

Manuscript Number:	GIGA-D-17-00297	
Full Title:	Improving the annotation of the Heterorhabditis bacteriophora genome	
Article Type:	Research	
Funding Information:	Wellcome Trust (204052/Z/16/Z)	Dr Florence McLean
Abstract:	<p>Genome assembly and annotation remains an exacting task. As the tools available for these tasks improve, it is useful to return to data produced with earlier instances to assess their credibility and correctness. The entomopathogenic nematode Heterorhabditis bacteriophora is widely used to control insect pests in horticulture. The genome sequence for this species was reported to encode an unusually high proportion of unique proteins and a paucity of secreted proteins compared to other related nematodes. We revisited the H. bacteriophora genome assembly and gene predictions to ask whether these unusual characteristics were biological or methodological in origin. We mapped an independent resequencing dataset to the genome and used the blobtools pipeline to identify potential contaminants. While present (0.2% of the genome span, 0.4% of predicted proteins), assembly contamination was not significant. Re-prediction of the gene set using BRAKER1 and published transcriptome data generated a predicted proteome that was very different from the published one. The new gene set had a much reduced complement of unique proteins, better completeness values that were in line with other related species' genomes, and an increased number of proteins predicted to be secreted. It is thus likely that methodological issues drove the apparent uniqueness of the initial H. bacteriophora genome annotation and that similar contamination and misannotation issues affect other published genome assemblies.</p>	
Corresponding Author:	Florence McLean, BM BCh, MSc(R) University of Edinburgh Edinburgh, UNITED KINGDOM	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	University of Edinburgh	
Corresponding Author's Secondary Institution:		
First Author:	Florence McLean	
First Author Secondary Information:		
Order of Authors:	Florence McLean Duncan Berger Dominik R Laetsch Hillel T Schwartz Mark Blaxter	
Order of Authors Secondary Information:		
Opposed Reviewers:		
Additional Information:		
Question	Response	
Are you submitting this manuscript to a special series or article collection?	No	
Experimental design and statistics	Yes	

<p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

Improving the annotation of the *Heterorhabditis bacteriophora* genome

Florence McLean¹, Duncan Berger¹, Dominik R. Laetsch¹, Hillel T. Schwartz² and Mark Blaxter¹

¹ Institute of Evolutionary Biology, University of Edinburgh, Edinburgh EH9 3JT, UK

² Division of Biology and Biological Engineering, California Institute of Technology, Pasadena, California, USA

email addresses and ORCID

Florence McLean f.e.mclean@sms.ed.ac.uk 0000-0001-8218-4904

Duncan Berger duncan.berger@bath.edu 0000-0002-2017-6623

Dominik R. Laetsch dominik.laetsch@gmail.com 0000-0001-7887-0186

Hillel Schwartz hillels@caltech.edu 0000-0002-3448-8652

Mark Blaxter mark.blaxter@ed.ac.uk 0000-0003-2861-949X

Competing interests: The authors have no competing interests in relation to this manuscript.

30 **Abstract**

31 Genome assembly and annotation remains an exacting task. As the tools available for these
32 tasks improve, it is useful to return to data produced with earlier instances to assess their
33 credibility and correctness. The entomopathogenic nematode *Heterorhabditis bacteriophora* is
34 widely used to control insect pests in horticulture. The genome sequence for this species
35 was reported to encode an unusually high proportion of unique proteins and a paucity of
36 secreted proteins compared to other related nematodes. We revisited the *H. bacteriophora*
37 genome assembly and gene predictions to ask whether these unusual characteristics were
38 biological or methodological in origin. We mapped an independent resequencing dataset to
39 the genome and used the blobtools pipeline to identify potential contaminants. While
40 present (0.2% of the genome span, 0.4% of predicted proteins), assembly contamination was
41 not significant. Re-prediction of the gene set using BRAKER1 and published transcriptome
42 data generated a predicted proteome that was very different from the published one. The
43 new gene set had a much reduced complement of unique proteins, better completeness
44 values that were in line with other related species' genomes, and an increased number of
45 proteins predicted to be secreted. It is thus likely that methodological issues drove the
46 apparent uniqueness of the initial *H. bacteriophora* genome annotation and that similar
47 contamination and misannotation issues affect other published genome assemblies.

48

49

50

51

52

53 Introduction

54 The sequencing and annotation of a species' genome is often but the first step in exploiting
55 these data for comprehensive biological understanding. As with all scientific endeavour,
56 genome sequencing technologies and the bioinformatics toolkits available for assembly and
57 annotation are being continually improved. It should come as no surprise therefore that first
58 estimates of genome sequences and descriptions of the genes they contain can be improved.
59 For example, the genome of the nematode *Caenorhabditis elegans* was the first animal
60 genome to be sequenced [1]. The genome sequence and annotations have been updated
61 many times since, as further exploration of this model organism revealed errors in original
62 predictions, such that today, with release WS260 (<http://www.wormbase.org/>) [2], very few
63 of the 19099 protein coding genes announced in the original publication [1] retain their
64 original structure and sequence. The richness of the annotation of *C. elegans* is driven by the
65 size of the research community that uses this model species. However for most species,
66 where the community using the genome data is small or less-well funded, initial genome
67 sequences and gene predictions are not usually updated.

68 *Heterorhabditis bacteriophora* is an entomopathogenic nematode which maintains a mutualistic
69 association with the bacterium *Photorhabdus luminescens*. Unlike many other parasitic
70 nematodes, it is amenable to *in vitro* culture [3] and is therefore of interest not only to
71 evolutionary and molecular biologists investigating parasitic and symbiotic systems, but also
72 to those concerned with the biological control of insect pests [4, 5]. *P. luminescens* colonises
73 the anterior intestine of the free-living infective juvenile stage (Ij). Ijs are attracted to insect
74 prey by chemical signals [6, 7]. On contacting a host, the Ijs invade the insect's haemocoel
75 and actively regurgitate *P. luminescens* into the haemolymph. The bacterial infection rapidly
76 kills the insect, and *H. bacteriophora* grow and reproduce within the cadaver. After 2-3 cycles

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

77 of replication, the nematode progeny develop into IJs, sequester *P. luminescens* and seek out
78 new insect hosts.

79 Axenic *H. bacteriophora* IJs are unable to develop past the L1 stage [8] , and *H. bacteriophora*
80 may depend on *P. luminescens* for secondary metabolite provision [9, 10]. Mutation of the
81 global post-transcriptional regulator Hfq in *P. luminescens* reduced the bacterium's secondary
82 metabolite production and led to failed nematode development, despite the bacterium
83 maintaining virulence against host (*Galleria mellonella*) larvae [11]. Together these symbionts
84 are efficient killers of pest (and other) insects, and understanding of the molecular
85 mechanisms of host killing could lead to new insecticides.

86 *H. bacteriophora* was selected by the National Human Genome Research Initiative as a
87 sequencing target [12]. Genomic DNA from axenic cultures of the inbred strain *H.*
88 *bacteriophora* TTO1 was sequenced using Roche 454 technology and a high quality 77 Mb
89 draft genome assembly produced [13]. This assembly was predicted (using JIGSAW [14]) to
90 encode 21250 proteins. Almost half of these putative proteins had no significant similarity to
91 entries in the GenBank non-redundant protein database, suggesting an explosion of novelty
92 in this nematode. The predicted *H. bacteriophora* proteome had fewer orthologues of Kyoto
93 Encyclopedia of Genes and Genomes loci in the majority of metabolic categories than nine
94 other nematodes. *H. bacteriophora* was also predicted to have a relative paucity of secreted
95 proteins compared to free-living nematodes, postulated to reflect a reliance on *P.*
96 *luminescens* for secreted effectors [13]. The 5.7 Mb genome of *P. luminescens* has also been
97 sequenced [15]. The *H. bacteriophora* proteome had fewer shared orthologues when
98 clustered and compared to other rhabditine (Clade V) nematodes (including *Caenorhabditis*
99 *elegans* and the many animal parasites of the Strongylomorpha) [16].

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

I00 In preliminary analyses we noted that while the genome sequence itself had high
I01 completeness scores when assessed with the Core Eukaryote Gene Mapping Approach
I02 (CEGMA) [17] (99.6% complete) and Benchmarking Universal Single-Copy Orthologs
I03 (BUSCO) [18] (80.9% complete and 5.6% fragmented hits for the BUSCO Eukaryota gene
I04 set), the predicted proteome scored poorly (47.8% complete and 34.7% fragmented by
I05 BUSCO; see below). Another unusual feature of the *H. bacteriophora* gene set was the
I06 proportion of non-canonical splice sites (i.e. those with a 5' GC splice donor site, as
I07 opposed to the normal 5' GT). Most nematode (and other metazoan) genomes have low
I08 proportions of non-canonical introns (less than 1%), but the published gene models had over
I09 9% non-canonical introns. This is more than double the proportion predicted for *Globodera*
I10 *rostochiensis*, a plant parasitic nematode where the unusually high proportion of non-
I11 canonical introns was validated *via* manual curation [19].

I12 If these unusual characteristics reflect a truly divergent proteome, the novel proteins in *H.*
I13 *bacteriophora* may be crucial in its particular symbiotic and parasitic relationships, and of
I14 great interest to development of improved strains for horticulture. However, it is also
I15 possible that contamination of the published assembly or annotation artefacts underpin
I16 these unusual features. We re-examined the *H. bacteriophora* genome and gene predictions,
I17 and used more recent tools to re-predict protein coding genes from the validated assembly.
I18 As the BRAKER1 predictions were demonstrably better than the original ones, we explored
I19 whether some of the unusual characteristics of the published protein set, in particular the
I20 level of novelty and the proportion of secreted proteins, were supported by the BRAKER1
I21 protein set.

I22
I23

124 Results

125 No evidence for substantial contamination of the *H. bacteriophora* genome assembly

126 We used BlobTools [20] to assess the published genome sequence [13] for potential
127 contamination. The raw read data from the published assembly was not available on the
128 trace archive or short read archive (SRA). We thus utilised new Illumina short-read re-
129 sequencing data generated from strain G2a1223, an inbred derivative of *H. bacteriophora*
130 strain "Gebre", isolated by Adler Dillman in Moldova. G2a1223 has about 1 single-nucleotide
131 change per ~2000 nucleotides compared to the originally-sequenced TT01 strain. G2a1223
132 was grown in culture on the non-colonising bacterium *Photorhabdus temperata*. The majority
133 of these data (96.3% of the reads) mapped as pairs to the assembly, suggesting completeness
134 of the published assembly with respect to the new raw read data. In addition, 99.96% of the
135 published assembly had at least 10-fold coverage from the new raw reads.

136 The assembly was explored using a taxon-annotated GC-coverage plot, with coverage taken
137 from the new Illumina data and sequence similarity from the NCBI nucleotide database (nt)
138 (Figure 1). *H. bacteriophora* was excluded from the database search used to annotate the
139 scaffolds to exclude self hits from the published assembly. All large scaffolds clustered
140 congruently with respect to read coverage and CG content. A few (57) scaffolds had best
141 BLASTn matches to phyla other than Nematoda (Table 1). A small amount (5 kb) of likely
142 remaining *P. luminescens* contamination was noted. We identified 100 kb of the genome of a
143 strain of the common culture contaminant bacterium *Stenotrophomonas maltophilia* [21].
144 Contamination of the assembly with *S. maltophilia* was acknowledged [13] but removal of
145 scaffolds before annotation was not discussed. Two high-coverage scaffolds that derived
146 from the *H. bacteriophora* mitochondrial genome were annotated as "undefined Eukaryota"
147 because of taxonomic misclassification in the NCBI nt database. Many scaffolds with

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

148 coverages close to that of the expected nuclear genome had best matches to two
149 unexpected sources: the platyhelminths *Echinostoma caproni* and *Dicrocoelium dendriticum*, and
150 several hymenopteran arthropods. Inspection of these matches showed that they were due
151 to high sequence similarity to a family of *H. bacteriophora* mariner-like transposons [22] and
152 thus these were classified as *bona fide* nematode nuclear sequences. A group of scaffolds
153 contained what appears to be a *H. bacteriophora* nuclear repeat with highest similarity to
154 histone H3.3 sequences from Diptera and Hymenoptera. The remaining scaffolds had low-
155 scoring nucleotide matches to a variety of chordate, chytrid and arthropod sequences from
156 deeply conserved genes (tubulin, kinases), but had coverages similar to other nuclear
157 sequences.

158 Scaffolds with average coverage of less than 10-fold were removed from the assembly (35
159 scaffolds spanning 132949 bases, 0.2% of the total span; see Supplementary File 1). This
160 removed all scaffolds aligning to *S. maltophilia* and to *Photorhabdus* spp. (104 kb). The origins
161 of the additional 28 kb were not investigated. In the published annotation [13], 76 genes
162 were predicted from these scaffolds.

163

164 **Improved gene predictions are biologically credible and have unexceptional novelty**

165 New gene predictions were generated from a soft-masked version of the filtered assembly
166 using the RNA-seq based annotation pipeline BRAKER1 [23], generating 16070 protein
167 predictions from 15747 protein coding genes (see Supplementary File 2). We compared the
168 soft-masked predictions to those from the published analysis [13] (Figure 2, Table 2). The
169 predicted proteins from the new BRAKER1/soft-masked gene set were, on average, longer
170 (Figure 2A). While the average number of introns per gene was the same in the
171 BRAKER1/soft-masked and published predictions, the BRAKER1/soft-masked gene set had

172 more single-exon genes (Figure 2B). Hard masking of the genome and re-prediction resulted
173 in fewer single exon genes, suggesting that many of these putative genes could be derived
174 from repetitive sequence (Supplementary Files 3 and 4), but only 316 of the single exon
175 genes from the BRAKER1/soft-masked assembly had similarity to transposases or
176 transposons. The BRAKER1/soft-masked annotations were taken forward for further
177 analysis.

178 Four-fifths (83.3%) of the published protein-coding gene predictions [13] overlapped to
179 some extent with the BRAKER1/soft-masked predictions at the genome level, with a mean
180 of 67% of the nucleotides of each BRAKER1/soft-masked gene covered by a published gene
181 (Figure 2C). Half (8061) of the 15747 BRAKER1/soft-masked gene predictions had an
182 overlap proportion of ≥ 0.9 with the published predictions. At the level of protein sequence
183 only 836 proteins were identical between the two predictions, and only 2099 genes had
184 identical genome start and stop positions.

185 The BRAKER1/soft-masked and published gene sets were checked for completeness using
186 BUSCO [18], based on the Eukaryota lineage gene set, and *Caenorhabditis* as the species
187 parameter for orthologue finding. The BRAKER1/soft-masked gene set contained a
188 substantially higher percentage of complete, and lower percentage of fragmented BUSCO
189 genes than the published set (Table 2). Two *H. bacteriophora* transcriptome datasets, publicly
190 available Roche 454 data and Sanger expressed sequence tags, were mapped to the
191 published and BRAKER1/soft-masked transcriptomes to assess gene set completeness. This
192 suggested that the BRAKER1/soft-masked transcriptome predictions were more complete
193 than the original (Table 2).

194 Nearly half (9893/20964; 47.2%) of the published proteins were reported to have no
195 significant matches in the NCBI non-redundant protein database (nr) [13]. This surprising

196 result could be due to a paucity of data from species closely related to *H. bacteriophora* in
197 the NCBI nr database at the time of the search, or inclusion of poor protein predictions in
198 the published set, or both. Targeted investigation of these 9893 orphan proteins here was
199 not possible due to inconsistencies in gene naming in the publically available files. The
200 published and BRAKER1/soft-masked proteomes were compared to the Uniref90 database
201 [24], using DIAMOND v0.9.5 [25] with an expectation value cut-off of $1e^{-5}$. In the published
202 proteome, 8962 proteins (42.7%) had no significant matches in Uniref90. Thus a relatively
203 poorly populated database was not the main driver for the high number of orphan proteins
204 reported in the published proteome. In the BRAKER1/soft-masked proteome, only 2889
205 proteins (18.3%) had no hits in the Uniref90 database (Table 2).

206 OrthoFinder v1.1.4 [26] was used to define orthologous groups in the proteomes of 23
207 rhabditine (Clade V) nematodes (Supplementary Files 5 and 6) and just the published *H.*
208 *bacteriophora* protein-coding gene predictions, or just the BRAKER/soft-masked proteome,
209 or both. All proteins <30 amino-acids long were excluded from clustering (see
210 Supplementary File 5). We identified 5442 singletons (26.8% of the proteome) when the
211 analysis included only the published *H. bacteriophora* protein set. An additional 248 proteins
212 formed *H. bacteriophora*-specific orthogroups. Orthology analysis including only the
213 BRAKER/soft-masked protein set predicted 1112 *H. bacteriophora* singletons (7.1% of the
214 proteome) with 167 proteins in *H. bacteriophora*-specific orthogroups (Figure 2D). In
215 comparison, when the orthology analysis included the BRAKER1/soft-masked predictions
216 there were 1858 *C. elegans* singletons (9.2% of the *C. elegans* proteome). Very few universal,
217 single copy orthologues were defined in either analysis. Exploring “fuzzy-1-to-1”
218 orthogroups (where true 1-to-1 orthology was found for greater than 75% of the 24 species
219 - i.e. 18 or more species), the published protein predictions had more missing fuzzy-1-to-1
220 orthologues than did the BRAKER1/soft-masked predictions (Table 2). In the clustering that

221 included both proteomes, 2019 clusters contained more proteins from the BRAKER1/soft-
222 masked than the published proteome, whereas 2714 contained a larger number contributed
223 from the published than the BRAKER1/soft-masked proteome (Supplementary File 6).

224 The published *H. bacteriophora* gene set had additional peculiarities. The published set of
225 gene models included 102274 introns, 9069 of which (8.9%) had non-canonical splice sites
226 (i.e. 5' GC – AG 3'). Some of the genes in the published gene set had up to nine
227 noncanonical introns (Figure 2E). In the BRAKER1/soft-masked gene set there were 109767
228 introns, 868 (0.8%) of which had non-canonical splice sites. This proportion is in keeping
229 with that found in most other rhabditine nematodes. For example, the extensively manually
230 annotated *C. elegans* has 2429 (0.6%) non-canonical (5' GC – AG 3') introns. In *C. elegans*
231 non-canonical introns are frequently found only in alternately spliced, and shorter isoforms,
232 and over 93-99% were in genes that had homologues in other species, depending on the
233 species used in the protein orthology clustering. However, in the published *H. bacteriophora*
234 gene set, 34-49% of the genes with GC – AG introns were in *H. bacteriophora*-unique
235 proteins.

236 A supermatrix maximum likelihood phylogeny was generated from the fuzzy-I-I
237 orthologues in the clustering that included both *H. bacteriophora* proteomes (Figure 3; see
238 Supplementary File 7). The phylogeny, rooted with *Pristionchus* spp., shows the *H.*
239 *bacteriophora* proteomes as sisters. However the BRAKER1/soft-masked proteome has a
240 shorter branch length to *Heterorhabditis*' most recent common ancestor with other Clade V
241 nematodes, suggesting that the published proteome includes uniquely divergent sequences.

242 The secretome of *H. bacteriophora* has been of particular interest as it may contain proteins
243 involved in symbiotic interactions with *P. luminescens*, and proteins crucial to invasion and
244 survival within the insect haemocoel. In the original publication, only 603 proteins (2.8% of

245 the proteome) were predicted to be secreted [13]. This proportion is much lower than in
246 free living nematodes such as *C. elegans* and it was postulated that *H. bacteriophora* relies on
247 *P. luminescens* for secreted effectors [13]. The signal peptide detection method used in the
248 original analyses was not described [13]. We used SignalP version 4.1 within Interproscan to
249 annotate proteins in both the BRAKER1 and published *H. bacteriophora* proteomes. Proteins
250 having a predicted signal peptide but no transmembrane domain were classified as secreted.
251 We identified 1023 (6.5%) putative secreted proteins in the BRAKER1/soft-masked
252 proteome and 1065 (5.1%) in the published proteome. By the same method other
253 rhabditine (Clade V) nematodes that do not have known symbiotic associations with
254 bacteria, such as *Teladorsagia circumcincta*, had comparable secretome sizes to *H.*
255 *bacteriophora* (Supplementary File 8). This suggests that *H. bacteriophora* does not have a
256 reduced secretome compared to other, related nematodes that do not have symbiont
257 partners.

258

259

260

261

262

263

264 Discussion

265 Assembly of, and gene finding in, new genomes is a challenging task, and especially so in
266 larger genomes and those phylogenetically distant from any previously analysed exemplar.
267 When applied *de novo* to datasets from extremely well-assembled and well-annotated model
268 species, even the best methods fail to recover fully contiguous assemblies and yield

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

269 predicted gene sets that have poor correspondence with the known truth [27]. A major
270 issue with primary assemblies and gene sets arises when exceptional findings are taken at
271 face value, and used to assert exceptional biology in a target species [28]. Where these
272 exceptions are in fact the result of methodological failings, the scientific record, including the
273 public databases, becomes contaminated. At best, erroneous assertions can be quickly
274 checked and corrected, but at worst they can mislead and inhibit subsequent work.

275 A second concern arises from the recognition that while no method can currently produce
276 perfect assemblies and perfect gene sets from raw data, analyses using the same toolsets will
277 resemble each other and reflect the successes and failings of the particulars of the
278 algorithms employed. However, when comparing genome assemblies and gene sets
279 produced by different pipelines, it may be that the disparity in output generated by different
280 pipelines dominates any signal from biology. Genomes assembled and annotated with the
281 same tools will look more similar, and in a pool of assemblies and protein sets the one
282 species that used a variant process will be flagged as exceptional. Again, the model
283 organisms show the way: as new data and new scrutiny is added to the genome, better and
284 better analyses are available. With additional analysis, and additional independent data,
285 genome and gene predictions can be improved markedly for any species [29].

286 Here we examined the “outlier” whole-genome protein predictions from the
287 entomopathogenic nematode *H. bacteriophora* [13]. The original publication noted that the
288 number of novel proteins (those restricted to *H. bacteriophora*) was particularly large, while
289 the number of secreted proteins was rather small, and suggested that these genome
290 features might be a result of evolution to the species’ novel lifestyle (which includes an
291 essential symbiosis with the bacterium *P. luminescens*). Overall we found that while the
292 published genome sequence had a small amount of bacterial contamination, and a small

293 number of “nematode” genes were predicted from these contaminants, the assembly itself
294 was of high quality. Our re-prediction of the gene set of *H. bacteriophora* however suggested
295 that the excess of unique genes, the lack of secreted proteins and several other surprising
296 features of the original gene set were likely to be artefacts of the gene prediction pipeline
297 chosen. While our gene set was by no means perfect (for example we identified an excess
298 of single exon genes that derive from likely repetitive sequence) it had better biological
299 completeness and credibility.

300 We used the RNA-seq based annotation pipeline BRAKER1 [23], not available to the
301 authors of the original genome publication, who used JIGSAW [14] (see Supplementary File
302 9). While JIGSAW achieved high sensitivity and specificity at the level of nucleotide, exon
303 and gene predictions in the nematode genome annotation assessment project, nGASP [27],
304 direct comparison of the sensitivity and specificity of JIGSAW and BRAKER1 has not been
305 published to the best of our knowledge. BRAKER1 has been shown to give superior
306 prediction results over *ab initio* GeneMark-ES, or *ab initio* AUGUSTUS alone [23]. In
307 particular, BRAKER1 is able to better use transcriptome data for gene finding. While we
308 supplied only low volumes of Sanger-sequenced ESTs and a partial Roche 454 transcriptome
309 to BRAKER1, the resulting gene set has much improved numerical and biological scores. In
310 particular we note that the biological completeness of the predicted gene set now matches
311 that of the genome sequence from which it was derived (Table 2).

312 The published gene set had an unusually high proportion (8.9%) of non-canonical (5' GC –
313 AG 3') introns. While most genomes have a low proportion of non-canonical introns
314 (usually approximately 0.5% of all introns), some species have markedly higher proportions
315 [19]. The high proportion found initially in *H. bacteriophora* could perhaps have been taken
316 as a warning that the prediction set was of concern. We note that gene predictors can be

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

317 set to disallow any predictions that require non-canonical splicing, and many published
318 genomes have zero non-canonical introns. These gene prediction sets are likely to
319 categorically miss true non-canonically spliced genes.

320 The new BRAKER1 gene prediction set had many fewer species-unique genes (7.1%) than
321 did the original (42.7%) when compared to 23 other related nematodes. We regard this
322 reduction in novelty as indicative of a better prediction, as, for example, *C. elegans*, the best-
323 annotated nematode genome, had only 9.2% of species unique genes in our analysis. Having
324 a large proportion of orphan proteins is not unique to the published *H. bacteriophora*
325 predictions. Nearly half (47%) of the gene predictions in *Pristionchus pacificus* were reported
326 to have no homologues in fifteen other nematode species [30]. Evaluation of proteomic and
327 transcriptomic evidence, as well as patterns of synonymous and non-synonymous
328 substitution, suggested that as many as 42-81% of these genes were in fact expressed [31].
329 Therefore the high proportion of orphan genes in *H. bacteriophora* is not *prima facie*
330 evidence of poor gene predictions. Expanded transcriptomic and comparative data are
331 needed to build on the work we have presented in affirming the true *H. bacteriophora* gene
332 set.

333 Biological pest control agents may become increasingly important for ensuring crop
334 protection in the future [32]. A number of factors currently limit the commercial
335 applicability of *H. bacteriophora*, including their short shelf life, susceptibility to environmental
336 stress and limited insect tropism [12, 33]. Accurate genome annotation will assist in the
337 analysis of *H. bacteriophora*, facilitating the exploration of genes involved in its parasitic and
338 symbiotic interactions, and supporting genetic manipulation to enhance its utility as a
339 biological control agent.

340 **Acknowledgements**

1
2
3 341 This project was supported by FMs Wellcome Trust-funded graduate programme
4
5 342 [204052/Z/16/Z]. Sujai Kumar, Lewis Stevens, Carlos Caurcel and Elisabeth Sjokvist offered
6
7
8 343 expert technical support and advice. Igor Antoshechkin of the Millard and Muriel Jacobs
9
10 344 Genetics and Genomics Laboratory at Caltech assisted with Illumina sequencing. Adler
11
12
13 345 Dillman provided the parental strain for the inbred *H. bacteriophora* strain G2a1223.
14

15 346

16
17
18 347

19
20
21
22 348

23
24
25 349

26
27
28 350

29
30
31 351

32
33
34
35 352

36
37
38 353

39
40
41
42 354

43 44 45 355 **Methods**

46 47 48 49 356 **Input data and data availability**

50
51
52 357 The *H. bacteriophora* genome and annotations [13] were downloaded from Wormbase
53
54
55 358 Parasite (WBPS8) [34]. ESTs [35, 36] were obtained from NCBI dbEST [37]. Roche 454
56
57 359 transcriptome data [13] were obtained from the Short Read Archive. *H. bacteriophora* strain
58
59
60 360 Gebre, a gift from Adler Dillman, was inbred by selfing single hermaphrodites for five

361 generations to generate the strain G2a1223. New Illumina HiSeq2000, paired end, 75 base
362 data were generated from *H. bacteriophora* G2a1223 genomic DNA by the Millard and
363 Muriel Jacobs Genetics and Genomics Laboratory at Caltech. They have been deposited in
364 SRA under XXXXXXXX [in process].

365

366 The revised gene annotations for *H. bacteriophora* have been submitted to the INSDC under
367 project XXXXXXXX [to be advised]. The Supplementary files for this manuscript are
368 additionally available at <https://github.com/DRL/mclean2017>. All custom scripts developed
369 for this manuscript are available at <https://github.com/DRL/mclean2017>.

370 **Contaminant screening and Removal of Low Coverage Scaffolds**

371 The assembly scaffolds were aligned to the NCBI nucleotide (nt) database, release 204, using
372 Nucleotide-Nucleotide BLAST v2.6.0+ (RRID:SCR_008419) in megablast mode, with an e-
373 value cut off of $1e^{-25}$ and a culling limit of 2 [38]. *H. bacteriophora* hits were excluded from
374 the search using a list of all *H. bacteriophora* associated gene identifiers downloaded from
375 NCBI GenBank nucleotide database, release 219. Raw, paired-end Illumina reads from the
376 re-sequencing project were mapped against the assembly, as paired, using Burrows-Wheeler
377 Aligner (BWA) v0.7.15 (RRID:SCR_010910) in mem mode with default options [39]. The
378 output was converted to a BAM file using Samtools v1.3.1 (RRID:SCR_002105) [40] and
379 overall mapping statistics generated in flagstat mode.

380 Blobtools v0.9.19 [20] was used to create taxon annotated GC-coverage plots for the
381 published assembly, using the Nucleotide-Nucleotide BLAST and raw read mapping results.
382 Scaffolds that did not have Nematoda as a top BLAST hit at the phylum level were identified,
383 and the species-level top BLAST hit, length of scaffold, and scaffold mean base coverage

384 were extracted from the Blobology output. Scaffolds with a mean base coverage of <10x
385 were identified from the output of the Blobology pipeline and removed from the assembly.
386 A list of excluded scaffolds is available in Supplementary File 1.

387 **Generation of BRAKER1 Gene Predictions**

388 Before annotation the published assembly was soft masked for known Nematoda repeats
389 from the RepeatMasker Library v4.0.6 using RepeatMasker v4.0.6 (RRID:SCR_012954) [41]
390 with default options. The two publicly available Roche 454 RNA-seq data files were adaptor
391 and quality-trimmed using BBDuk v36.92 (unpublished toolkit from Joint Genome Institute,
392 n.d.). Reads below an average quality of 10 or shorter than 25 nucleotides were discarded.
393 Regions with average quality below 20 were trimmed. The cleaned reads were mapped to
394 the soft masked assembly using STAR v2.5 (RRID:SCR_005622) with default options [42,
395 43]. The soft masked assembly was annotated with BRAKER1 [23] with guidance from the
396 mapping output from STAR. An identical annotation method was applied to a hard masked
397 version of the assembly. Hard masking was for known Nematoda repeats from the
398 RepeatMasker Library v4.0.6 using RepeatMasker v4.0.6 with default options. The published
399 and BRAKER1 proteomes were compared using DIAMOND v0.9.5 [25] in BLASTP mode to
400 the Uniref90 database (release 03/2017) [24] with an expectation value cut-off of $1e^{-5}$ and
401 no limit on the number of target sequences. Hits to *H. bacteriophora* proteins were
402 removed using its TaxonID.

403 **Gene Prediction Statistics**

404 Gene-level statistical summaries were calculated including only the longest isoforms of the
405 BRAKER1 gene predictions. The longest isoform for each gene in the BRAKER1 *H.*
406 *bacteriophora* annotation was identified from the general feature format file, and then

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

407 selected from the protein FASTA files. The general feature format file (GFF) for the
408 published gene predictions did not contain any isoforms and was analysed in its entirety.
409 Mean protein lengths were calculated from the amino-acid protein sequence files. Introns
410 were inferred for the published GFF file using GenomeTools v1.5.9 in -addintrons mode
411 [44]. Intron frequencies were then calculated for the published and BRAKER1 annotations
412 from their respective GFF files. Exon frequencies were calculated for the published
413 annotations directly from the GFF file. For the BRAKER1 annotations, exon frequency per
414 gene was assumed to be equivalent to coding DNA sequence (CDS) frequency and inferred
415 from the general feature format file as exon features were not included in the GFF. Intron
416 frequency histograms and bar plots were generated in Rstudio v1.0.136 (RRID:SCR_005622)
417 with R v3.3.2 (RRID:SCR_001905) and in some instances the package ggplot2 v2.2.1. As
418 intron frequency lists did not contain single exon genes (those with no introns), these were
419 added manually to the intron frequency lists in Microsoft Excel before importing the data
420 into Rstudio.

421 The proportion of introns with GC – AG splice junctions was assessed for the gene models
422 of *C. elegans* (WS258), and the published and BRAKER1/soft-masked gene models of *H.*
423 *bacteriophora*. Intronic features were added to GFF3 files using GenomeTools v1.5.9 [44] ('gt
424 gff3 -sort -tidy -retainids -fixregionboundaries -addintrons') and splice sites were
425 extracted using the script extractRegionFromCoordinates.py [19]. Results were visualised
426 using the script plot_GCAG_counts.R (see <https://github.com/DRL/mclean2017>).

427 Gene features, extracted from the GFF files, were assessed for overlap using bedtools v2.26
428 (RRID:SCR_006646) in intersect mode [45]. Only genes on the same strand were
429 considered to be overlapping. To calculate the number of identical proteins shared between
430 the published and BRAKER1 proteomes non-redundant protein fasta files were generated

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

431 using cd-hit v4.6.1 (RRID:SCR_007105) [46] for the BRAKER1 and published predictions.

432 The files were concatenated, sorted and unique sequences counted using unix command line
433 tools.

434 BUSCO v2.0.1 (RRID:SCR_015008) [18], with Eukaryota as the lineage dataset, and

435 *Caenorhabditis* as the species parameter for orthologue finding was applied to both

436 proteomes and the published assembly to calculate BUSCO scores. CEGMA

437 (RRID:SCR_015055) [17] was run on the published genome sequence. BWA was used with

438 default settings to map the RNA-seq datasets to the CDS transcripts from the published and

439 BRAKER1 annotations and the summary statistics obtained with Samtools v1.3.1 in flagstat

440 mode.

441 **Protein orthology analyses**

442 OrthoFinder v1.1.4 [26] with default settings was used to identify orthologous groups in the

443 proteomes of 23 Clade V nematodes with the addition of either the BRAKER1/soft-masked

444 and published *H. bacteriophora* proteomes separately or simultaneously. The proteomes for

445 the 23 Clade V nematodes were downloaded from WBPS8 (available at:

446 <http://parasite.wormbase.org/index.html>) or GenomeHubs.org (available at

447 <http://ensembl.caenorhabditis.org/index.html>), and detailed source information is available in

448 Supplementary File 5. All proteomes were filtered to contain only the longest isoform of

449 each gene, and for all proteomes (except the BRAKER1/soft-masked *H. bacteriophora*

450 protein set), proteins less than 30 amino-acids in length were excluded before clustering.

451 For the *H. bacteriophora* BRAKER1/soft-masked protein set, proteins less than 30 amino-

452 acids (SF5.2) were removed manually from the orthofinder clustering statistics after

453 clustering. None of these proteins seeded new clusters and are therefore will not have

454 influenced the clustering results. Kinfin v0.9 [47], was used with default settings to identify

1
2 455 true and fuzzy 1-to-1 orthologues, and their associated species specific statistics. Fuzzy 1-to-
3
4 456 1 orthologues are true 1-to-1 orthologues for greater than 75% of the species clustered.
5
6 457 For the clustering analysis presented in Supplementary File 3, the BRAKER1/soft masked and
7
8 458 published proteomes were clustered simultaneously to the 23 other Clade V nematode
9
10 459 proteomes, and singletons, and species-specific clusters were excluded.

11 12 13 460 **Interproscan and search for transposons**

14
15
16 461 Interproscan v5.19-58.0 (RRID:SCR_005829) [48] was used in protein mode to identify
17
18 462 matches with the BRAKER1 and published *H. bacteriophora* predicted proteomes in the
19
20 463 following databases: TIGRFAM v15.0, ProDom v2006.1, SMART-7.1, SignalP-EUK v4.1,
21
22 464 PrositePatterns v20.119, PRINTS v42.0, SuperFamily v1.75, Pfam v29.0, and PrositeProfiles
23
24 465 v20.119. InterProScan was run with the option for all match calculations to be run locally
25
26 466 and with gene ontology annotation activated. The number of single exon genes with
27
28 467 similarity to transposons or transposases in the BRAKER1/soft masked predictions was
29
30 468 calculated by searching the full InterProScan results for the strings 'Transposon',
31
32 469 'transposon', 'Transposase', or 'transposase' and the number of single exon gene
33
34 470 InterProScan results containing these terms counted. InterProScan results from searching
35
36 471 the SignalP-EUK-4.1 database were queried to identify putative secreted proteins. Those
37
38 472 with a predicted signal peptide but no transmembrane region were considered to be
39
40 473 secreted.

41 42 43 44 45 46 47 48 49 474 **Phylogenetic Analyses**

50
51
52
53 475 Both *H. bacteriophora* proteomes were clustered simultaneously with the 23 Clade V
54
55 476 nematode proteomes into orthologous groups using Orthofinder v1.0 [26]. The fuzzy 1-to-1
56
57 477 orthologues were extracted and processed using GNU parallel [49]. They were aligned

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

478 using MAFFT v7.267 (RRID:SCR_011811) [50], and the alignments trimmed with NOISY
479 v1.5.12. A maximum likelihood gene tree was generated for each orthologue using RaXML
480 v8.1.20 (RRID:SCR_006086) [51] with a PROTGAMMAGTR amino-acid substitution model.
481 Rapid Bootstrap analysis and search for the best -scoring ML tree within one program run
482 with 100 rapid bootstrap replicates was used. The trees were pruned using
483 PhyloTreePruner v1.0 [52] to remove paralogues, with 0.5 as the bootstrap cutoff and a
484 minimum of 20 species in the orthogroup after pruning for inclusion in the supermatrix.
485 Where species had more than one putative orthologue in an orthogroup the longest was
486 selected. The remaining 897 orthogroups were re-aligned using MAFFT v7.267, trimmed
487 with NOISY v1.5.12 and concatenated into a supermatrix using FASconCAT v.1.0 [53]. A
488 supermatrix maximum-likelihood tree was generated using RAxML with the rapid hill
489 climbing algorithm (default), with a PROTGAMMAGTR amino-acid substitution model and
490 100 bootstrap replicates. *Pristionchus* spp. were designated as the outgroup. The tree was
491 visualised in Dendroscope v3.5.9 [54].

492

493

494

495 **Bibliography**

496 1. *C. elegans* Sequencing Consortium. Genome sequence of the nematode *C. elegans*: a
497 platform for investigating biology. *Science*. 1998;282:2012–8.

498 2. Howe KL, Bolt BJ, Cain S, Chan J, Chen WJ, Davis P, et al. WormBase 2016: expanding to
499 enable helminth genomic research. *Nucleic Acids Res*. 2016;44:D774-80.

500 doi:10.1093/nar/gkv1217.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
- 501 3. Gil GH, Choo HY, Gaugler R. Enhancement of entomopathogenic nematode production
502 in in-vitro liquid culture of *Heterorhabditis bacteriophora* by fed-batch culture with glucose
503 supplementation. Appl Microbiol Biotechnol. 2002;58:751–5. doi:10.1007/s00253-002-0956-
504 1.
 - 505 4. Memari Z, Karimi J, Kamali S, Goldansaz SH, Hosseini M. Are Entomopathogenic
506 Nematodes Effective Biological Control Agents Against the Carob Moth, *Ectomyelois*
507 *ceratoniae*? J Nematol. 2016;48:261–7.
 - 508 5. Rezaei N, Karimi J, Hosseini M, Goldani M, Campos-Herrera R. Pathogenicity of Two
509 Species of Entomopathogenic Nematodes Against the Greenhouse Whitefly, *Trialeurodes*
510 *vaporariorum* (Hemiptera: Aleyrodidae), in Laboratory and Greenhouse Experiments. J
511 Nematol. 2015;47:60–6.
 - 512 6. Dillman AR, Guillermin ML, Lee JH, Kim B, Sternberg PW, Hallem EA. Olfaction shapes
513 host-parasite interactions in parasitic nematodes. Proc Natl Acad Sci U S A.
514 2012;109:E2324-33. doi:10.1073/pnas.1211436109.
 - 515 7. Anbesse S, Ehlers RU. Attraction of *Heterorhabditis* sp. toward synthetic (E)-beta-
516 caryophyllene, a plant SOS signal emitted by maize on feeding by larvae of *Diabrotica virgifera*
517 *virgifera*. Commun Agric Appl Biol Sci. 2010;75:455–8.
 - 518 8. Han R, Ehlers RU. Pathogenicity, development, and reproduction of *Heterorhabditis*
519 *bacteriophora* and *Steinernema carpocapsae* under axenic in vivo conditions. J Invertebr Pathol.
520 2000;75:55–8. doi:10.1006/jjpa.1999.4900.
 - 521 9. Ciche TA, Bintrim SB, Horswill AR, Ensign JC. A Phosphopantetheinyl transferase
522 homolog is essential for *Photorhabdus luminescens* to support growth and reproduction of
523 the entomopathogenic nematode *Heterorhabditis bacteriophora*. J Bacteriol. 2001;183:3117–
524 26. doi:10.1128/JB.183.10.3117-3126.2001.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
- 525 10. Bennett HPJ, Clarke DJ. The pbgPE operon in *Photorhabdus luminescens* is required for
526 pathogenicity and symbiosis. *J Bacteriol.* 2005;187:77–84. doi:10.1128/JB.187.1.77-84.2005.
- 527 11. Tobias NJ, Heinrich AK, Eresmann H, Wright PR, Neubacher N, Backofen R, et al.
528 *Photorhabdus*-nematode symbiosis is dependent on hfq-mediated regulation of secondary
529 metabolites. *Environ Microbiol.* 2017;19:119–29. doi:10.1111/1462-2920.13502.
- 530 12. Ciche T. The biology and genome of *Heterorhabditis bacteriophora*. *WormBook.* 2007;1–
531 9. doi:10.1895/wormbook.1.135.1.
- 532 13. Bai X, Adams BJ, Ciche TA, Clifton S, Gaugler R, Kim K, et al. A lover and a fighter: the
533 genome sequence of an entomopathogenic nematode *Heterorhabditis bacteriophora*. *PLoS*
534 *ONE.* 2013;8:e69618. doi:10.1371/journal.pone.0069618.
- 535 14. Allen JE, Salzberg SL. JIGSAW: integration of multiple sources of evidence for gene
536 prediction. *Bioinformatics.* 2005;21:3596–603. doi:10.1093/bioinformatics/bti609.
- 537 15. Duchaud E, Rusniok C, Frangeul L, Buchrieser C, Givaudan A, Taourit S, et al. The
538 genome sequence of the entomopathogenic bacterium *Photorhabdus luminescens*. *Nat*
539 *Biotechnol.* 2003;21:1307–13. doi:10.1038/nbt886.
- 540 16. Blaxter M, Koutsovoulos G. The evolution of parasitism in Nematoda. *Parasitology.*
541 2015;142 Suppl 1:S26-39. doi:10.1017/S0031182014000791.
- 542 17. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in
543 eukaryotic genomes. *Bioinformatics.* 2007;23:1061–7. doi:10.1093/bioinformatics/btm071.
- 544 18. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:
545 assessing genome assembly and annotation completeness with single-copy orthologs.
546 *Bioinformatics.* 2015;31:3210–2. doi:10.1093/bioinformatics/btv351.
- 547 19. Eves-van den Akker S, Laetsch DR, Thorpe P, Lilley CJ, Danchin EGJ, Da Rocha M, et al.

- 1
2 549 The genome of the yellow potato cyst nematode, *Globodera rostochiensis*, reveals insights into
3 the basis of parasitism and virulence. *Genome Biol.* 2016;17:124. doi:10.1186/s13059-016-
4 0985-1.
5 550
6
7 551 20. Laetsch DR, Blaxter ML. BlobTools: Interrogation of genome assemblies [version 1;
8 referees: 1 approved with reservations]. *F1000Res.* 2017;6:1287.
9
10 552 doi:10.12688/f1000research.12232.1.
11
12 553
13
14
15 554 21. Fierst JL, Murdock DA, Thanthiriwatte C, Willis JH, Phillips PC. Metagenome-Assembled
16 Draft Genome Sequence of a Novel Microbial *Stenotrophomonas maltophilia* Strain Isolated
17 from *Caenorhabditis remanei* Tissue. *Genome Announc.* 2017;5.
18 555
19
20 556 doi:10.1128/genomeA.01646-16.
21
22 557
23
24
25 558 22. Grenier E, Abadon M, Brunet F, Capy P, Abad P. A *mariner*-like transposable element in
26 the insect parasite nematode *Heterorhabditis bacteriophora*. *J Mol Evol.* 1999;48:328–36.
27 559
28
29
30
31 560 23. Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. BRAKER1: Unsupervised RNA-
32 Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics.*
33 561
34 2016;32:767–9. doi:10.1093/bioinformatics/btv661.
35
36 562
37
38
39 563 24. Boutet E, Lieberherr D, Tognolli M, Schneider M, Bansal P, Bridge AJ, et al.
40 UniProtKB/Swiss-Prot, the Manually Annotated Section of the UniProt KnowledgeBase:
41 How to Use the Entry View. *Methods Mol Biol.* 2016;1374:23–54. doi:10.1007/978-1-4939-
42 3167-5_2.
43 564
44 565
45
46 566
47
48
49 567 25. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND.
50 *Nat Methods.* 2015;12:59–60. doi:10.1038/nmeth.3176.
51 568
52
53
54 569 26. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome
55 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol.*
56 570
57 2015;16:157. doi:10.1186/s13059-015-0721-2.
58 571
59
60
61

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- 572 27. Coghlan A, Fiedler TJ, McKay SJ, Flicek P, Harris TW, Blasiar D, et al. nGASP--the
573 nematode genome annotation assessment project. *BMC Bioinformatics*. 2008;9:549.
574 doi:10.1186/1471-2105-9-549.
- 575 28. Koutsovoulos G, Kumar S, Laetsch DR, Stevens L, Daub J, Conlon C, et al. No evidence
576 for extensive horizontal gene transfer in the genome of the tardigrade *Hypsibius dujardini*.
577 *Proc Natl Acad Sci U S A*. 2016;113:5053–8. doi:10.1073/pnas.1600338113.
- 578 29. Yoshida Y, Koutsovoulos G, Laetsch DR, Stevens L, Kumar S, Horikawa DD, et al.
579 Comparative genomics of the tardigrades *Hypsibius dujardini* and *Ramazzottius varieornatus*.
580 *PLoS Biol*. 2017;15:e2002266. doi:10.1371/journal.pbio.2002266.
- 581 30. Baskaran P, Rödelsperger C, Prabh N, Serobyán V, Markov GV, Hirsekorn A, et al.
582 Ancient gene duplications have shaped developmental stage-specific expression in
583 *Pristionchus pacificus*. *BMC Evol Biol*. 2015;15:185. doi:10.1186/s12862-015-0466-2.
- 584 31. Prabh N, Rödelsperger C. Are orphan genes protein-coding, prediction artifacts, or non-
585 coding RNAs? *BMC Bioinformatics*. 2016;17:226. doi:10.1186/s12859-016-1102-x.
- 586 32. Kergunteuil A, Bakhtiari M, Formenti L, Xiao Z, Defosse E, Rasmann S. Biological
587 Control beneath the Feet: A Review of Crop Protection against Insect Root Herbivores.
588 *Insects*. 2016;7. doi:10.3390/insects7040070.
- 589 33. Ali F, Wharton DA. Cold tolerance abilities of two entomopathogenic nematodes,
590 *Steinernema feltiae* and *Heterorhabditis bacteriophora*. *Cryobiology*. 2013;66:24–9.
591 doi:10.1016/j.cryobiol.2012.10.004.
- 592 34. Howe KL, Bolt BJ, Shafie M, Kersey P, Berriman M. WormBase ParaSite - a
593 comprehensive resource for helminth genomics. *Mol Biochem Parasitol*. 2017;215:2–10.
594 doi:10.1016/j.molbiopara.2016.11.005.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
- 595 35. Sandhu SK, Jagdale GB, Hogenhout SA, Grewal PS. Comparative analysis of the
596 expressed genome of the infective juvenile entomopathogenic nematode, *Heterorhabditis*
597 *bacteriophora*. Mol Biochem Parasitol. 2006;145:239–44.
598 doi:10.1016/j.molbiopara.2006.01.002.
- 599 36. Hao Y-J, Montiel R, Lucena MA, Costa M, Simoes N. Genetic diversity and comparative
600 analysis of gene expression between *Heterorhabditis bacteriophora* Az29 and Az36 isolates:
601 uncovering candidate genes involved in insect pathogenicity. Exp Parasitol. 2012;130:116–25.
602 doi:10.1016/j.exppara.2011.12.001.
- 603 37. Boguski MS, Lowe TM, Tolstoshev CM. dbEST--database for “expressed sequence tags”.
604 Nat Genet. 1993;4:332–3. doi:10.1038/ng0893-332.
- 605 38. Boratyn GM, Camacho C, Cooper PS, Coulouris G, Fong A, Ma N, et al. BLAST: a more
606 efficient report with usability improvements. Nucleic Acids Res. 2013;41 Web Server
607 issue:W29-33. doi:10.1093/nar/gkt282.
- 608 39. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
609 transform. Bioinformatics. 2009;25:1754–60. doi:10.1093/bioinformatics/btp324.
- 610 40. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
611 Alignment/Map format and SAMtools. Bioinformatics. 2009;25:2078–9.
612 doi:10.1093/bioinformatics/btp352.
- 613 41. Tarailo-Graovac M, Chen N. Using RepeatMasker to identify repetitive elements in
614 genomic sequences. Curr Protoc Bioinformatics. 2009;Chapter 4:Unit 4.10.
615 doi:10.1002/0471250953.bi0410s25.
- 616 42. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
617 universal RNA-seq aligner. Bioinformatics. 2013;29:15–21.
618 doi:10.1093/bioinformatics/bts635.

- 619 43. Dobin A, Gingeras TR. Optimizing RNA-Seq Mapping with STAR. *Methods Mol Biol.*
2016;1415:245–62. doi:10.1007/978-1-4939-3572-7_13.
- 621 44. Gremme G, Steinbiss S, Kurtz S. GenomeTools: a comprehensive software library for
efficient processing of structured genome annotations. *IEEE/ACM Trans Comput Biol
Bioinform.* 2013;10:645–56. doi:10.1109/TCBB.2013.68.
- 624 45. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
features. *Bioinformatics.* 2010;26:841–2. doi:10.1093/bioinformatics/btq033.
- 626 46. Huang Y, Niu B, Gao Y, Fu L, Li W. CD-HIT Suite: a web server for clustering and
comparing biological sequences. *Bioinformatics.* 2010;26:680–2.
doi:10.1093/bioinformatics/btq003.
- 629 47. Laetsch DR, Blaxter ML. KinFin: Software for Taxon-Aware Analysis of Clustered
Protein Sequences. *G3 (Bethesda).* 2017. doi:10.1534/g3.117.300233.
- 631 48. Finn RD, Attwood TK, Babbitt PC, Bateman A, Bork P, Bridge AJ, et al. InterPro in
2017-beyond protein family and domain annotations. *Nucleic Acids Res.* 2017;45:D190–9.
doi:10.1093/nar/gkwl107.
- 634 49. Tange O. GNU Parallel - The Command-Line Power Tool. *login: The USENIX
Magazine.* 2011;36:42–7.
- 636 50. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
improvements in performance and usability. *Mol Biol Evol.* 2013;30:772–80.
doi:10.1093/molbev/mst010.
- 639 51. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with
thousands of taxa and mixed models. *Bioinformatics.* 2006;22:2688–90.
doi:10.1093/bioinformatics/btl446.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
- 642 52. Kocot KM, Citarella MR, Moroz LL, Halanych KM. PhyloTreePruner: A Phylogenetic
643 Tree-Based Approach for Selection of Orthologous Sequences for Phylogenomics. *Evol*
644 *Bioinform Online*. 2013;9:429–35. doi:10.4137/EBO.S12813.
- 645 53. Kück P, Meusemann K. FASconCAT: Convenient handling of data matrices. *Mol*
646 *Phylogenet Evol*. 2010;56:1115–8. doi:10.1016/j.ympev.2010.04.024.
- 647 54. Huson DH, Scornavacca C. Dendroscope 3: an interactive tool for rooted phylogenetic
648 trees and networks. *Syst Biol*. 2012;61:1061–7. doi:10.1093/sysbio/sys062.

649
650
651

652
653

654
655

656
657

658 **Figures and Legends**

659 **Figure 1. Taxon-annotated GC-coverage plot of the *H. bacteriophora* assembly.**

660 Bottom left panel: Each scaffold or contig is represented by a single filled circle. Each scaffold
661 is placed in the main panel based on its GC proportion (X axis) and coverage by reads from
662 the Illumina re-sequencing project (Y axis). The fill colour of the circle indicates the taxon of
663 the top BLASTn hit in the NCBI nt database for that scaffold. The colours are annotated in

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

664 the top right hand key, which indicates taxon assignment and (in brackets) the number of
665 contigs and scaffolds so assigned, their total span, and their N50 length. The circles are
666 scaled to scaffold length, as indicated in the key at the base of the main panel.

667 Right panel: Nucleotide span in kb at each coverage level.

668 Top panel: Nucleotide span in kb at each GC proportion.

669

670 **Figure 2. Comparisons of BRAKER1/soft-masked and original gene predictions**
671 **from *H. bacteriophora***

672 (A, B) Frequency histograms of intron count (A) and protein length (B) in BRAKER1/soft-
673 masked (blue) and published (yellow) protein coding gene predictions. Outlying proteins
674 longer than >2500 amino-acids (n=40) or genes containing >60 introns (n=20) are not
675 shown.

676 (C) Frequency histogram of the proportion of each BRAKER1 gene prediction overlapped
677 by a published gene prediction at the nucleotide level.

678 (D) Comparison of singleton, proteome-specific, and shared proteins in the published and
679 BRAKER1/soft-masked protein sets.

680 (E) Counts of non-canonical GC/AG introns in gene predictions from the published and
681 BRAKER1 *H. bacteriophora* gene sets, and the model nematode *Caenorhabditis elegans*
682 (WS258). Counts are of genes containing at least one non-canonical GC/AG intron with the
683 specified number of non-canonical introns.

684

685 **Figure 3. Maximum likelihood phylogeny of selected rhabditine (Clade V)**
686 **nematodes.**

687 A supermatrix of aligned amino acid sequences from orthologous loci from both *H.*
688 *bacteriophora* predictions and a set of 23 rhabditine (Clade V) nematodes (see
689 Supplementary Table 3) were aligned and analysed with RaxML using a PROTGAMMAGTR
690 amino-acid substitution model. *Pristionchus* spp. were designated as the outgroup. Bootstrap
691 support values (100 bootstraps performed) were 100 for all branches except one.

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710 **Tables**

711 **Table I. Contamination screening of the *H. bacteriophora* assembly**

Number of scaffolds	Sum of scaffold spans (bp)	Mean coverage *	Best matches in NCBI nt database	Assignment
12	99556	2.8	<i>Stenotrophomonas maltophilia</i> genome	bacterial culture contaminant **
4	4709	0.1	<i>Photorhabdus</i> sp. genomes	symbiont culture contaminant **
2	2144	756.0	poorly annotated mitochondrial matches	<i>H. bacteriophora</i> mitochondrial fragments
22	3051844	69.6	mariner transposons in Metazoa, especially Hymenoptera and Platyhelminthes	<i>H. bacteriophora</i> nuclear genome mariner transposon family (highest coverage 960-fold)
10	334100	76.6	low score match to several histone H3.3 across Metazoa	<i>H. bacteriophora</i> nuclear sequence
7	713932	56.5	chance nucleotide matches to conserved genes in other taxa	<i>H. bacteriophora</i> nuclear sequences

712

713 * The average read coverage of the whole assembly was 85.3.

714 ** These scaffolds were removed by the low-coverage filter.

715

716

717

718

719

720

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

721 **Table 2. Comparison of the published and BRAKER1/soft-masked protein coding**
 722 **gene predictions.**

<i>Prediction set</i>	<i>Published [13]</i>	<i>BRAKER1/soft-masked</i>
Number of protein coding genes predicted	20964	1,747
Mean protein length (amino acids)	218.8	344.5
Number of single exon genes	1728	2326
Mean number of exons per gene*	5.9	7.8
Proportion of non-canonical (GC-AG) introns	8.87%	0.79%
Percentage mapping to publicly available transcriptome reads		
<i>Sanger ESTs</i>	80.45%	84.26%
<i>Roche 454 reads</i>	37.18%	58.03%
BUSCO score for proteome		
<i>Complete</i>	47.8%	94%
<i>Fragmented</i>	34.7%	4.3%
Number of proteins with no hits in Uniref90	8,962	2,889
Protein singletons in clustering	5442	1112
Conserved, single-copy orthologues**		
<i>Total</i>	2089	2330
<i>Missing</i>	377	141
<i>Expanded</i>	184	84

723
 724 * Number of exons: number of coding DNA sequence (CDS) entries per gene for
 725 BRAKER1 predictions. CDS features, not exons are outputted by AUGUSTUS in general
 726 feature format (GFF).

727 ** The list of strict one-to-one orthologues was augmented with protein clusters where
 728 75% of species had single copy representatives (“fuzzy-1-to-1” orthologues identified by
 729 KinFin).

730

731

732

733

734 **Supplementary Files**

735 The supplementary files for this work are described below. All Supplementary files are
736 available at <https://github.com/DRL/mclean2017>.

737 **Supplementary file 1: Scaffolds and contigs removed from the *Heterorhabditis***
738 ***bacteriophora* assembly because of low coverage in the new whole genome**
739 **sequencing dataset**

740 Text file.

741 **Supplementary File 2: BRAKER1/soft-masked annotations of *Heterorhabditis***
742 ***bacteriophora*.**

743 A zipped archive (14.1 Mb) of the BRAKER1/soft-masked annotations of *Heterorhabditis*
744 *bacteriophora*. The archive contains three text files: the GFF format file, the GTF format file
745 and the amino acid sequences of the protein predictions in FASTA format.

746 **Supplementary File 3: BRAKER1/hard-masked annotations of *Heterorhabditis***
747 ***bacteriophora*.**

748 A zipped archive (13.4 Mb) of the BRAKER1/hard-masked annotations of *Heterorhabditis*
749 *bacteriophora*. The archive contains three text files: the GFF format file, the GTF format file
750 and the amino acid sequences of the protein predictions in FASTA format.

751 **Supplementary File 4: Comparison of the BRAKER1/soft-masked and**
752 **BRAKER1/hard-masked gene predictions from *Heterorhabditis bacteriophora*.**

753 Tab-delimited text file.

754 **Supplementary File 5: OrthoFinder analyses of *Heterorhabditis bacteriophora***

755 **predicted proteomes.**

756 A zipped archive (20.3 Mb) of the OrthoFinder analyses of *Heterorhabditis bacteriophora*

757 predicted proteomes with 23 other nematode species. The archive contains the following

758 files:

759 SF5.1 A list of the proteomes included in the OrthoFinder analyses (text format file)

760 SF5.2 List of *Heterorhabditis bacteriophora* proteins of length <30 amino acids excluded from

761 the OrthoFinder analyses (text format file).

762 SF5.3 The OrthoFinder output files. A zipped archive of the three OrthoFinder clustering

763 result files (published *H. bacteriophora* + 23 species; BRAKER1/soft-masked + 23 species:

764 published + soft-masked + 23 species).

765 SF5.4 Table with count of orthogroups at each contribution ratio from the BRAKER1/soft-

766 masked and published proteomes after clustering with 23 other Clade V nematodes.

767 **Supplementary File 6: KinFin analyses of *Heterorhabditis bacteriophora* predicted**

768 **proteomes.**

769 A zipped archive (27.8 Mb) of the KinFin analyses from the OrthoFinder analyses of

770 *Heterorhabditis bacteriophora* predicted proteomes.

1
2 **772** **Supplementary File 7: Phylogenetic analyses of *Heterorhabditis bacteriophora***
3 **predicted proteomes.**

4
5
6 **773** A zipped archive (11.2 Mb) of the supermatrix alignment and the phylogenetic trees
7
8 **774** produced for the the analyses of the *Heterorhabditis bacteriophora* proteomes. The archive
9
10 **775** contains the following files:

11
12
13
14 **776** SF7.1 Alignments of orthogroups used to build the supermatrix (directory of aligned
15
16 **777** sequences in fasta format).

17
18
19 **778** SF7.2 Supermatrix of aligned sequences (FASTA .fas format file).

20
21
22
23 **779** SF7.3 Phylogenetic analysis output files (NEWICK format text file).

24
25
26 **780** **Supplementary File 8: Secretome analyses of *Heterorhabditis bacteriophora***
27 **predicted proteomes.**

28
29
30
31
32 **782** Secretome analyses of *Heterorhabditis bacteriophora* and other rhabditine nematodes. The
33
34 **783** zipped archive (8 kb) contains the following text format files.

35
36
37
38 **784** SF8.1 Secretome predictions from the published Bai *et al.* (2013) protein predictions.

39
40
41
42 **785** SF8.2 Secretome predictions from the BRAKER1/soft-masked predictions.

43
44
45 **786** **Supplementary File 9: BRAKER1 and JIGSAW annotation pipelines.**

46
47
48
49 **787** Figure illustrating the differences between the BRAKER1 and the Bai *et al* 2013 JIGSAW
50
51 **788** prediction methods used for *Heterorhabditis bacteriophora*. PDF file.

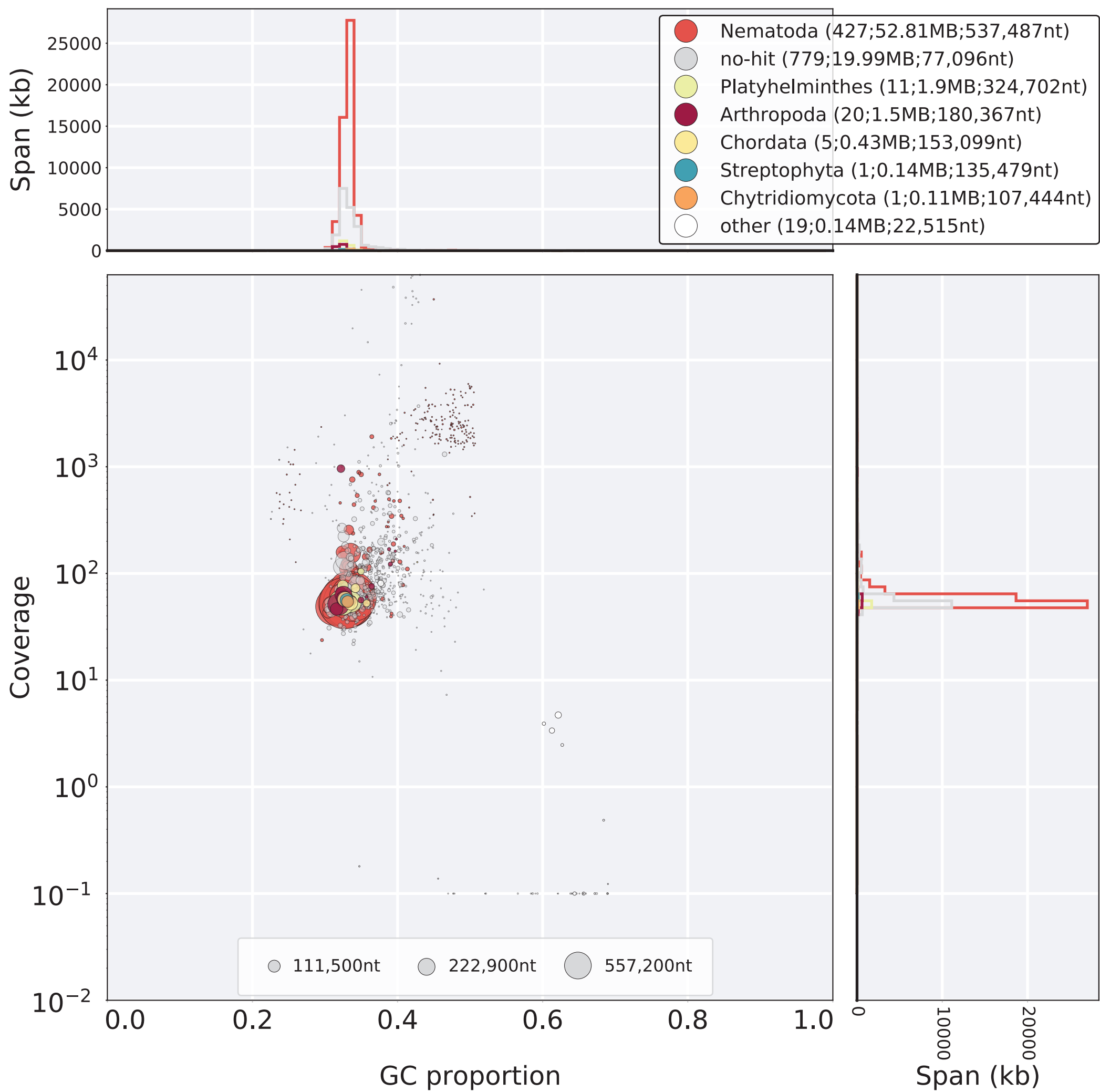


Figure 2

[Click here to download Figure Figure_2.pdf](#)