized on SLURM with a custom package

981    (github.com/JohnUrban/battery).

982

983    **Scaffolding**

984    For hybrid scaffolding, optical maps >150 kb were assembled into consensus maps

985    (CMAPs) using BioNano Pipeline Version 2884 and RefAligner Version 2816 (BioNano

986    Genomics). Each selected assembly was used with the BNG CMAPs to create genome-wide

987    hybrid scaffolds using hybridScaffold.pl version 4741 (BioNano Genomics). Quiver and PBJelly

988    (English et al. 2012) were used to polish and gap-fill the scaffolds. PBJelly was used additionally

989    to join more scaffolds with long-read evidence into "meta-scaffolds", and Quiver and Pilon were

990    used for final polishing.

991

992    **Assembly anchoring**

993    Haplotigs were identified using Minimap2 (Li 2018) and purge_haplotigs (Roach et al.

994    2018). To anchor contigs into chromosomes, sequences that were previously mapped to

995    chromosomes experimentally were mapped to the assemblies using BLAST (Altschul et al. 1990).

996    Differentiating between autosomal and X-linked contigs was performed by requiring haploid

997    coverage levels across at least 80% of a contig to be called as X-linked (else autosomal), using

998    Minimap2 and BEDTools (Quinlan and Hall 2010).

999

1000    **Transcriptome assemblies**

1001    For strand-specific RNA-sequencing libraries, poly-A RNA was prepared from a given sex

1002    and stage using TRIzol (Invitrogen/Thermofisher), DNase (Qiagen), RNeasy columns (Qiagen),

1003    and Oligo-dT DynaBeads (Life Technologies). RNA integrity was evaluated on 1.1%

1004    formaldehyde 1.2% agarose gels. RNA purity and quantity were measured with the NanoDrop

1005    (ThermoScientific) and Qubit (ThermoFisher) throughout. Libraries were prepared from poly-A

1006    RNA using NEB's Magnesium Fragmentation Module, SSIII (Invitrogen) first strand synthesis with

1007    random primers, NEBNext Second Strand Synthesis module with ACGU nucleotide mix (10 mM

1008    each of dATP, dCTP, dGTP, and 20 mM of dUTP), NEBNext End Repair and dA-Tailing (NEB),

1009    and ligation (NEB: NEBNext Quick Ligation Reaction Buffer, NEB Adaptor, Quick T4 Ligase). The

1010    libraries were size-selected with AMPure beads (Beckman Coulter). Uracil-cutting for strand-

1011    specificity (and hairpin adapter cutting) was performed with NEBNext USER enzyme, followed by

1012    PCR using NEBNext High-Fidelity 2X PCR Master Mix and NEBNext indexed and universal

1013    primers for 12 cycles. A final size-selection of PCR products was performed with AMPure beads.

1014    Purity, quantity, and size of the libraries were checked with NanoDrop, Qubit and Fragment

1015    Analyzer (Agilent). Traces suggested the mean estimated fragment sizes was around 420 bp,

1016    indicating mean insert sizes near 300 bp. Libraries were sequenced to yield 100 bp paired-end

1017    reads using the Illumina HiSeq 2000. The strand-specific RNA-seq datasets were combined and

1018    assembled with Trinity (Grabherr et al. 2011) or using HiSat2 (Kim et al. 2019) and StringTie

1019    (Pertea et al. 2015). Transcriptome assemblies were evaluated with BUSCO (Simão et al. 2015),

1020    RSEM-Eval (Li et al. 2014), and TransRate (Smith-Unna et al. 2016).

1021

1022    **Repeat and gene annotation**

1023    Species-specific repeat libraries were built using RepeatModeler (Smit and Hubley 2008).

1024    These were combined with previously known repeat sequences from *Bradysia coprophila* as well

1025    as all Arthropod repeats in the RepeatMasker Combined Database: Dfam_Consensus-20181026

1026    (Hubley et al. 2016), RepBase-20181026 (Bao et al. 2015). To predict protein-coding genes,

1027    Maker2 (Holt and Yandell 2011) was used with transcriptome evidence described above,

1028    transcript and protein sequences from related species for homology evidence, Augustus (Hoff

1029    and Stanke 2019), SNAP (Korf 2004), GeneMark-ES (Ter-Hovhannisyan et al. 2008), and

1030    RepeatMasker (Smit et al. 2013) with repeat libraries described above. InterProScan (Quevillon

1031    et al. 2005) was used to identify Pfam domains and GO terms from predicted protein sequences,

1032    and BLASTp was to find the best matches to curated proteins in the entire UniProtKB/Swiss-Prot

1033    database (The UniProt Consortium 2019). Maker2 transcriptomes were evaluated using

1034    annotation edit distances, BUSCO, RSEM-Eval, and TransRate.

1035

1036    **DNA modification analyses**

1037    PBalign (github.com/PacificBiosciences/pbalign) with BLASR v2 (Chaisson and Tesler

1038    2012) was used to align PacBio reads to the entire unfiltered assembly to avoid forcing incorrect

1039    mappings. Pbh5tools (github.com/PacificBiosciences/pbh5tools) was used to merge and sort the

1040    mapped        reads.        ipdSummary        from        kineticsTools        v0.6.0

1041    (github.com/PacificBiosciences/kineticsTools) was used to predict base modifications across the

1042    Canu genome assembly (--pvalue 0.01 --minCoverage 3 --methylMinCov 10 --identifyMinCov 5).

1043    AgIn (Suzuki et al. 2016) was also used to look at CpG methylation. For all analyses on predicted

1044    DNA modifications, we used only primary contigs labeled as Arthopoda. Kmer enrichment scores

1045    for dimers and trimers were obtained from the Chi-square standardized residuals found when

1046    comparing the distribution of kmers that had a specific modification at a fixed position with the

1047    genome-wide distribution of kmers with the target base at that position. We also used this

1048    approach to define enriched 7-mers for position weight matrix motifs using WebLogo (Crooks et

1049 al. 2004). In addition, the 9 bp sequences centered on the top 500 or 5000 scoring specific

1050 modification calls were used with MEME (Bailey and Elkan 1994) to identify motifs using a second

1051 order Markov model background file trained on the *Sciara* genome assembly (fasta-get-markov -

1052 m 2 -dna). We determined if DNA modifications were enriched/depleted in various genomic

1053 regions using binomial models. When separating genes by expression level for this analysis,

1054 Salmon (Patro et al. 2017) was used to quantify expression over our Maker2 protein-coding gene

1055 annotation using male embryo RNA-seq. BEDtools was used to obtain spacing distances between

1056 modified bases as well as between random bases of the same type (e.g. m6A vs random A).

1057 Although 10 bp periodicities were obvious by visual examination, we formally determined the

1058 periodicities observed in counts of inter-modification distances between 0-200 bp by running a

1059 discrete Fourier transform (DFT) analysis using the Fast Fourier Transform (FFT) from Python's

1060 Numpy package.

1061

1062 For the MinION analysis, only datasets generated from the MkI, SQK-MAP006 kit, and

1063 R7.3 70 bps 6mer pore model were used, and only reads that aligned to primary contigs annotated

1064 as Arthropoda. We compared the signal distributions for each kmer in our *Sciara* dataset to the

1065 expected ONT kmer models, and to a MinION dataset generated from whole genome PCR on *E.*

1066 *coli* genomic DNA using the same kit and pore model (BioProject PRJEB13021; Run

1067 ERR1309547; www.ebi.ac.uk/ena; Simpson et al. 2017). MinION reads were aligned with BWA.

1068 Nanopolish (Simpson et al. 2017) was used to learn updated kmer models from the native *Sciara*

1069 and E. coli PCR MinION datasets. MEME was used to identify short motifs in all 6mers that

1070 differed from the expected ONT model.

1071

1072 **<u>Further bioinformatics</u>**

1073 The Supplemental Methods contains software versions, as well as further details and

1074 exact commands for: read processing, genome assembling, polishing, evaluating, scaffolding,

54

1075    gap filling, bacterial filtering, haplotig filtering, anchoring, transcriptome assemblies and

1076    evaluations, repeat library construction, repeat-masking, training gene predictors, alternative

1077    transcript and protein evidence, Maker2 iterations and evaluations, and the PacBio and MinION

1078    DNA modification analyses. Bioinformatics analyses were largely aided by custom scripts located

1079    at          github.com/JohnUrban/sciara-project-tools,          github.com/JohnUrban/fast5tools,

1080    github.com/JohnUrban/battery,                    github.com/JohnUrban/lave,                    and

1081    github.com/JohnUrban/fftDnaMods.

1082

1083                                              **DATA ACCESS**

1084    Raw Illumina (DNA and RNA-seq), PacBio, MinION, and BioNano data generated in this study as

1085    well as BioNano CMAPs and PacBio kinetics and DNA modification results have been submitted

1086    to the NCBI BioProject database (http://www.ncbi.nlm.nih.gov/bioproject/) under accession

1087    number PRJNA123456. This Whole Genome Shotgun project has been deposited at

1088    DDBJ/ENA/GenBank under the accession VSDI00000000, and the Canu assembly version

1089    selected as the first draft genome release in this paper (Bcop v1.0) is version VSDI01000000.

1090    The automated Bcop_v1.0 annotation for the Canu assembly is available at the i5k Workspace

1091    (i5k.nal.usda.gov) where manual curation updates will be made.

1092

1093                                    **DISCLOSURE DECLARATION**

1094    JMU and SAG were members of the MinION Access Program and received free reagents from

1095    ONT. JMU was also a member of the MinION Access and Reference Consortium (MARC) that

1096    conducts experiments partially funded by ONT.

1097

1098                                       **ACKNOWLEDGMENTS**

**AUTHOR CONTRIBUTIONS**

John Urban (JMU) collected all embryos, larvae, pupae, and adult Sciara needed for all experiments. JMU prepared all MinION libraries and performed all MinION sequencing and

56

1127    analyses. JMU wrote the suites of tools for working with MinION data

1128    (https://github.com/JohnUrban/fast5tools), automating the battery of assembly evaluations

1129    (https://github.com/JohnUrban/battery), genome alignment visualizations

1130    (https://github.com/JohnUrban/lave), and all general bioinformatics over the course of this project

1131    (https://github.com/JohnUrban/sciara-project-tools). JMU obtained high molecular weight

1132    genomic DNA and delivered it to the Technology Development Group at the Institute of Genomics

1133    & Multiscale Biology at the Icahn School of Medicine at Mount Sinai, where PacBio sequencing

1134    libraries were prepared and sequenced. JMU performed all short- and long-read assemblies,

1135    genome polishing, assembly evaluations, repeat modeling and annotation, and gene annotation.

1136    JMU did all RNA work and library preparations for all RNA-seq samples representing replicates

1137    from both sexes at different stages, and performed all transcriptome assemblies and RNA-seq

1138    data analysis. JMU performed DNA modification analyses with PacBio single molecule kinetics

1139    data and MinION single molecule ionic current data. CMC made DNA plugs from *Sciara* pupae

1140    collected and sent to her by JMU, and performed the BioNano preparations and imaging on Irys

1141    platform. RM and NL performed BioNano hybrid scaffolding with selected assemblies sent to

1142    them. SJB provided guidance in our acquisition of BioNano data and provided oversight to CMC,

1143    RM, and NL. MSF prepared the Illumina DNA library. JEB did all *Sciara* mass matings. ACS

1144    provided support and guidance on this work. SAG pioneered and guided the *Sciara* genome effort.

1145    JMU conceived the experiments and analyses. JMU and SAG wrote the manuscript.

**REFERENCES**

Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol* **215**: 403–410.

Andersson SGE, Zomorodipour A, Andersson JO, Sicheritz-Pontén T, Alsmark UCM, Podowski RM, Näslund AK, Eriksson A-S, Winkler HH, Kurland CG. 1998. The genome sequence of Rickettsia prowazekii and the origin of mitochondria. *Nature* **396**: 133–140.

Armstrong MJ, Jin Y, Allen EG, Jin P. 2019. Diverse and dynamic DNA modifications in brain and diseases. *Hum Mol Genet* **28**: R241–R253.

Bailey TL, Elkan C. 1994. Fitting a mixture model by expectation maximization to discover motifs in biopolymers. *Proc Int Conf Intell Syst Mol Biol* **2**: 28–36.

Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Prjibelski AD, et al. 2012. SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *J Comput Biol* **19**: 455–477.

Bao W, Kojima KK, Kohany O. 2015. Repbase Update, a database of repetitive elements in eukaryotic genomes. *Mob DNA* **6**: 11.

Bewick AJ, Vogel KJ, Moore AJ, Schmitz RJ. 2017. Evolution of DNA methylation across insects. *Mol Biol Evol* **34**: 654–665.

Bienz-Tadmor B, Smith HS, Gerbi SA. 1991. The promoter of DNA puff gene II/9-1 of Sciara coprophila is inducible by ecdysone in late prepupal salivary glands of Drosophila melanogaster. *Cell Regul* **2**: 875–88.

Bird AP. 1980. DNA methylation and the frequency of CpG in animal DNA. *Nucleic Acids Res* **8**: 1499–1504.

Blackmon H, Ross L, Bachtrog D. 2017. Sex determination, sex chromosomes, and karyotype evolution in insects. In *Journal of Heredity*, Vol. 108 of, pp. 78–93, Oxford University Press.

Boivin A, Vendrely R, C Vendrely. 1948. L'acide désoxyribonucléique du noyau cellulaire, dépositaire des caractères héréditaires; arguments d'ordre analytique. *Comp trend Acad Sci* **226**: 1061–1063.

Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**: 2114–20.

Capuano F, Mülleder M, Kok R, Blom HJ, Ralser M. 2014. Cytosine DNA Methylation Is Found in Drosophila melanogaster but Absent in Saccharomyces cerevisiae, Schizosaccharomyces pombe, and Other Yeast Species. *Anal Chem* **86**: 3697–3702.

Chaisson MJ, Tesler G. 2012. Mapping single molecule sequencing reads using basic local alignment with successive refinement (BLASR): application and theory. *BMC Bioinformatics* **13**: 238.

Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, Clum A, Copeland A, Huddleston J, Eichler EE, et al. 2013. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods* **10**: 563–9.

Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R, Figueroa-Balderas R, Morales-Cruz A, et al. 2016. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat Methods*.

Chodavarapu RK, Feng S, Bernatavichute Y V., Chen PY, Stroud H, Yu Y, Hetzel JA, Kuo F, Kim J, Cokus SJ, et al. 2010. Relationship between nucleosome positioning and DNA methylation.

1189    *Nature* **466**: 388–392.

1190    Clark ME, Anderson CL, Cande J, Karr TL. 2005. Widespread prevalence of Wolbachia in
1191        laboratory stocks and the implications for Drosophila research. *Genetics* **170**: 1667–1675.

1192    Clark SC, Egan R, Frazier PI, Wang Z. 2013. ALE: a generic assembly likelihood evaluation
1193        framework for assessing the accuracy of genome and metagenome assemblies.
1194        *Bioinformatics* **29**: 435–43.

1195    Clark TA, Murray IA, Morgan RD, Kislyuk AO, Spittle KE, Boitano M, Fomenkov A, Roberts RJ,
1196        Korlach J. 2012. Characterization of DNA methyltransferase specificities using single-
1197        molecule, real-time DNA sequencing. *Nucleic Acids Res* **40**: e29.

1198    Collings CK, Anderson JN. 2017. Links between DNA methylation and nucleosome occupancy in
1199        the human genome. *Epigenetics and Chromatin* **10**.

1200    Crooks GE, Hon G, Chandonia J-M, Brenner SE. 2004. WebLogo: a sequence logo generator.
1201        *Genome Res* **14**: 1188–90.

1202    Crouse HV, Brown A, Mumford BC. 1971. L-chromosome inheritance and the problem of
1203        chromosome "imprinting" in Sciara (Sciaridae, Diptera). *Chromosoma* **34**: 324–339.

1204    Crouse H V., Gerbi SA, Liang CM, Magnus L, Mercer IM. 1977. Localization of ribosomal DNA
1205        within the proximal X heterochromatin of Sciara coprophila (Diptera, Sciaridae).
1206        *Chromosoma* **64**: 305–318.

1207    Crouse H V. 1960. The Controlling Element in Sex Chromosome Behavior in Sciara. *Genetics* **45**:
1208        1429–43.

1209    de Saint Phalle B, Sullivan W. 1996. Incomplete sister chromatid separation is the mechanism of
1210        programmed chromosome elimination during early Sciara coprophila embryogenesis.
1211        *Development* **122**: 3775–84.

1212    de Saint Phalle B, Sullivan W. 1998. Spindle assembly and mitosis without centrosomes in
1213        parthenogenetic Sciara embryos. *J Cell Biol* **141**: 1383–91.

1214    Deshmukh S, Ponnaluri VKC, Dai N, Pradhan S, Deobagkar D. 2018. Levels of DNA cytosine
1215        methylation in the Drosophila genome. *PeerJ* **2018**.

1216    DiBartolomeis SM, Gerbi SA. 1989. Molecular characterization of DNA puff II/9A genes in Sciara
1217        coprophila. *J Mol Biol* **210**: 531–40.

1218    Eastman EM, Goodman RM, Erlanger BF, Miller OJ. 1980. 5-Methylcytosine in the DNA of the
1219        polytene chromosomes of the dipteraSciara coprophila,Drosophila melanogaster andD.
1220        persimilis. *Chromosoma* **79**: 225–239.

1221    Eid J, Fehr A, Gray J, Luong K, Lyle J, Otto G, Peluso P, Rank D, Baybayan P, Bettman B, et al.
1222        2009. Real-Time DNA Sequencing from Single Polymerase Molecules. *Science (80- )* **323**.

1223    El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, Qureshi M, Richardson LJ, Salazar
1224        GA, Smart A, et al. 2019. The Pfam protein families database in 2019. *Nucleic Acids Res* **47**:
1225        D427–D432.

1226    English AC, Richards S, Han Y, Wang M, Vee V, Qu J, Qin X, Muzny DM, Reid JG, Worley KC, et al.
1227        2012. Mind the gap: upgrading genomes with Pacific Biosciences RS long-read sequencing
1228        technology. *PLoS One* **7**: e47768.

1229    Escribá MC, Greciano PG, Méndez-Lago M, De Pablos B, Trifonov VA, Ferguson-Smith MA,
1230        Goday C, Villasante A. 2011. Molecular and cytological characterization of repetitive DNA
1231        sequences from the centromeric heterochromatin of Sciara coprophila. *Chromosoma* **120**:
1232        387–397.

1233    Felsheim RF, Kurtti TJ, Munderloh UG. 2009. Genome Sequence of the Endosymbiont Rickettsia
1234        peacockii and Comparison with Virulent Rickettsia rickettsii: Identification of Virulence
1235        Factors ed. N. Ahmed. *PLoS One* **4**: e8361.
1236    Flusberg BA, Webster DR, Lee JH, Travers KJ, Olivares EC, Clark TA, Korlach J, Turner SW. 2010.
1237        Direct detection of DNA methylation during single-molecule, real-time sequencing. *Nat*
1238        *Methods* **7**: 461–465.
1239    Foulk MS, Liang C, Wu N, Blitzblau HG, Smith H, Alam D, Batra M, Gerbi SA. 2006. Ecdysone
1240        induces transcription and amplification in Sciara coprophila DNA puff II/9A. *Dev Biol* **299**:
1241        151–63.
1242    Foulk MS, Waggener JM, Johnson JM, Yamamoto Y, Liew GM, Urnov FD, Young Y, Lee G, Smith
1243        HS, Gerbi SA. 2013. Isolation and characterization of the ecdysone receptor and its
1244        heterodimeric partner ultraspiracle through development in Sciara coprophila.
1245        *Chromosoma* **122**: 103–19.
1246    Fu Y, Luo G-Z, Chen K, Deng X, Yu M, Han D, Hao Z, Liu J, Lu X, Doré LC, et al. 2015. N6-
1247        Methyldeoxyadenosine Marks Active Transcription Start Sites in Chlamydomonas. *Cell* **161**:
1248        879–892.
1249    Gabrusewycz-Garcia N. 1964. Cytological and autoradiographic studies in Sciara coprophila
1250        salivary gland chromosomes. *Chromosoma* **15**: 312–44.
1251    Gerbi SA. 1971. Localization and characterization of the ribosomal RNA cistrons in Sciara
1252        coprophila. *J Mol Biol* **58**: 499–511.
1253    Gerbi SA. 1986. Unusual chromosome movements in sciarid flies. *Results Probl Cell Differ* **13**:
1254        71–104.
1255    Gerbi SA, Strezoska Z, Waggener JM. 2002. Initiation of DNA replication in multicellular
1256        eukaryotes. *J Struct Biol* **140**: 17–30.
1257    Ghodsi M, Hill CM, Astrovskaya I, Lin H, Sommer DD, Koren S, Pop M. 2013. De novo likelihood-
1258        based measures for comparing genome assemblies. *BMC Res Notes* **6**: 334.
1259    Goll MG, Bestor TH. 2005. Eukaryotic Cytosine Methyltransferases. *Annu Rev Biochem* **74**: 481–
1260        514.
1261    Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
1262        Raychowdhury R, Zeng Q, et al. 2011. Full-length transcriptome assembly from RNA-Seq
1263        data without a reference genome. *Nat Biotechnol* **29**: 644–52.
1264    Greciano PG, Ruiz MF, Kremer L, Goday C. 2009. Two new chromodomain-containing proteins
1265        that associate with heterochromatin in Sciara coprophila chromosomes. *Chromosoma* **118**:
1266        361–376.
1267    Greer EL, Blanco MA, Gu L, Sendinc E, Liu J, Aristizábal-Corrales D, Hsu C-H, Aravind L, He C, Shi
1268        Y. 2015. DNA Methylation on N6-Adenine in C. elegans. *Cell* **161**: 868–878.
1269    Hennig W. 1973. Diptera (Two-winged Flies). *Handb Zool*.
1270    Hoff KJ, Stanke M. 2019. Predicting Genes in Single Genomes with AUGUSTUS. *Curr Protoc*
1271        *Bioinforma* **65**: e57.
1272    Holt C, Yandell M. 2011. MAKER2: an annotation pipeline and genome-database management
1273        tool for second-generation genome projects. *BMC Bioinformatics* **12**: 491.
1274    Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, Bao W, Smit AFA, Wheeler TJ. 2016. The Dfam
1275        database of repetitive DNA families. *Nucleic Acids Res* **44**: D81–D89.
1276    Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M, Otto TD. 2013. REAPR: a universal tool

1277      for genome assembly evaluation. *Genome Biol* **14**: R47.

1278    Ip CLC, Loose M, Tyson JR, de Cesare M, Brown BL, Jain M, Leggett RM, Eccles DA, Zalunin V,
1279      Urban JM, et al. 2015. MinION Analysis and Reference Consortium: Phase 1 data release
1280      and analysis. *F1000Research* **4**.

1281    Jia D, Jurkowska RZ, Zhang X, Jeltsch A, Cheng X. 2007. Structure of Dnmt3a bound to Dnmt3L
1282      suggests a model for de novo DNA methylation. *Nature* **449**: 248–251.

1283    Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M,
1284      Nagayasu E, Maruyama H, et al. 2014. Efficient de novo assembly of highly heterozygous
1285      genomes from whole-genome shotgun short reads. *Genome Res* **24**: 1384–95.

1286    Kato Y, Kaneda M, Hata K, Kumaki K, Hisano M, Kohara Y, Okano M, Li E, Nozaki M, Sasaki H.
1287      2007. Role of the Dnmt3 family in de novo methylation of imprinted and repetitive
1288      sequences during male germ cell development in the mouse. *Hum Mol Genet* **16**: 2272–
1289      2280.

1290    Kerrebrock AW, Srivastava R, Gerbi SA. 1989. Isolation and characterization of ribosomal DNA
1291      variants from Sciara coprophila. *J Mol Biol* **210**: 1–13.

1292    Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. 2019. Graph-based genome alignment and
1293      genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol* **37**: 907–915.

1294    Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. 2017. Canu: Scalable and
1295      accurate long-read assembly via adaptive κ-mer weighting and repeat separation. *Genome*
1296      *Res* **27**: 722–736.

1297    Korf I. 2004. Gene finding in novel genomes. *BMC Bioinformatics* **5**: 59.

1298    Laetsch DR, Blaxter ML. 2017. BlobTools: Interrogation of genome assemblies. *F1000Research*
1299      **6**: 1287.

1300    Lam ET, Hastie A, Lin C, Ehrlich D, Das SK, Austin MD, Deshpande P, Cao H, Nagarajan N, Xiao M,
1301      et al. 2012. Genome mapping on nanochannel arrays for structural variation analysis and
1302      sequence assembly. *Nat Biotechnol* **30**: 771–776.

1303    Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**:
1304      357–9.

1305    Lawson ET, Mousseau TA, Klaper R, Hunter MD, Werren JH. 2001. Rickettsia associated with
1306      male-killing in a buprestid beetle. *Heredity (Edinb)* **86**: 497–505.

1307    Li B, Fillmore N, Bai Y, Collins M, Thomson JA, Stewart R, Dewey CN. 2014. Evaluation of de
1308      novo transcriptome assemblies from RNA-Seq data. *Genome Biol* **15**: 553.

1309    Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. 2015. MEGAHIT: an ultra-fast single-node solution
1310      for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*
1311      **31**: 1674–1676.

1312    Li E, Beard C, Jaenisch R. 1993. Role for DNA methylation in genomic imprinting. *Nature* **366**:
1313      362–5.

1314    Li H. 2016. Minimap and miniasm: fast mapping and de novo assembly for noisy long
1315      sequences. *Bioinformatics* **32**: 2103–10.

1316    Li H. 2018. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094–
1317      3100.

1318    Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
1319      *Bioinformatics* **25**: 1754–60.

1320    Lin Y, Yuan J, Kolmogorov M, Shen MW, Chaisson M, Pevzner PA. 2016. Assembly of long error-

1321  prone reads using de Bruijn graphs. *Proc Natl Acad Sci U S A* **113**: E8396–E8405.

1322  Liu J, Zhu Y, Luo G-Z, Wang X, Yue Y, Wang X, Zong X, Chen K, Yin H, Fu Y, et al. 2016. Abundant
1323  DNA 6mA methylation during early embryogenesis of zebrafish and pig. *Nat Commun* **7**:
1324  13052.

1325  Loman NJ, Quick J, Simpson JT. 2015. A complete bacterial genome assembled de novo using
1326  only nanopore sequencing data. *Nat Methods* **advance on**.

1327  Luo G-Z, Blanco MA, Greer EL, He C, Shi Y. 2015. DNA N(6)-methyladenine: a new epigenetic
1328  mark in eukaryotes? *Nat Rev Mol Cell Biol* **16**: 705–10.

1329  Luo G-Z, Hao Z, Luo L, Shen M, Sparvoli D, Zheng Y, Zhang Z, Weng X, Chen K, Cui Q, et al. 2018.
1330  N6-methyldeoxyadenosine directs nucleosome positioning in Tetrahymena DNA. *Genome*
1331  *Biol* **19**: 200.

1332  Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y, et al. 2012. SOAPdenovo2:
1333  an empirically improved memory-efficient short-read de novo assembler. *Gigascience* **1**:
1334  18.

1335  Matthews BJ, Dudchenko O, Kingan SB, Koren S, Antoshechkin I, Crawford JE, Glassford WJ,
1336  Herre M, Redmond SN, Rose NH, et al. 2018. Improved reference genome of Aedes aegypti
1337  informs arbovirus vector control. *Nature* **563**: 501–507.

1338  Mavrich TN, Jiang C, Ioshikhes IP, Li X, Venters BJ, Zanton SJ, Tomsho LP, Qi J, Glaser RL,
1339  Schuster SC, et al. 2008. Nucleosome organization in the Drosophila genome. *Nature* **453**:
1340  358–362.

1341  McAlpine JF, Wood DM. 1989. Manual of Nearctic Diptera. *Agric Canada Monogr* **3**.

1342  Mendelowitz LM, Schwartz DC, Pop M. 2015. Maligner: a fast ordered restriction map aligner.
1343  *Bioinformatics* **32**: 1016–1022.

1344  Mok EH, Smith HS, DiBartolomeis SM, Kerrebrock AW, Rothschild LJ, Lange TS, Gerbi SA. 2001.
1345  Maintenance of the DNA puff expanded state is independent of active replication and
1346  transcription. *Chromosoma* **110**: 186–96.

1347  Mondo SJ, Dannebaum RO, Kuo RC, Louie KB, Bewick AJ, LaButti K, Haridas S, Kuo A, Salamov A,
1348  Ahrendt SR, et al. 2017. Widespread adenine N6-methylation of active genes in fungi. *Nat*
1349  *Genet* **49**: 964–968.

1350  Nigro RG, Campos MCC, Perondini ALP. 2007. Temperature and the progeny sex-ratio in Sciara
1351  ocellaris (Diptera, Sciaridae). *Genet Mol Biol* **30**: 152–158.

1352  Nikolenko SI, Korobeynikov AI, Alekseyev MA. 2013. BayesHammer: Bayesian clustering for
1353  error correction in single-cell sequencing. *BMC Genomics* **14 Suppl 1**: S7.

1354  Panikar CS, Rajpathak SN, Abhyankar V, Deshmukh S, Deobagkar DD. 2015. Presence of DNA
1355  methyltransferase activity and CpC methylation in Drosophila melanogaster. *Mol Biol Rep*
1356  **42**: 1615–1621.

1357  Pardue M Lou, Gerbi SA, Eckhardt RA, Gall JG. 1970. Cytological localization of DNA
1358  complementary to ribosomal RNA in polytene chromosomes of Diptera. *Chromosoma* **29**:
1359  268–290.

1360  Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. 2017. Salmon provides accurate, fast, and
1361  bias-aware transcript expression. *Nat Methods* **14**: 417–419.

1362  Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. 2015. StringTie
1363  enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*
1364  **33**: 290–295.

Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, Lopez R. 2005. InterProScan: protein domains identifier. *Nucleic Acids Res* **33**: W116-20.

Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**: 841–2.

Raddatz G, Guzzardo PM, Olova N, Fantappié MR, Rampp M, Schaefer M, Reik W, Hannon GJ, Lyko F. 2013. Dnmt2-dependent methylomes lack defined DNA methylation patterns. *Proc Natl Acad Sci U S A* **110**: 8627–8631.

Rasch EM. 1970b. DNA cytophotometry of salivary gland nuclei and other tissue systems in dipteran larvae. In *In Introduction to Quantitative Cytochemistry* (eds. G.L. Wied and G.F. Bahr), pp. 357–397, Academic Press, New York.

Rasch EM. 2006. Genome size and determination of DNA content of the X chromosomes, autosomes, and germ line-limited chromosomes of Sciara coprophila. *J Morphol* **267**: 1316–25.

Rasch EM. 1970a. Two-wavelength cytophotometry of Sciara salivary gland chromosomes. In *Introduction to Quantitative Cytochemistry* (eds. G.L. Wied and G.F. Bahr), Vol. 2 of, pp. 335–355, Academic Press, New York.

Rausch C, Hastert FD, Cardoso MC. 2020. DNA Modification Readers and Writers and Their Interplay. *J Mol Biol*.

Roach MJ, Schmidt SA, Borneman AR. 2018. Purge Haplotigs: allelic contig reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics* **19**: 460.

Salzberg SL, Phillippy AM, Zimin A, Puiu D, Magoc T, Koren S, Treangen TJ, Schatz MC, Delcher AL, Roberts M, et al. 2012. GAGE: A critical evaluation of genome assemblies and assembly algorithms. *Genome Res* **22**: 557–67.

Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, von Haeseler A, Schatz MC. 2018. Accurate detection of complex structural variations using single-molecule sequencing. *Nat Methods* **15**: 461–468.

Serbus LR, Casper-Lindley C, Landmann F, Sullivan W. 2008. The Genetics and Cell Biology of *Wolbachia* -Host Interactions. *Annu Rev Genet* **42**: 683–707.

Sformo T, Kohl F, McIntyre J, Kerr P, Duman JG, Barnes BM. 2009. Simultaneous freeze tolerance and avoidance in individual fungus gnats, Exechia nugatoria. *J Comp Physiol B Biochem Syst Environ Physiol* **179**: 897–902.

Sharakhova M V., Hammond MP, Lobo NF, Krzywinski J, Unger MF, Hillenmeyer ME, Bruggner R V., Birney E, Collins FH. 2007. Update of the Anopheles gambiae pest genome assembly. *Genome Biol* **8**.

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

Simpson JT, Durbin R. 2010. Efficient construction of an assembly string graph using the FM-index. *Bioinformatics* **26**: i367-73.

Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJM, Birol I. 2009. ABySS: a parallel assembler for short read sequence data. *Genome Res* **19**: 1117–23.

Simpson JT, Workman RE, Zuzarte PC, David M, Dursi LJ, Timp W. 2017. Detecting DNA cytosine methylation using nanopore sequencing. *Nat Methods* **14**: 407–410.

Smit A, Hubley R. 2008. RepeatModeler Open-1.0. http://www.repeatmasker.org.

1409    Smit A, Hubley R, Green P. 2013. RepeatMasker Open-4.0. http://www.repeatmasker.org.

1410    Smith-Unna R, Boursnell C, Patro R, Hibberd JM, Kelly S. 2016. TransRate: reference-free quality
1411    assessment of de novo transcriptome assemblies. *Genome Res* **26**: 1134–44.

1412    Stuart JJ, Chen M-S, Shukle R, Harris MO. 2012. Gall Midges (Hessian Flies) as Plant Pathogens.
1413    *Annu Rev Phytopathol* **50**: 339–357.

1414    Suzuki Y, Korlach J, Turner SW, Tsukahara T, Taniguchi J, Qu W, Ichikawa K, Yoshimura J, Yurino
1415    H, Takahashi Y, et al. 2016. AgIn: measuring the landscape of CpG methylation of individual
1416    repetitive elements. *Bioinformatics* **32**: 2911–9.

1417    Takayama S, Dhahbi J, Roberts A, Mao G, Heo S-J, Pachter L, Martin DIK, Boffelli D. 2014.
1418    Genome methylation in D. melanogaster is found at specific short motifs and is
1419    independent of DNMT2 activity. *Genome Res* **24**: 821–830.

1420    The UniProt Consortium. 2019. UniProt: A worldwide hub of protein knowledge. *Nucleic Acids
1421    Res* **47**: D506–D515.

1422    Ter-Hovhannisyan V, Lomsadze A, Chernoff YO, Borodovsky M. 2008. Gene prediction in novel
1423    fungal genomes using an ab initio algorithm with unsupervised training. *Genome Res* **18**:
1424    1979–90.

1425    Urban JM, Bliss J, Lawrence CE, Gerbi SA. 2015. Sequencing ultra-long DNA molecules with the
1426    Oxford Nanopore MinION. *bioRxiv* doi: 10.1101/019281.

1427    Urnov FD, Liang C, Blitzblau HG, Smith HS, Gerbi SA. 2002. A DNase I hypersensitive site flanks
1428    an origin of DNA replication and amplification in Sciara. *Chromosoma* **111**: 291–303.

1429    Vaser R, Sović I, Nagarajan N, Šikić M. 2017. Fast and accurate de novo genome assembly from
1430    long uncorrected reads. *Genome Res* **27**: 737–746.

1431    Vezzi F, Narzisi G, Mishra B, Nagarajan N, Pop M, Vezzi F, Narzisi G, Mishra B, Lander E, Linton L,
1432    et al. 2012. Reevaluating Assembly Evaluations with Feature Response Curves: GAGE and
1433    Assemblathons ed. A. Rzhetsky. *PLoS One* **7**: e52210.

1434    Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q, Wortman
1435    J, Young SK, et al. 2014. Pilon: an integrated tool for comprehensive microbial variant
1436    detection and genome assembly improvement. *PLoS One* **9**: e112963.

1437    Wang X, Li Z, Zhang Q, Li B, Lu C, Li W, Cheng T, Xia Q, Zhao P. 2018. DNA methylation on N6-
1438    adenine in lepidopteran Bombyx mori. *Biochim Biophys Acta - Gene Regul Mech* **1861**:
1439    815–825.

1440    Wang Y, Chen X, Sheng Y, Liu Y, Gao S. 2017. N6-adenine DNA methylation is associated with
1441    the linker DNA of H2A.Z-containing well-positioned nucleosomes in Pol II-transcribed
1442    genes in Tetrahymena. *Nucleic Acids Res* **45**: 11594–11606.

1443    Werren JH, Windsor DM. 2000. Wolbachia infection frequencies in insects: Evidence of a global
1444    equilibrium? *Proc R Soc B Biol Sci* **267**: 1277–1285.

1445    White MJ. 1949. Cytological evidence on the phylogeny and classification of the Diptera.
1446    *Evolution* **3**: 252–261.

1447    Wiegmann BM, Richards S. 2018. Genomes of Diptera. *Curr Opin Insect Sci* **25**: 116–124.

1448    Wiegmann BM, Trautwein MD, Winkler IS, Barr NB, Kim JW, Lambkin C, Bertone MA, Cassel BK,
1449    Bayless KM, Heimberg AM, et al. 2011. Episodic radiations in the fly tree of life. *Proc Natl
1450    Acad Sci U S A* **108**: 5690–5695.

1451    Wu N, Liang C, DiBartolomeis SM, Smith HS, Gerbi SA. 1993. Developmental progression of DNA
1452    puffs in Sciara coprophila: amplification and transcription. *Dev Biol* **160**: 73–84.

64

1453 Xiao C-L, Zhu S, He M, Chen D, Zhang Q, Chen Y, Yu G, Liu J, Xie S-Q, Luo F, et al. 2018. N6-
1454      Methyladenine DNA Modification in the Human Genome. *Mol Cell* **71**: 306-318.e7.
1455 Ye C, Hill CM, Wu S, Ruan J, Ma Z (Sam). 2016. DBG2OLC: Efficient Assembly of Large Genomes
1456      Using Long Erroneous Reads of the Third Generation Sequencing Technologies. *Sci Rep* **6**:
1457      31900.
1458 Zemach A, McDaniel IE, Silva P, Zilberman D. 2010. Genome-wide evolutionary analysis of
1459      eukaryotic DNA methylation. *Science (80- )* **328**: 916–919.
1460 Zerbino DR, Birney E. 2008. Velvet: algorithms for de novo short read assembly using de Bruijn
1461      graphs. *Genome Res* **18**: 821–9.
1462 Zhang G, Huang H, Liu D, Cheng Y, Liu X, Zhang W, Yin R, Zhang D, Zhang P, Liu J, et al. 2015. N6-
1463      Methyladenine DNA Modification in Drosophila. *Cell* **161**: 893–906.
1464 Zhou C, Wang C, Liu H, Zhou Q, Liu Q, Guo Y, Peng T, Song J, Zhang J, Chen L, et al. 2018.
1465      Identification and analysis of adenine N6-methylation sites in the rice genome. *Nat Plants*
1466      **4**: 554–563.
1467