

1 **TITLE:** Chromosome-scale scaffolding of the fungus gnat genome (Diptera: *Bradysia*
2 *coprophila*).

3
4 **AUTHORS:**

5 John M. Urban¹, Susan A. Gerbi^{2#}, Allan C. Spradling^{1#}
6
7

8 ¹ Carnegie Institution for Science, Department of Embryology, Howard Hughes Medical
9 Institute Research Laboratories, 3520 San Martin Drive, Baltimore, MD, 21218, USA
10

11 ² Brown University Division of Biology and Medicine, Department of Molecular Biology,
12 Cell Biology and Biochemistry, Providence, RI 02912 USA
13

14 # Co-senior authors contributed equally to this work.
15

16 **EMAILS:**

17 John M. Urban : jurban@carnegiescience.edu

18 Susan A. Gerbi : susan_gerbi@brown.edu

19 Allan C. Spradling : spradling@carnegiescience.edu
20
21

22 **CORRESPONDING AUTHOR:**

23 John M. Urban
24

25 **ABSTRACT:**

26 Background: The lower dipteran fungus gnat, *Bradysia (Sciara) coprophila*, has unusual
27 chromosome biology. Chromosome imprinting was first discovered in this system. All
28 paternal chromosomes are eliminated during spermatogenesis whereas both maternal X
29 sister chromatids are retained. Embryos start out with three copies of the X chromosome,
30 but 1-2 copies are eliminated from all somatic cells as a part of sex determination, and
31 one is eliminated in the germline to restore diploidy. These developmentally normal
32 features present opportunities to study unusual chromosome movements that may occur
33 as rare/abnormal events in other systems. To help with such studies, we previously
34 generated a highly contiguous optical-map-scaffolded long-read assembly (Bcop_v1) of
35 the male somatic genome that contains four chromosomes. However, the scaffolds were
36 not chromosome-scale, the majority of the assembly lacked chromosome assignments,
37 and the order and orientation of the contigs along chromosomes remained unknown.
38

39 Findings: Male pupae chromatin conformation capture (Hi-C) data was collected and used
40 to first produce a corrected Bcop_v1 assembly by breaking contigs at mis-joined regions,
41 and then used to order and orient the corrected contigs into chromosome-scale scaffolds.
42 The resulting assembly (Bcop_v2) had four chromosome-scale scaffolds as expected.
43 Previously known chromosomal locations of locus-specific sequences were used to (i)
44 identify the corresponding chromosome for each scaffold, and (ii) orient the chromosome
45 scaffolds in the same order as published polytene chromosome maps. Finally, the gene

46 annotations produced for Bcop_v1 were lifted over to Bcop_v2 to allow a seamless
47 transition in adopting the updated chromosome-scale assembly.

48

49

50 Conclusions: Studies of the unusual chromosome movements in *Bradysia coprophila* will
51 benefit from the updated assembly (Bcop_v2) where each somatic chromosome is
52 represented by a single scaffold.

53

54

55 **KEYWORDS:** *Bradysia coprophila*, *Sciara coprophila*, fungus gnats, Diptera, insects,
56 genome assembly, chromosomes, scaffolding, Hi-C

57

58

59

60 **Introduction:**

61

62 The lower Dipteran dark-winged fungus gnat, *Bradysia (Sciara) coprophila*, is an
63 important model system for studying chromosome biology due to its dynamic genome.
64 Most notable is its chromosome cycle. Rather than having either a diploid or haploid set
65 of the same chromosomes (X, II, III, IV, and L) in every cell, the chromosome constitution
66 of a cell varies based on whether it is somatic or germline, male or female, and early
67 embryo or late embryo [1]. Regarding somatic or germline, there are chromosomes
68 referred to as the L chromosomes that are only found in germ cells [1]. Regarding male
69 or female, half of all females, but no males, contain a variant of the X chromosome called
70 the X' (X prime) [1]. Regarding embryo age, embryos start out with three copies of the X
71 chromosome, two of which are paternally derived, and with a variable number of the
72 germline-limited L chromosomes, but later embryos have only 1-2 X chromosomes and
73 no L chromosomes in somatic cells, for example. The L chromosomes are eliminated
74 from all somatic nuclei not set aside for the germline in the 5th-6th nuclear division [1]. The
75 difference in X chromosome copy number arises because embryos from XX mothers are
76 destined to be male and eliminate two paternal X chromosomes from somatic nuclei in
77 the 7th-9th nuclear division whereas embryos produced by X'X mothers are fated to be
78 female and eliminate only one paternal X [1]. There is a later point in development in
79 which both the male and female germline are diploid for the X after eliminating one
80 paternal copy [1]. In addition, one or more L chromosomes are eliminated from the
81 germline around the same time, which prevents their accumulation [1]. Thereafter, the
82 meiotic events in oogenesis appear to proceed normally whereas those in
83 spermatogenesis do not. In meiosis I of spermatogenesis, all paternal somatic
84 chromosomes are eliminated in a bud of cytoplasm while the L chromosomes appear to
85 escape this imprinting effect and are retained with the maternal set of somatic
86 chromosomes at the single pole of a naturally occurring monopolar spindle [1]. In meiosis
87 II, both maternally-derived sister chromatids of the X chromosome undergo
88 developmentally programmed nondisjunction and are retained at one pole of a bipolar
89 spindle [1]. Rather than four products, these unusual events lead to one product of male
90 meiosis: a sperm that contains only maternally derived somatic chromosomes (X, II, III,
91 and IV), that is diploid for X, haploid for the autosomes (II, III, and IV), and variable for L

92 chromosomes. In addition to the unusual chromosome cycle, *Bradysia coprophila* larval
93 salivary glands have highly polytene chromosomes with over 8000 copies in each nucleus
94 [2], and developmentally-programmed gene amplification [3,4]. Overall, fungus gnats
95 present many opportunities to study chromosome biology.

96
97 Previously, we published a highly contiguous assembly of the somatic
98 chromosomes (X, II, III, and IV) [5]. We targeted the somatic genome in males as it
99 presented the lowest complexity version of the genome, lacking both the L and X'
100 chromosomes. Specifically, we sequenced genomic DNA from male embryos using long-
101 read (PacBio RS II SMRT and Oxford Nanopore Technologies MinION) and short-read
102 (Illumina) technologies [5]. As part of that process, we assembled the sequencing data
103 ~100 different ways, and used many evaluations to determine the best assembly and
104 polishing pipelines for our datasets [5]. The assembly with the best evaluation scores was
105 scaffolded with optical maps from BioNano Genomics [5]. The resulting assembly,
106 Bcop_v1, had a contig NG50 of ~2.4 Mb [5] and scaffold NG50 of ~8.2 Mb (Table 1) [5].
107 Since, in *Bradysia coprophila*, late male embryos have a single copy of the X, but are
108 diploid for the autosomes, it was possible to classify each contig as either X-linked or
109 autosomal [5]. Moreover, using the known chromosomal locations of a handful of
110 sequences, we were able to anchor ~28-33% of the autosomal sequences into
111 chromosomes II, III, and IV [5]. Overall, nearly 50% of the expected genome size was
112 anchored into chromosomes based on coverage and known sequence locations [5].
113 Nevertheless, chromosome assignments are not available for the majority of that
114 assembly (Bcop_v1). Moreover, Bcop_v1 lacks any information regarding the ordering
115 and orientation of the contigs along chromosome sequences. To better facilitate studies
116 into gene amplification and the unusual chromosome cycle in *Bradysia coprophila*, our
117 goal was to develop better models of the chromosome sequences. We approached this
118 by using chromosome conformation capture with deep sequencing (Hi-C) to scaffold the
119 chromosomes [6], as has been successfully done by many others [7-9]. Specifically, Hi-
120 C is so powerful at putting chromosome sequences together from a collection of sub-
121 chromosomal sequences because of two major attributes of in vivo chromosome
122 interactions captured by Hi-C data: (i) because intra-chromosomal interactions are much
123 more frequent than inter-chromosomal interactions, it is possible to use interaction
124 frequencies to cluster contigs into 'chromosome groups', and (ii) because interaction
125 frequencies decay with distance, it is possible to use them to order and orient the contigs
126 in a chromosome group in the order they most likely appear along a chromosome.

127
128 Overall, we present here the first chromosome-scale genome assembly for
129 *Bradysia coprophila*, specifically reporting the first full-length scaffolds for the somatic
130 chromosomes (X, II, III, IV). In addition, we lifted over the annotations from the former
131 genome assembly, allowing studies in progress to seamlessly transition to the updated
132 genome assembly. We expect that the chromosome-scale genome and gene annotations
133 will be immediately useful to the growing research community interested in the unique
134 biology of *Bradysia coprophila* as well as the broader research community interested in
135 Dipteran evolution and comparative genomics.

136
137

138 **Methods and Results:**

139

140 **Male Pupae Collection**

141 The dark-winged fungus gnat, *Bradysia coprophila* [10], was previously referred to
142 as *Sciara coprophila* (Lintner, 1895) in chromosomal and molecular biology research
143 papers since the 1920s [1,11], and is also known by other names such as *Sciara tillicola*
144 (Loew, 1850) and *Sciara amoena* (Winnertz, 1867). In this study, the fungus gnats were
145 from the HoLo2 line maintained in the International *Sciara* Stock Center at Brown
146 University (<https://sites.brown.edu/sciara/>). *B. coprophila* females are monogenic,
147 meaning they have either only male or only female offspring, which is determined by
148 whether they harbor a variant of the X chromosome (X') or not. Specifically, X'X females
149 are female producers and XX females are male producers. The Holo2 line has phenotypic
150 marker gene on the X' chromosome called *Wavy*. Female producers (X'X) have the *Wavy*
151 wing phenotype whereas male producers (XX) have wild-type straight wings. To collect
152 only male pupae, crosses between straight-winged females (XX) and males (XO) were
153 used to obtain strictly male progeny from which pupae were collected at the appropriate
154 time approximately 4 weeks later. The collected male pupae were snap frozen in liquid
155 nitrogen and stored at -80°C until needed.

156

157 **Hi-C Data Collection**

158 The frozen male pupae were ground into a fine powder. Chromatin conformation
159 capture data was then generated using a Phase Genomics (Seattle, WA) Proximo Hi-C
160 2.0 Kit, which is a commercially available version of the Hi-C protocol [6]. Following the
161 manufacturer's instructions, intact cells from two samples of finely ground frozen male
162 pupae were crosslinked using a formaldehyde solution, digested using the Sau3A1
163 restriction enzyme, and proximity-ligated with biotinylated nucleotides. This creates
164 chimeric molecules made from DNA fragments that came close together *in vivo*, but that
165 may not be close together along the genome sequence. As instructed by the
166 manufacturer's protocol, the chimeric molecules were pulled down with streptavidin beads
167 and processed into an Illumina-compatible sequencing library. Sequencing was
168 performed on an Illumina NovaSeq, generating a total of 115.2 million 2x101 bp paired-
169 end reads.

170

171 **Hi-C assembly correction and scaffolding**

172 The Hi-C data was processed, as described below, according to Phase Genomics
173 recommendations found here: [https://phasegenomics.github.io/2019/09/19/hic-
174 alignment-and-qc.html](https://phasegenomics.github.io/2019/09/19/hic-alignment-and-qc.html). The paired-end Hi-C reads were first used to produce a corrected
175 assembly. Briefly, the reads were aligned to the primary contigs of Bcop_v1 [5] using
176 BWA-MEM [12] with the -5SP and -t 8 options specified, and all other options default. The
177 output was piped into SAMBLASTER [13] to flag PCR duplicates that were later excluded
178 from analyses. That output was subsequently piped into SAMtools [14] using "-F 2304" to
179 remove non-primary and secondary alignments. Juicebox was used to identify and break
180 contigs at 46 putative mis-joined regions based on disruptions in the expected pattern
181 from Hi-C alignments along the contigs [7,15]. This transformed the input assembly
182 (Bcop_v1) into the corrected assembly (Bcop_v1_corrected). The paired-end Hi-C reads

183 were then aligned to Bcop_v1_corrected using the same alignment procedure as above,
184 and these alignments were the input to the Hi-C-based scaffolding described next.

185
186 The Phase Genomics' Proximo Hi-C genome scaffolding platform was used to
187 create chromosome-scale scaffolds from the Bcop_v1_corrected as described in Bickhart
188 et al. [8]. As in the LACHESIS method [9], the Proximo scaffolding process computes a
189 contact frequency matrix from the aligned Hi-C read pairs that is normalized by the
190 number of DPNII restriction sites (GATC) on each contig. Proximo then constructs
191 scaffolds by optimizing the expected contact frequencies and other statistical patterns in
192 Hi-C data. Using a brute-force approach, approximately 20,000 separate Proximo runs
193 were performed to optimize the number of scaffolds and the concordance with the
194 observed Hi-C data. Finally, Juicebox was again used to correct scaffolding errors,
195 resulting in a final set of four primary chromosome-scale scaffolds each spanning 58-71
196 Mb (Table 1, Figure 1), a ~363 kb short scaffold, and 57 unscaffolded sequences with 46
197 marked as “debris”. The “debris” was largely comprised of optical-map gap sequences
198 (NNNNNN) from the input assembly. Specifically, six unscaffolded debris sequences
199 ranging from ~25-50 kb had 100% N content. These sequences were removed from the
200 assembly. Furthermore, 24 sequences had either leading or trailing NNNNNs ranging
201 from 1,489-36,461 bp, making up ~6-96% of the unscaffolded sequences they resided
202 on. The leading and trailing N content was trimmed off these sequences. The remaining
203 52 cleaned up unscaffolded sequences summed to ~1.8 Mb, with just ~19.5 kb N content.
204 Thus, the vast majority of the ~299 Mb of primary sequence from Bcop_v1 (Table 1) input
205 into the Hi-C scaffolding process was integrated into the four primary chromosome
206 scaffolds that summed to ~297 Mb. The 539 associated contigs from Bcop_v1 that were
207 not included during the Hi-C scaffolding process were also added back to the assembly.
208 These sequences summed to 10.8 Mb, and were ~20 kb on average. Note that we also
209 tried including the Bcop_v1 associated sequences during the Hi-C scaffolding process,
210 but preferred the results when using only the primary assembly. We now refer to the four
211 chromosome-scale scaffolds as the primary assembly, and the other 591 sequences as
212 the associated contigs. Most of the associated contigs are haplotigs, and researchers
213 interested mainly in the chromosome-scale scaffolds can opt to use or ignore them.

214
215 Overall, Phase Genomics Proximo Hi-C genome scaffolding yielded the updated
216 genome assembly, Bcop_v2, that has four primary chromosome scaffolds of 58.3, 70.5,
217 71.0, and 97 Mb (Table 2, Figure 1), which make up ~96% of the entire assembly. The
218 four primary scaffolds together with the 591 associated sequences have an NG50 of ~71
219 Mb pertaining to only 2 scaffolds, and an expected size of ~81 Mb (Table 1). The
220 assembly retained 209 gaps from the original BioNano scaffolding process, that ranged
221 from 25 bp to ~662 kb with a gap N50 of ~101.5 kb. The optical map gaps are represented
222 as capital N's in the assembly and are of variable size. Hi-C scaffolding introduced 241
223 gaps, 5 of which are on the short ~363 kb scaffold, the remainder being in the
224 chromosome scaffolds. The Hi-C gaps are represented as lower-case n's and are always
225 100 bp in length.

226
227
228

229 **Anchoring and orienting the scaffolds**

230 Sequences with known chromosomal locations based on previous *in situ*
231 hybridization work were used with BLAST [16] to identify the corresponding chromosome
232 for each scaffold (Table 3). The chromosome assignments were in general agreement
233 with the expected sizes of each chromosome (Table 2) with chromosomes II and X being
234 the smallest, followed by chromosomes III and IV. The total length of autosomes (226
235 Mb) was also concordant with expectations (218-231 Mb) (Table 2).
236

237 Since the locations of multiple sequences for each chromosome are known, it was
238 also possible to orient the chromosome scaffold sequences in the same direction as locus
239 numbers have been assigned on polytene maps [17]. There were no conflicts among the
240 multiple anchors. In other words, each scaffold corresponded to one and only one
241 chromosome, the 4 chromosome-scale scaffolds corresponded to different
242 chromosomes, and when multiple anchors were present for a given chromosome, they
243 all mapped in the expected order along the same scaffold. In addition to unique locus-
244 specific landmarks, the centromere-associated sequence, Sccr, (*Sciara coprophila*
245 centromeric repeat [18]) were in the expected positions. Specifically, chromosomes X, II,
246 and III are acrocentric with their centromeres very close to the beginning of the
247 chromosomes (as oriented). We found Sccr mapped within the first few megabases of
248 the corresponding scaffolds. In contrast, chromosome IV is metacentric with a centrally-
249 located centromere. We found Sccr mapped to more medial positions on the scaffold
250 corresponding to IV. Overall, the agreement with expected sizes, the lack of anchor
251 conflicts, the concordance with the sequence of loci on polytene maps, and the relative
252 positioning of centromere-associated sequence suggests that the Hi-C scaffolding
253 process accurately clustered contigs into chromosome groups and put them in the correct
254 order.
255
256

257 **Liftover of gene annotations to the chromosome-scale assembly (Bcop_v2).**

258 There are presently two high quality gene annotation sets produced using the
259 previous genome assembly (Bcop_v1) that was used as the input for Hi-C scaffolding.
260 One was constructed by us [5] using Maker2 [19], which can be found at USDA Ag Data
261 Commons [20] and the USDA i5k Workspace [21]. The other was made by the NCBI
262 Eukaryotic Genome Annotation Pipeline for RefSeq, is called NCBI *Bradysia coprophila*
263 Annotation Release 100, and can be found at NCBI [22-24]. While the FASTA files
264 containing the transcript and protein sequences are in no need of updating, the GFF files
265 that detail the coordinates of the gene models, including their exons, introns, CDS, and
266 UTRs, on Bcop_v1 are not useful for the updated assembly, Bcop_v2. Since Bcop_v2
267 was produced from the Bcop_v1 sequences, a simple liftover procedure was adequate
268 for creating new GFF files detailing the coordinates of the genes on Bcop_v2. Specifically,
269 we used the program designed for this task, LiftOff [25] (version v1.6.3) with Minimap2
270 [26] (version 2.24-r1122). Along with the Bcop_v1 and Bcop_v2 assembly FASTA files
271 and the Bcop_v1 GFF file, the following parameters were given to LiftOff: -p 16 -u
272 unmapped.txt -polish -flank 0.25 -a 0.9 -s 0.9.
273

274 Using the Maker2 annotation as an example of results, of the 23,117 genes and
275 28,870 transcripts, only 2 genes with 4 transcripts were not mapped to Bcop_v2. Of the
276 28,866 transcript models successfully lifted over to Bcop_v2, 28865 (all but one) had
277 perfect full-length alignments to the original transcript sequences, although 53 (<0.2%)
278 were a bit shorter than the original. We can visually confirm in a genome browser that the
279 gene models at the same loci look the same between assemblies and that the lifted over
280 gene models are highly concordant with RNA-seq datasets from various life stages
281 mapped to Bcop_v2 (data not shown) [5].

282
283

284 **Conclusions:**

285 Our goal was to upgrade Bcop_v1 [5] to chromosome-scale scaffolds. Bcop_v1 is
286 a high quality, high contiguity assembly of the somatic genome (chromosomes X, II, III,
287 IV) produced from PacBio RS II and Oxford Nanopore MinION long read data from male
288 embryos, and subsequently scaffolded with optical maps of male pupae DNA molecules
289 from BioNano Genomics. Bcop_v1 was produced through an extensive evaluation
290 process comparing nearly 100 assemblies. Moreover, as *Bradysia coprophila* males have
291 only one copy of the X chromosome, but two copies of autosomes, all contigs in Bcop_v1
292 were classified as X-linked or autosomal. Finally, two gene annotation sets were
293 constructed for Bcop_v1. We produced one using Maker2 [5,19-21] and NCBI produced
294 one for RefSeq [5,22-24]. As studies are ongoing with these gene sets, we also sought
295 to lift over these gene annotation sets directly onto the chromosome-scale scaffolds so
296 researchers can seamlessly transition current studies to the new assembly.

297
298 The updated assembly of the *Bradysia coprophila* somatic genome (Bcop_v2)
299 represents each chromosome (X, II, III, IV) as its own scaffold for the first time, and is
300 oriented in the same direction as polyene maps produced in the 1960s [17] to be
301 consistent with historical research. In the future, the polytene zones in those chromosome
302 maps can roughly be mapped onto their corresponding Bcop_v2 chromosome-scaffold
303 sequences using a combination of landmark sequences identified by FISH, genomic
304 profiling techniques, and assumptions about the sizes of the polytene zones. Indeed, it
305 may be possible to re-purpose the Hi-C data reported here to identify bands and
306 interbands, which have been associated with topologically associated domains (TADs) in
307 *Drosophila* polytene chromosomes using Hi-C data [27-28], across the Bcop_v2
308 chromosome sequences. Regardless, the male (XO) pupae Hi-C paired-end reads will
309 continue to be useful for correcting mis-assemblies and scaffolding future *Bradysia*
310 *coprophila* genome assemblies produced with newer long-read datasets (e.g. PacBio HiFi
311 or ultra-long Q20 nanopore reads).

312

313 We have also provided the Bcop_v1 gene annotations directly lifted over to
314 Bcop_v2 to enable researchers to seamlessly shift ongoing studies from Bcop_v1 to
315 Bcop_v2. The chromosome-level nature of Bcop_v2 will allow researchers to design
316 experiments and analyze data in a chromosome-specific way. The location and
317 sequences of the telomeres and centromeres are now more obvious and open to future
318 investigation. The X chromosome scaffold will prove useful in identifying the Controlling
319 Element, an rDNA-associated sequence embedded in the short arm of the X upstream of

320 the centromere. Moreover, the X chromosome scaffold will benefit studies of the X
321 chromosome variant, X-prime (X'), which differs from the X by a large paracentric
322 inversion [1]. There are also interesting biological signals in the Hi-C maps that can be
323 explored further. The notable example here are the three dots on the X chromosome Hi-
324 C map (see the third chromosome cluster in Figure 1) that appear in the expected pattern
325 of the three foldback repeats identified by Crouse [29]. Finally, the *B. coprophila*
326 chromosome scaffolds will allow proper studies of synteny and chromosome evolution
327 across closely and distantly related flies, and may prove useful for researchers working
328 with closely related species (such as *B. odoraphaga*, *B. ocellaris*, *B. impatiens*) when they
329 are working with fragmented genome assemblies by letting them cluster the contigs into
330 likely chromosome groups, and even ordering and orienting them into pseudo-
331 chromosome scaffolds.

332

333 **Data Availability**

334 This Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank
335 under the accession VSDI00000000. The version described in this paper, referred to
336 throughout as Bcop_v2, is version VSDI02000000. This project is associated with
337 BioProject PRJNA291918 and BioSample SAMN12533751. All data, including Hi-C
338 sequencing data and lifted over gene models, will be available upon publication of this
339 manuscript.

340

341 **Abbreviations**

342 Bcop: *Bradysia coprophila*; bp: base pairs; BUSCO: Benchmarking Universal Single-
343 Copy Orthologs; Gb: gigabase pairs; Hi-C = High dimensional Chromosome conformation
344 capture; i5k: initiative to sequence 5000 insect (and other arthropod) genomes
345 (<http://i5k.github.io/about>) ; kb: kilobase pairs; Mb: megabase pairs; NCBI: National
346 Center for Biotechnology Information; PacBio: Pacific Biosciences; PCR: polymerase
347 chain reaction; Sccr: *Sciara coprophila* centromeric repeat ; SMRT: single-molecule real-
348 time. ONT: Oxford Nanopore Technologies.

349

350 **Competing Interests**

351 Authors declare no competing interests.

352

353 **Funding**

354 Financial support has come from National Institutes of Health (NIH) NIH/GM121455 to
355 SAG. ACS is an Investigator of the Howard Hughes Medical Institute (HHMI) in the HHMI
356 unit at the Carnegie Institute for Science (CIS). JMU is an HHMI/CIS Research Associate
357 under ACS.

358

359 **Authors' Contributions**

360 JMU, SAG, and ACS conceived the project. SAG collected male pupae. Hi-C library
361 preparation, sequencing, and scaffolding was performed as a service provided by Phase
362 Genomics. JMU performed the original Bcop_v1 assembly, and corresponded with Phase
363 Genomics during the optimization process for Bcop_v2 scaffolding. JMU performed all
364 subsequent analyses and operations, including chromosome identification, scaffold
365 orientation, assembly evaluations, further Hi-C visualizations, annotation liftover, repeat

366 annotation, and synteny analyses. All authors participated intellectually in the
367 development and execution of this project. All authors wrote, read, and edited the
368 manuscript.

369

370 **Acknowledgements**

371 We thank Jacob Bliss who helped with *Bradysia coprophila* maintenance during data
372 collection. We thank the correspondents we interacted with at Phase Genomics, including
373 Hayley Mangelson, Kayla Young, Shawn Sullivan, Brian Fan, Andrew Wiser, and Ivan
374 Liachko. We would like to acknowledge Rob Baird and Laura Ross for being early
375 adopters of Bcop_v2. We thank Jennifer Urban for proof-reading the manuscript.

376

377 **References**
378

- 379 1. Gerbi SA. Unusual chromosome movements in Sciarid flies. In: Hennig W, editor.
380 Results and problems in cell differentiation, vol. 13 Germ Line - Soma
381 Differentiation. Berlin, Heidelberg: Springer-Verlag; 1986. p. 71–104.
382 2. Rasch EM. DNA cytophotometry of salivary gland nuclei and other tissue systems
383 in dipteran larvae. In: Wied BGF, editor. Introduction to quantitative cytochemistry,
384 vol. 2. New York: Academic Press; 1970b. p. 357–97.
385 3. Gerbi SA, Strezoska Z, Waggener JM. Initiation of DNA replication in multicellular
386 eukaryotes. J Struct Biol. 2002;140(1-3):17–30. [https://doi.org/10.1016/S1047-](https://doi.org/10.1016/S1047-8477(02)00538-5)
387 [8477\(02\)00538-5](https://doi.org/10.1016/S1047-8477(02)00538-5).
388 4. Rasch EM. Two-wavelength cytophotometry of *Sciara* salivary gland
389 chromosomes. In: Wied BGF, editor. Introduction to quantitative cytochemistry,
390 vol. 2. New York: Academic Press; 1970a. p. 335–55.
391 5. Urban JM, et al., High contiguity de novo genome assembly and DNA modification
392 analyses for the fungus fly, *Sciara coprophila*, using single-molecule sequencing.
393 BMC Genomics 22, 1-13 (2021).
394 6. Lieberman-Aiden, E., Van Berkum, N.L., Williams L., Imakaev M., Ragoczy T.,
395 Telling A., Amit, I., Lajoie, B.R., Sabo, P.J., Dorschner, M.O., Sandstrom, R.,
396 Bernstein, B., Bender, M.A., Groudine, M., Gnirke, A., Stamatoyannopoulos, J.,
397 Mirny, L.A., Lander, E.S., Dekker, J. (2009). Comprehensive mapping of long-range
398 interactions reveals folding principles of the human genome. Science, 326, 289-
399 293.
400 7. Suhas S.P. Rao*, Miriam H. Huntley*, Neva C. Durand, Elena K. Stamenova, Ivan
401 D. Bochkov, James T. Robinson, Adrian L. Sanborn, Ido Machol, Arina D. Omer,
402 Eric S. Lander, Erez Lieberman Aiden. “A 3D Map of the Human Genome at
403 Kilobase Resolution Reveals Principles of Chromatin Looping.” Cell 159, 2014.
404 8. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, Lee J, Lam ET,
405 Liachko I, Sullivan ST, Burton JN, Huson HJ, Nystrom JC, Kelley CM, Hutchison
406 JL, Zhou Y, Sun J, Crisà A, Ponce de León FA, Schwartz JC, Hammond JA,
407 Waldbieser GC, Schroeder SG, Liu GE, Dunham MJ, Shendure J, Sonstegard TS,
408 Phillippy AM, Van Tassell CP, Smith TP. Single-molecule sequencing and
409 chromatin conformation capture enable de novo reference assembly of the
410 domestic goat genome. Nat Genet. 49(4):643-650 (2017).
411 9. Burton, JN; Adey, A; Patwardhan, RP; Qiu, R; Kitzman, JO; Shendure, J.
412 Chromosome-scale scaffolding of de novo genome assemblies based on
413 chromatin interactions. Nat. Biotech. 31, 1119 (2013).
414 10. Steffan WA. A generic revision of the family Sciaridae (Diptera) of America north
415 of 892 Mexico. University of California publications in Entomology; 1966.
416 11. Metz CW. Chromosomes and sex in *Sciara*. Science. 1925;61(1573):212-214.
417 12. Li H. and Durbin R. (2010) Fast and accurate long-read alignment with Burrows-
418 Wheeler transform. Bioinformatics, 26, 589-595. [PMID: 20080505]
419 13. Gregory G. Faust, Ira M. Hall; SAMBLASTER: fast duplicate marking and structural
420 variant read extraction, Bioinformatics, Volume 30, Issue 17, 1 September 2014,
421 Pages 2503– 2505, <https://doi.org/10.1093/bioinformatics/btu314>

- 422 14. Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., Marth G.,
423 Abecasis G., Durbin R. and 1000 Genome Project Data Processing Subgroup
424 (2009) The Sequence alignment/map (SAM) format and SAMtools. *Bioinformatics*,
425 25, 2078-9. [PMID: 19505943]
- 426 15. Neva C. Durand*, James T. Robinson*, Muhammad S. Shamim, Ido Machol, Jill
427 P. Mesirov, Eric S. Lander, and Erez Lieberman Aiden. “Juicebox provides a
428 visualization system for Hi-C contact maps with unlimited zoom.” *Cell Systems*,
429 July 2016.
- 430 16. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search
431 tool. *J Mol Biol.* 1990;215(3):403–10. [https://doi.org/10.1016/S0022-](https://doi.org/10.1016/S0022-2836(05)80360-2)
432 [2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2).
- 433 17. Gabrusewycz-Garica N. Cytological and autoradiographic studies in *Sciara*
434 *coprophila* salivary gland chromosomes. *Chromosoma.* 1964;15(3):312–44.
435 <https://doi.org/10.1007/BF00321517>.
- 436 18. Escribá MC, Greciano PG, Méndez-Lago M, De Pablos B, Trifonov VA, Ferguson-
437 Smith MA, et al. Molecular and cytological characterization of repetitive DNA
438 sequences from the centromeric heterochromatin of *Sciara coprophila*.
439 *Chromosoma.* 2011;120(4):387–97. <https://doi.org/10.1007/s00412-011-0320-2>.
- 440 19. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database
441 management tool for second-generation genome projects. *BMC Bioinformatics.*
442 2011;12(1):491. <https://doi.org/10.1186/1471-2105-12-491>.
443 <https://doi.org/10.15482/USDA.ADC/1522618>
- 444 21. <https://i5k.nal.usda.gov/content/data-downloads>
- 445 22. https://ftp.ncbi.nlm.nih.gov/genomes/all/annotation_releases/38358/100/
- 446 23. https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Bradysia_coprophila/100/
- 447 24. [https://ftp.ncbi.nlm.nih.gov/genomes/all/annotation_releases/38358/100/GCF_01](https://ftp.ncbi.nlm.nih.gov/genomes/all/annotation_releases/38358/100/GCF_014529535.1_BU_Bcop_v1/)
448 [4529535.1 BU Bcop v1/](https://ftp.ncbi.nlm.nih.gov/genomes/all/annotation_releases/38358/100/GCF_014529535.1_BU_Bcop_v1/)
- 449 25. Shumate, Alaina, and Steven L. Salzberg. 2020. “Liftoff: Accurate Mapping of
450 Gene Annotations.” *Bioinformatics*, December.
451 <https://doi.org/10.1093/bioinformatics/btaa1016>.
- 452 26. Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences.
453 *Bioinformatics*, 34:3094-3100. doi:10.1093/bioinformatics/bty191
- 454 27. Eagen KP, Hartl TA, and Kornberg RD. Stable chromosome condensation
455 revealed by chromosome conformation capture. *Cell.* 2015 November 5; 163(4):
456 934–946. doi:10.1016/j.cell.2015.10.026.
- 457 28. Spradling AC. Symposium - Germ Cells, Imprinting, Gene Dosage, and
458 Regulation: Polytene Chromosome Structure and Somatic Genome
459 Instability. *Cold Spring Harb Symp Quant Biol* 2017 82: 293-304; Published in
460 Advance November 22, 2017, doi:10.1101/sqb.2017.82.033670
- 461 29. Crouse HV. X heterochromatin subdivision and cytogenetic analysis in *Sciara*
462 *coprophila* (Diptera, Sciaridae): Centromere localization. *Chromosoma.*
463 1977;63(1):39-55.
464
465
466
467

468 **Table 1 – Assembly Contiguity Statistics comparing the input (Bcop_v1) into Hi-C**
 469 **scaffolding to the resulting chromosome-scale assembly (Bcop_v2).**

	Bcop_v1 (Urban et al 2021)	Bcop_v1 Primary Only (Urban et al 2021)	Bcop_v2 (This paper)	Bcop_v2 Primary Only (This paper)
Number of sequences	744	205	595	4
Total length	309,775,056	298,965,442	309,636,011	296,980,291
Max sequence length	23,039,227	23,039,227	97,081,274	97,081,274
Min sequence length	671	4,264	671	58,343,183
Mean sequence length	416,364	1,458,368	520,397	74,245,073
Median sequence length	21,428	94,665	17,801	70,777,917
Sequence length N50	6,790,317	6,790,317	71,047,972	71,047,972
Sequence length L50	14	14	2	2
Expected sequence length	7,982,431	8,269,986	73,791,194	76,933,973
Sequence length NG50*	8,288,951	8,288,951	71,047,972	71,047,972
Sequence length LG50*	12	12	2	2
Normalized expected sequence length	8,831,279	8,830,144	81,601,468	81,599,549

470 **Expected genome size used rather than assembly length: 280 Mb (Urban et al 2021).*

471
472

473 **Table 2**

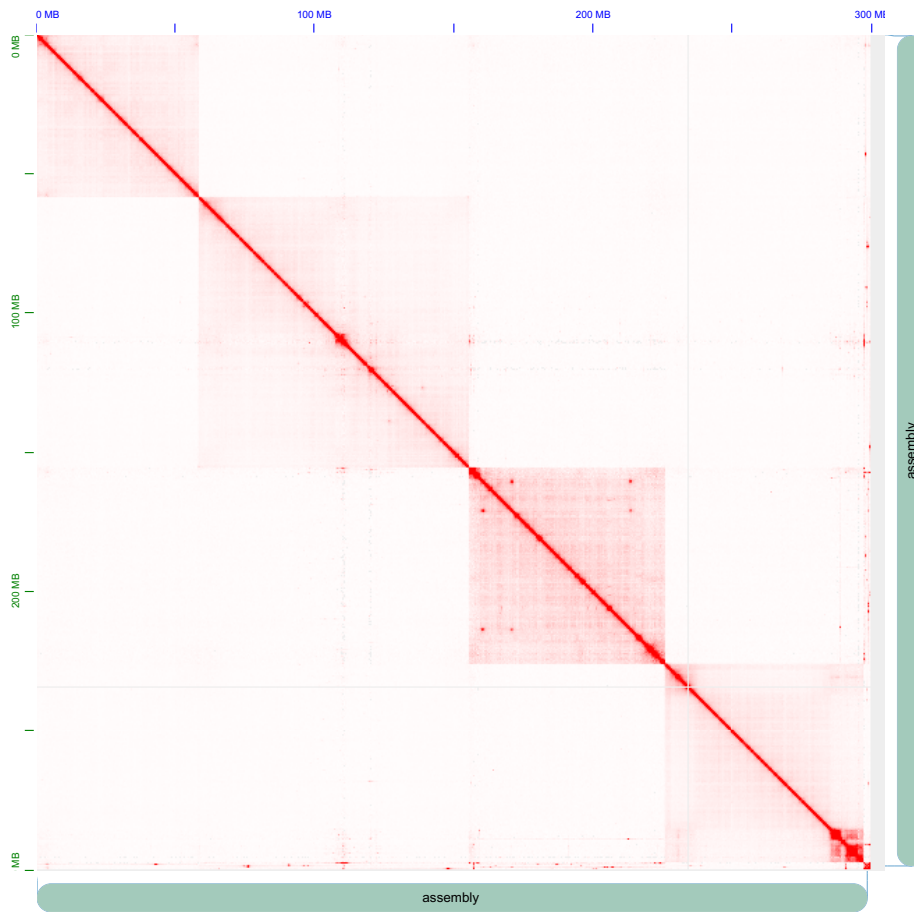
	Bcop_v1 (Urban et al 2021)	Bcop_v2 (This paper)	Expected size*
Total primary length classified as autosomal	223.8-232.8 Mb	226.5 Mb	218-231 Mb
Total length classified as X	62-71 Mb	70.5 Mb	48.9-62.2 Mb
Total length classified as II	13.1-28.5 Mb	58.3 Mb	62.2-66 Mb
Total length classified as III	5.4-12.5 Mb	71.0 Mb	66.6-70.7 Mb
Total length classified as IV	34.7-46.4 Mb	97.1 Mb	88.8-94.2 Mb

474 * *Expected size calculations are from Supplemental Table S1 in Urban et al 2021*

475

476

477 **Figure 1 – Hi-C interaction frequency map shows four chromosomes**



478