

1 Genome analysis of the fatal tapeworm *Sparganum proliferum* unravels the cryptic lifecycle and
2 mechanisms underlying the aberrant larval proliferation

3

4 Taisei Kikuchi^{1*}, Mehmet Dayi^{1,2}, Vicky L. Hunt³, Atsushi Toyoda⁴, Yasunobu Maeda¹, Yoko
5 Kondo⁵, Belkisyole Alarcon de Noya⁶, Oscar Noya⁶, Somei Kojima⁷, Toshiaki Kuramochi⁸,
6 Haruhiko Maruyama¹

7

8 ¹Faculty of Medicine, University of Miyazaki, Miyazaki, 889-1692 Japan

9 ²Forestry Vocational School, Duzce University, 81620, Duzce, Turkey

10 ³Department of Biology and Biochemistry, University of Bath, Bath, BA27AY UK

11 ⁴Comparative Genomics Laboratory, Department of Genomics and Evolutionary Biology, National
12 Institute of Genetics, Mishima, Shizuoka, 411-8540 Japan

13 ⁵Division of Medical Zoology, Department of Microbiology and Immunology, Faculty of Medicine,
14 Tottori University, Yonago 683-8503, Japan

15 ⁶Institute of Tropical Medicine, Central University of Venezuela, Venezuela, 2101 Maracay,
16 Caracas, Venezuela

17 ⁷Department of Clinical Laboratory Medicine, Chiba-Nishi General Hospital, Matsudo City, Chiba,
18 270-2251 Japan

19 ⁸Department of Zoology, National Museum of Nature and Science, 4-1-1 Amakubo, Tsukuba,
20 Ibaraki, 305-0005 Japan

21

22

23

24 *Corresponding author;

25 Taisei Kikuchi, Division of Parasitology, Faculty of Medicine, University of Miyazaki, Miyazaki 889-
26 1692, Japan, Tel: +81-985850990, Fax: +81-985843887, email: [taisei_kikuchi@med.miyazaki-](mailto:taisei_kikuchi@med.miyazaki-u.ac.jp)
27 [u.ac.jp](mailto:taisei_kikuchi@med.miyazaki-u.ac.jp)

28

29

30

31 **Abstract**

32 Background: The cryptic parasite *Sparganum proliferum* proliferates in humans and invades
33 tissues and organs. Only scattered cases have been reported, but *S. proliferum* infection is always
34 fatal. However, the *S. proliferum* phylogeny and lifecycle are still an enigma.

35 Results: To investigate the phylogenetic relationships between *S. proliferum* and other cestode
36 species, and to examine the underlying mechanisms of pathogenicity, we sequenced the entire *S.*
37 *proliferum* genome. Additionally, *S. proliferum* plerocercoid larvae transcriptome analyses were
38 performed to identify genes involved in asexual reproduction in the host. The genome sequences

39 confirmed that the *S. proliferum* genetic sequence is distinct from that of the closely related
40 *Spirometra erinaceieuropaei*. Moreover, nonordinal extracellular matrix coordination allows for
41 asexual reproduction in the host and loss of sexual maturity in *S. proliferum* is related to its fatal
42 pathogenicity in humans.

43 Conclusions: The high-quality reference genome sequences generated should prove valuable for
44 future studies of pseudophyllidean tapeworm biology and parasitism.

45

46 **Keywords:** pseudophyllidean tapeworm, gene family evolution, relaxed selection, extracellular
47 matrix coordination, asexual reproduction, oncogenes, homeobox, fibronectin, cadherin

48

49 **Background**

50 The cryptic parasite *Sparganum proliferum* was first identified in a 33-year-old woman in Tokyo in
51 1904. The patient's skin was infected with numerous residing cestode larva of unknown taxonomy.
52 Ijima *et al* [1] originally designated the parasite as *Plerocercoides prolifer*, and considered it a
53 pseudophyllidean tapeworm in the plerocercoid larval stage. In 1907, an extremely similar human
54 case was reported by Stiles in Florida, USA, and the responsible parasite was renamed *S.*
55 *proliferum* [2]. Clinical symptoms and post-mortem findings indicate that *S. proliferum* proliferates
56 in humans and invades various organs and tissues, including the skin, body walls, lungs,
57 abdominal viscera, lymph nodes, blood vessels, and the central nervous system, leading to
58 miserable disease prognosis [3, 4]. Not many cases have been reported to date but the infection
59 was fatal in all reported cases (reviewed in [5]).

60 There was a postulation about the origin of this parasite. Some parasitologists considered it to be
61 a new species of pseudophyllidean tapeworm, whereas others suspected that *S. proliferum* was a
62 virus-infected or aberrant form of *Spirometra erinaceieuropaei*, based on morphological similarities
63 [6, 7]. Recent DNA sequence analyses of mitochondrial NADH dehydrogenase subunit III,
64 mitochondrial tRNA, cytochrome oxidase subunit I, and nuclear succinate dehydrogenase iron-
65 sulfur protein subunit (sdhB) genes suggested that *S. proliferum* is a closely related but distinct
66 species of *S. erinaceieuropaei* [8, 9]. However, the adult stage of *S. proliferum* has not been
67 observed and the precise taxonomic relationships of *S. proliferum* with other worms remain unclear
68 because few genes have been sequenced.

69 In addition to taxonomic considerations, the pathogenicity of *S. proliferum* and its mechanisms of
70 proliferation and invasion in mammalian hosts are of considerable interest. In principle,
71 plerocercoids of pseudophyllidean tapeworms (spargana), including those of *S. erinaceieuropaei*
72 and other *Spirometra* species, do not proliferate asexually, but migrate through subcutaneous
73 connective tissues, causing only non-life threatening sparganosis (non-proliferative sparganosis).
74 Other organs, such as the lungs and liver or the central nervous system, may be niches for these
75 worms but are not commonly described. Symptoms of non-proliferative sparganosis are mainly
76 caused by the simple mass effect [5].

77 Asexual proliferation of larvae and the destruction of host tissues are characteristic of
78 cyclophyllidean tapeworms such as *Echinococcus*, which proliferates asexually by generating a
79 peculiar germinative layer in a hydatid cyst form [10]. In another cyclophyllidean tapeworm,
80 *Mesocestoides*, asexual multiplication is achieved by longitudinal fission [10, 11]. In contrast, the
81 pseudophyllidean *Sparganum* plerocercoid undergoes continuous branching and budding after
82 invading the human body by an unidentified route, and produces vast numbers of progeny
83 plerocercoids.

84 To clarify the phylogenetic relationship of the enigmatic parasite *S. proliferum* with other cestode
85 species and investigate the underlying pathogenic mechanisms, we sequenced its entire genome
86 as well as that of newly isolated *S. erinaceieuropaei*. We also performed transcriptome analyses
87 of *S. proliferum* plerocercoid larvae to identify genes that are involved in asexual reproduction in
88 the host. Those analyses revealed its phylogeny and gene evolution that contribute to the
89 proliferation and pathogenicity of *S. proliferum*.

90

91 **Results**

92 Genomic features of *S. proliferum* and *S. erinaceieuropaei*

93 We sequenced the *S. proliferum* genome using multiple insert-length sequence libraries
94 (Additional Table S1) and compiled a 653.4-Mb assembly of 7388 scaffolds with N50 of 1.2 Mb.
95 The *S. erinaceieuropaei* genome was assembled into 796 Mb comprising 5723 scaffolds with N50
96 of 821 kb. These assembly sizes were 51.9% and 63.2% of the previously published *S.*
97 *erinaceieuropaei* genome (UK isolate) [12]. CEGMA and BUSCO report the percentage of highly
98 conserved eukaryotic gene families that are present as full or partial genes in assemblies and
99 nearly 100% of core gene families are expected in most eukaryote genomes. BUSCO analyses
100 showed that 88.1% and 88.5% of core gene families were represented in *S. proliferum* and *S.*
101 *erinaceieuropaei* genomes, respectively, higher than or comparable to other previously published
102 tapeworm genomes (Table 1). CEGMA completeness values for *S. proliferum* and *S.*
103 *erinaceieuropaei* were slightly lower than those from BUSCO analyses. Low CEGMA
104 completeness was also seen in other pseudophyllidea tapeworm genomes, including *S.*
105 *erinaceieuropaei* UK isolate, *Diphyllobothrium latum*, and *Schistocephalus solidus* (Table 1). Low
106 CEGMA completeness values of these two genome assemblies, therefore, indicate
107 pseudophyllidean-specific loss or high divergence of the genes that are conserved in other
108 eukaryotic taxa. The average numbers of CEGs (hits for 248 single-copy eukaryotic core genes)
109 for *S. proliferum* and *S. erinaceieuropaei* were 1.2 and 1.3, respectively, indicating that the
110 assembly sizes roughly represent the haploid genome sizes of these tapeworms. However, in K-
111 mer analyses of Illumina short reads, we estimated haploid genome sizes of 582.9 and 530.1 Mb
112 for *S. proliferum* and *S. erinaceieuropaei*, respectively (Additional Fig S1a), indicating that the
113 assemblies contain heterozygous haplotypes and/or overestimated gap sizes. Ploidies were
114 inferred from heterozygous K-mer pairs and were diploid for both species (Figure S1b).

115 The genomes of *S. proliferum* and *S. erinaceieuropaei* are highly repetitive, with about 55.0%
116 repetition of the total genome length in both genomes (Additional Fig S2 and Table 2). Long
117 interspersed nuclear elements (LINEs) occupy 26.3% and 31.9% of the total genomes of *S.*
118 *proliferum* and *S. erinaceieuropaei*, respectively. These LINEs predominantly comprise the three
119 types (Penelope, RTE-BovB, and CR1), which are also abundant in other pseudophyllidea
120 genomes (Additional Fig S2).

121 A total of 25627 genes were predicted in *S. proliferum* assemblies, about 5000 fewer than for *S.*
122 *erinaceieuropaei* (30751), but more numerous than for other cestode genomes. In studies of the
123 *S. erinaceieuropaei* UK isolate {Bennett, 2014 #39}, the gene number (> 39000) was likely
124 overestimated due to fragmentation and redundancy in the assembly.

125

126 Phylogenetic placement of *S. proliferum*

127 Phylogenetic relationships of *S. proliferum* with other cestode species were inferred from 205
128 single-copy orthologues (Figure 1). A clear separation was identified between pseudophyllidea and
129 cyclophilidea clades. In the pseudophyllidea clade, *S. proliferum* occupied the basal position of the
130 *Spirometra* cluster, in which two *S. erinaceieuropaei* isolates (Japan and UK isolates) were placed
131 beside each other.

132 Phylogenetic tree topology based on mitogenomes of the 14 cestodes and all available
133 mitogenome data of *Spirometra* in the GenBank, was similar to that of the nuclear genome
134 (Additional Fig S3). Yet in contrast with the nuclear genome tree, the *S. erinaceieuropaei* UK
135 isolate was located in a basal position of the *Spirometra* cluster, placing *S. proliferum* in the middle
136 of *Spirometra* species, albeit with a long branch. These inconsistencies between nuclear and
137 mitogenome trees may reflect uncertainties of species classification in the genus *Spirometra* [13,
138 14]. Moreover, mitochondrial sequences can give poor inferences of species trees [15].
139 Cumulatively, these results suggest that *S. proliferum* has a close phylogenetic relationship with
140 *Spirometra* but is clearly distinguished by genomic features and gene contents.

141

142 Gene family evolution

143 Protein family (Pfam) analyses revealed highly similar protein domain distributions of *S. proliferum*
144 and *Spirometra* genomes ($r = 0.99$; Figure 2, Additional Table S2). Few domains differed
145 significantly in abundance between the two species. Among these, the *S. proliferum* genome was
146 underrepresented in zinc-finger families (Zf-C2H2, Zf-C2H2_4, Zf-C2H2_6, Zf-C2H2_jaz and Zf-
147 met), reverse transcriptase (RVT_1), exo/endonuclease/phosphatase, galactosyltransferase, and
148 alpha/beta hydrolase (abhydrolase_6). Overrepresented Pfam domains in *S. proliferum* included
149 a distinct type of zinc-finger domain (zf-3CxxC), fibronectin type III (fn3), trypsin, RNA polymerase
150 III RPC4 (RNA_pol_Rpc4), and an ADP-specific phosphofructokinase/glucokinase conserved
151 region (ADP_PFK_GK).

152 We performed gene family analysis using OrthoFinder with the predicted proteomes of *S.*
153 *proliferum*, *S. erinaceieuropaei*, and other selected cestode genomes. A total of 234522 proteins
154 from 14 cestode species were placed into 39174 gene families (Figure 1). The *S. proliferum*
155 proteome (25627 proteins) was encoded by 9136 gene families, among which 7364 were shared
156 by all 14 cestodes and 2550 proteins were specific to the species or singleton. The *S.*
157 *erinaceieuropaei* proteome (30751 proteins) was clustered into 9008 gene families, 3806 of which
158 were species specific or singletons. Only four gene families were specific to both *Spirometra* and
159 *Sparganum*.

160 We used computational analysis of gene family evolution (CAFE) to estimate gene family
161 expansion and contraction, and identified gene families with significantly higher than expected
162 rates of gains and losses (Figure 3, Additional Table S3). Twenty-one gene families were
163 significantly expanded in the *S. proliferum* lineage, and these included annotations for fibronectin,
164 reverse transcriptase, zinc-finger C2H2 type, and core histone (Additional Table S4). Significantly
165 contracted gene families (43 families) had annotations relating to signal transduction proteins, such
166 as phosphatases and kinases, and ion channels and ABC transporters (Additional Table S5).
167 Fibronectin, reverse transcriptase, zinc-finger C2H2 type, and peptidases were present in
168 expanded and contracted families.

169 In the *S. erinaceieuropaei* lineage, 63 and 15 gene families were significantly expanded or
170 contracted (Additional Table S6 and S7), respectively. Among them, highly lineage specific
171 expansion was found for 7 families (i.e. 10 or more genes in *S. erinaceieuropaei*, whereas one or
172 no genes in *S. proliferum*. For example, the Orthogroup OG0000184 contains one *S. proliferum*
173 gene and 44 *S. erinaceieuropaei* genes, encoding biphenyl hydrolase-like protein (BPHL), which
174 harbors the Pfam domain abhydrolase_6 (Figure 2). Although the other gene families mostly
175 encode proteins of unknown function, they were likely expanded after speciation from *S. proliferum*
176 and *S. erinaceieuropaei* and may have specific roles in the *S. erinaceieuropaei* lifecycle or
177 parasitism.

178

179 Conserved developmental pathway genes

180 Homeobox transcription factors are involved in patterning of body plans in animals. The homeobox
181 gene numbers are much fewer in parasitic flatworms than in most other bilaterian invertebrates,
182 which have a conserved set of approximately 100 homeobox genes. Genome severance of four
183 cyclophyllid cestodes revealed that out of 96 homeobox gene families that are thought to have
184 existed at the origin of the bilateria, 24 are not present in cestodes [16]. The pseudophyllid
185 cestodes *S. proliferum* and *S. erinaceieuropaei* have similar homeobox class repertoires as those
186 in cyclophyllid cestodes, in which class ANTP was the most abundant, followed by classes PRD
187 and TALE; Table 3). The total numbers of homeobox domains identified in *S. proliferum* and *S.*
188 *erinaceieuropaei* are 64 and 71, respectively, and because these were fewer than in the
189 cyclophyllids *Echinococcus multilocularis* and *Taenia solium* (Table 3), they are the most reduced

190 of any studied bilaterian animal. The three homeobox families Pou/Pou6, ANTP/Bsx, and
191 ANTP/Meox were not present in *S. proliferum* and *S. erinaceieuropaei*, whereas the homeobox
192 family ANTP/Ro was found in *S. proliferum* and *S. erinaceieuropaei* but not in *E. multilocularis* and
193 *T. solium* (Additional Fig S4).

194 Comparisons between *S. proliferum* and *S. erinaceieuropaei* showed that the homeobox families
195 TALE/Pknox, ANTP/Hox1, ANTP/Msxlx, and POU/Pou-like are missing in *S. proliferum*, despite
196 being present in the other cestodes. In contrast, the homeobox families ANTP/Dbx and PRD/Alx
197 were found in *S. proliferum* but not in *S. erinaceieuropaei*.

198 Other conserved genes with roles in flatworm developmental pathways, such as Hedgehog and
199 Notch, were conserved in *S. proliferum* and *S. erinaceieuropaei*. But in the Wnt pathway, whose
200 complement is much smaller than the ancestral complement in tapeworms [16], two further genes
201 (Axin and LEF1/TCF) were missing in *S. proliferum* and *S. erinaceieuropaei* (Table S8).

202

203 Horizontally transferred genes

204 To determine whether the present genomes contained horizontally transferred genes (HTGs) from
205 other organisms, we used a genome-wide prediction method based on a lineage probability index
206 using the software Darkhorse2 identified 19 and 33 putative HTGs in *S. proliferum* and *S.*
207 *erinaceieuropaei*, respectively (Additional Table S9 and S10). For these transfers, all possible host
208 organisms were bacteria except for one *Spirometra* gene that has high similarity to a chlorella virus
209 gene. Orthologues of most *S. proliferum* putative HTGs were also detected as horizontally
210 transferred in *S. erinaceieuropaei*. Moreover, possible host bacteria, including *Marinifilum breve*,
211 *Aphanizomenon flos-aquae*, *Alcanivorax* sp., and *Vibrio* sp., were shared by the two cestode
212 species and were aquatic or marine bacteria, indicating that these genes were acquired by a
213 common ancestor of the two tapeworms which had aquatic phase in the life cycle.

214

215 Positive selection of the *S. proliferum* lineage

216 Positive selection is a mechanism by which new advantageous genetic variants sweep through a
217 population and drive adaptive evolution. To investigate the roles of positive selection in the
218 evolution of *S. proliferum*, we performed dN/dS branch-site model analyses with single-copy
219 orthologous genes from 12 tapeworms and identified a total of 35 positively selected genes in the
220 *S. proliferum* lineage (Additional Table S11). Evolutionary pressures were identified for some
221 genes that are essential to cellular processes, including transcription/RNA processing/translation
222 genes encoding DNA-directed RNA polymerase II subunit, polypyrimidine tract-binding protein,
223 adenylate kinase, ribosomal protein L21, snu13 NHP2-like protein, and eukaryotic translation
224 initiation factors. Other identified genes were related to transportation (dynein intermediate chain
225 2) and mitochondrial processes (Rieske). Genes involved in stress and immune responses, such
226 as DNAJ/Hsp40, HIKESHI protein, Toll-like receptor, and Ig_3/Ig, were also positively selected in
227 the *S. proliferum* lineage, along with the RAS oncogene *Rab-4A*.

228 Environmental change often eliminates or weakens selective pressures that were formerly
229 important for the maintenance of a particular trait [17]. We detected 9 genes that were subject to
230 these circumstances of “relaxed selection” in the *S. proliferum* lineage, relative to the other
231 tapeworm lineages (Additional Table 12). These genes encode proteins with putative roles in
232 developmental regulation and cell differentiation. In particular, the receptor roundabout (ROBO)
233 and secreted molecules of the SLIT family, together, play important roles in guiding axons and
234 proper morphogenesis [18]. The Rho GTPase-activating protein is also highly expressed in highly
235 differentiated tissues and affects cell differentiation by negatively regulating Rho-GTPase signaling
236 [19]. Delta-like protein (DLL) is an inhibitory ligand of the Notch receptor pathway and is expressed
237 during brain development [20]. Vascular endothelial growth factor receptor is also known to
238 regulate stem cell homeostasis and repopulation in planarian species [21]. Hence, these instances
239 of relaxed selection indicate that the worm has long since used certain developmental pathways.
240 We also identified two genes encoding cadherin (protocadherin) that were subject to relaxed
241 selection. Cadherein is a transmembrane protein that mediates cell–cell adhesion in animals and
242 those relaxed selections indicate diverging cell adhesion process in the worm.

243

244 Differential gene expression involved in asexual proliferation and parasitism

245 We maintained *S. proliferum* via serial infection of mice and found that some plerocercoid worms
246 exhibit a highly branching structure (medusa-head form; Figure 4a), which was observed
247 frequently in heavily infected mice. In contrast, in mice with low worm burdens, most worms had
248 unadorned non-branching morphology (wasabi-root form). Worms with the medusa-head form are
249 considered the main sources of new plerocercoid worms in the host, and their proliferation is highly
250 related to their pathogenicity. We, therefore, identified genes with expression levels that
251 distinguished medusa-head and wasabi-root forms.

252 RNAseq analysis revealed 357 differentially expressed genes (DE genes) between medusa-head
253 and wasabi-root forms (246 upregulated and 111 downregulated in medusa-head) (Figure 4b).
254 The upregulated set in medusa-head forms were dominated by genes encoding peptidases and
255 peptidase inhibitors, such as tolloid-like proteins (19 genes), chymotrypsin-like proteins (6 genes)
256 and CAP domain-containing proteins (12 genes) as well as transposon-related proteins such as
257 gag-pol polyproteins and reverse transcriptases (30 genes) (Additional Table S13). This set of DE
258 genes was enriched in the GO categories for metalloendopeptidase activity and proteolysis
259 (Additional Table S14). Downregulated genes also encoded a variety of peptidases and peptidase
260 inhibitors, including leucyl aminopeptidase (5 genes), chymotrypsin-like elastase (7 genes), and
261 kunitz bovine pancreatic trypsin inhibitor domain protein (3 genes), with high representation under
262 the GO terms metalloexopeptidase, aminopeptidase, and manganese ion binding (Additional
263 Table S14). Peptidases and peptidase inhibitors are secreted by many types of pathogens,
264 including bacteria, fungi, and parasites, and often play critical roles in survival and virulence [22-
265 24]. Other genes known to be involved in pathogenicity in other pathogens were also upregulated

266 in the medusa-head form, including genes encoding multidrug resistance-associated proteins [25]
267 and tetraspanins. The latter proteins have four transmembrane domains and not only play roles in
268 a various aspects of cell biology but also are used by several pathogens for infection and regulate
269 cancer progression [26].

270 Genes that are involved in cell-growth and cancer development were also upregulated in the
271 medusa-head form, including those encoding proteins from wnt (wnt-111 and wnt-5) and ras/rab
272 (ras-0b, ras-2 and Rasef) pathways, transcription factors/receptors (sox1a, fibroblast growth factor
273 receptor) and homeobox proteins (prospero, PAX, orthopedia ALX and ISL2).

274 It has been shown that expansions of gene families and changes in expression levels have been
275 associated with the evolution of parasitism in previous studies [27, 28]. An upregulation of genes
276 from expanded gene families was also found in *S. proliferum*. For instance, 15 genes were
277 identified as upregulated from an expanded gene family (OrthoGroup OG0000040). The
278 orthogroup OG0000044 includes genes encoding mastin precursors, and six of these were
279 upregulated and another six were downregulated in the medusa-head form (Additional Table S13).
280 Phylogenetic analyses of those gene families indicate that some of these orthogroups are
281 conserved across flatworms, while others are specific to the Pseudophyllidea clade of flatworms
282 (Additional Fig S5).

283 Among the present DE genes, 85 that were upregulated in medusa-head forms have no known
284 functions. These included 17, 10, 3, 2, and 2 genes from orthogroups OG0000083, OG0003096,
285 OG0010117, OG0011363, and OG0011373, respectively. These orthogroups were expanded in
286 the *S. proliferum* lineage and the DE genes had extremely high fold changes (Figure 4c). Because
287 their products predominantly harbour secretion signal peptides (Additional Table S13), they are
288 likely to be secreted by the parasite into the host and play important roles in parasitism, aberrant
289 larval proliferation in the host, and/or modulation of host immunity.

290

291

292 **Discussion**

293 *S. proliferum* is a cryptic parasite with fatal consequences, but its phylogeny and lifecycle are
294 poorly understood. In this study, we sequenced the *S. proliferum* genome and performed
295 comparative genomics with other tapeworm species, including the newly-sequenced *S.*
296 *erinaceieuropaei* genome. The *S. erinaceieuropaei* genome was sequenced previously [12], with
297 an estimated genome size of more than 1.2 Gb, but because the source material was from a biopsy
298 the assembled sequence was highly fragmented. Hence, the *S. erinaceieuropaei* genome
299 presented herein provides a more reliable estimate of the size and contents of this parasite
300 genome. The new genome assembly was about two thirds of the size of the previous assembly
301 but remains the largest genome among sequenced tapeworms. Compared to cyclophyllidean
302 tapeworms, including *Echinococcus* and *Taenia* spp., for which high-quality genome references
303 are available [16, 29, 30], genome information for pseudophyllidean tapeworms is limited [31]. The

304 genomes presented in this study could, therefore, serve as a powerful resource for more
305 comprehensive studies of tapeworm genomics and will facilitate the understanding of
306 pseudophyllidean tapeworm biology and parasitism.

307

308 There have been three big knowledge gaps for the present cryptic tapeworm: 1) its phylogenetic
309 relationship with *Spirometra* species, 2) its lifecycle including the definitive and intermediate hosts,
310 and 3) genetic and physiological differences with non-proliferating *Spirometra* species that enable
311 the worm to reproduce asexually in non-definitive hosts, such as humans and mice.

312 To determine phylogenetic relationships, we confirmed that the genetic sequence of *S. proliferum*
313 is distinct from that of *S. erinaceieuropaei*, despite the close relationship between these species.
314 Specifically, the *S. proliferum* genome is about 150-Mb smaller and contains 5000 fewer protein
315 coding genes than in *S. erinaceieuropaei*. Both genomes, nonetheless, showed diploidy. These
316 data suggest that *S. proliferum* is not an aberrant form of *Spirometra* worm by virus infection or by
317 small mutations [6, 7] and not a hybrid origin of multiple *Spirometra* species. In agreement, no
318 virus-like sequences were detected in *S. proliferum* DNA or RNA raw reads.

319 We were unable to identify definitive or intermediate hosts of *S. proliferum* in the current study.
320 Recent horizontal transfers of genes or mobile elements can indicate phylogenetic relationships,
321 because HGT events occur between closely associated organisms. Well-known examples include
322 HGT from Wolbachia symbionts to their host insect [32, 33] and transfer of BovB retrotransposons
323 between ruminants and snakes via parasitic ticks [34, 35]. We found that RTE/BovB repeats are
324 abundant in the *S. proliferum* genome, but were likely acquired by an ancestral pseudophyllidea,
325 as indicated by their abundance in *D. latum* and *S. solidus*. Moreover, our HGT screening analyses
326 indicate several genes that were likely acquired from bacteria but these HGTs likely have occurred
327 before specification of *S. proliferum* and *S. erinaceieuropaei*. The high-quality reference genomes
328 presented herein, however, provide valuable resources for further attempts to identify vestigial *S.*
329 *proliferum* sequences in other organisms or to perform analyses of protein–protein interactions
330 between hosts and parasites.

331 Loss of genes that are involved in the development of multicellular organisms and nervous systems,
332 including homeobox genes and genes for zinc-finger domain containing proteins, and relaxed
333 selection of some developmental genes (ROBO, Slit, RHOGAP, etc.) suggests that *S. proliferum*
334 has lost the ability to undergo proper development and complete the sexual lifecycle. Although the
335 precise functions of homeobox genes in tapeworms remain elusive, proteins of homeobox families
336 that are missing in *S. proliferum* (TALE/Pknox, ANTP/Hox1, ANTP/Msxlx and POU/Pou-like)
337 appear to have important roles in the development of embryos and adult body plans. For example,
338 Hox1 of the HOX gene family specifies regions of the body plan of embryos and the head–tail axis
339 of animals [36]. Products of the Pknox gene family, also known as the PREP gene family, are
340 implicated as cofactors of Hox proteins [37]. Msxlx homeobox gene was highly upregulated in the
341 ovaries and was continually expressed in fertilized ova in the uterus in *Hymenolepis microstoma*.

342 This gene was related to the female reproductive system in this tapeworm [38]. POU class genes
343 are present in all animals and are extensively in nervous system development and the regulation
344 of stem cell pluripotency in vertebrates [39]. Specific loss of Pou-like genes and relaxed selection
345 of Pou3 suggest that *S. proliferum* has low dependency on POU genes.

346 We contend that the loss of sexual maturity of this parasite is related to its fatal pathogenicity in
347 humans, because survival of the parasite is dependent on asexual reproductive traits of budding
348 and branching, which lead to 100% lethality in infected humans. Accordingly, we identified genes
349 that are upregulated in vigorously budding worms using transcriptome analyses and then selected
350 genes that are putatively important for asexual proliferation, such as a variety of peptidase genes
351 and oncogene-like genes. Among them, groups of secreted proteins with unknown functions were
352 of great interest. They were expanded in the *S. proliferum* genome and showed more than 10-fold
353 changes in expression levels. Recently, an *S. erinaceieuropaei* gene belonging to one of those
354 groups (orthogroup OG0000083) was cloned and named plerocercoid-immunosuppressive factor
355 (P-ISF) (Yoko Kondo, under review). P-ISF is a cysteine-rich glycoprotein abundant in plerocercoid
356 excretory/secretory products and likely involved in immunomodulation of its hosts by suppressing
357 osteoclastogenesis including the gene expression of TNF- α and IL-1 β , and nitric oxide production
358 in macrophages [40, 41]. Upregulation of P-ISF genes in *S. proliferum* proliferating worms is
359 therefore reasonable and the expansion of the gene family in *S. proliferum* indicates the
360 considerable contribution to the specific lifecycle. The other upregulated gene families of unknown
361 function are also expanded in *S. proliferum* suggesting possible important roles in the hosts,
362 therefore, future studies of these novel genes are required to fully understand the mechanism
363 underlying the *S. proliferum* parasitism.

364 Fibronectin is an extracellular matrix (ECM) glycoprotein that controls the deposition of other ECM
365 proteins, including collagens and latent TGF-beta binding protein [42]. During branching
366 morphogenesis, accumulations of fibronectin fibrils promote cleft formation by suppressing
367 cadherin localization, leading to loss of cell–cell adhesion [43]. The present observations of the *S.*
368 *proliferum* lineage show specific expansions of three gene families containing fibronectin type III
369 domains. *S. proliferum* also had fewer cadherin genes than *S. erinaceieuropaei* and three of them
370 are subject to relaxed selection in *S. proliferum*. These results collectively suggest nonordinal ECM
371 coordination in *S. proliferum*, allowing the formation of highly branching structures and enabling
372 asexual proliferation in the host.

373

374

375 **Methods**

376 **Biological materials**

377 *S. proliferum* strain Venezuela was used for the genome analyses. The parasite was originally
378 isolated from a Venezuelan patient in 1981 and has been maintained by serial passages using
379 BALB/c mice via intraperitoneal injections of the plerocercoids in National Science Museum as

380 described in Noya et al [44, 45]. *S. erinaceieuropaei* was isolated from a Japanese four-striped rat
381 snake (*Elaphe quadrivirgata*) collected in Yamaguchi prefecture, Japan in 2014.

382

383 DNA and RNA extraction and sequencing

384 *S. proliferum* worms were collected from the abdominal cavity of infected mice and washed
385 thoroughly with 1x PBS. Plerocercoids of *S. erinaceieuropaei* were isolated from the subcutaneous
386 tissues of the snake. Genomic DNA was extracted using Genomic-tip (Qiagen) following the
387 manufacturer's instructions.

388 Paired-end sequencing libraries (Additional Table 1) were prepared using the TruSeq DNA Sample
389 Prep kit (Illumina) according to the manufacturer's instructions. Multiple mate-paired libraries (3, 8,
390 12 and 16 kb) were also constructed using the Nextera Mate-Paired Library Construction kit
391 (Illumina). Libraries were sequenced on the Illumina HiSeq 2500 sequencer using the Illumina
392 TruSeq PE Cluster kit v3 and TruSeq SBS kit v3 (101, 150 or 250 cycles x 2) or the Illumina MiSeq
393 sequencer with the v3 kit (301 cycles x 2) (Additional Table S1). The raw sequence data were
394 analysed using the RTA 1.12.4.2 analysis pipeline and were used for genome assembly after
395 removal of adapter, low quality, and duplicate reads.

396 RNA was extracted from individual worms using TRI reagent according to the manufacturer's
397 instructions. Total RNA samples were qualified using Bioanalyzer 2100 (Agilent Technology, Inc.).
398 Only samples with an RNA integrity value (RIN) greater than 8.0 were used for library construction.
399 One hundred ng of total RNA was used to construct an Illumina sequencing library using the
400 TruSeq RNA-seq Sample Prep kit according to the manufacturer's recommended protocols
401 (Illumina, San Diego, USA). The libraries were sequenced for 101 or 151 bp paired-ends on an
402 Illumina HiSeq2500 sequencer using the standard protocol (Illumina).

403

404 K-mer Analysis

405 A k-mer count analysis was performed using K-mer Counter (KMC) [46], on the paired-end Illumina
406 data. Only the first read was used to avoid counting overlapping k-mers. Genome size and ploidy
407 estimations were performed using Genomescope [47] and Smudgeplot, respectively [48].

408

409 Genome assembly

410 Illumina reads from multiple paired-end and mate-pair libraries (Additional Table 1) were
411 assembled using the Platanus assembler [49] with the default parameter. Haplomerger2 [50] was
412 then used to remove remaining haplotypic sequences in the assembly and contigs were further
413 scaffolded using Illumina mate-pair reads using SSPACE [51]. CEGMA v2 [52] and BUSCO [53]
414 were used to assess the completeness of the assemblies.

415 Mitochondrial genomes (mitogenomes) were reconstructed from Illumina reads with MITObim
416 version 1.6 [54]. Mitochondrial fragments in the nuclear genome assembly were identified by
417 BLASTX using *S. mansoni* mitochondrial genes as queries and those fragments were extended

418 by iterative mappings of Illumina short reads using MITObim. Assembled mitogenomes were
419 annotated for protein-coding, tRNA and rRNA genes using the MITOS web server [55]. Assemblies
420 and annotations were manually curated using the Artemis genome annotation tool [56] with based
421 on evidence supports from sequence similarity to other published mitogenomes.

422

423 Repeat analysis

424 Repeats within the genome assemblies were identified using RepeatModeler (v1.0.4,
425 <http://www.repeatmasker.org/RepeatModeler.html>) and RepeatMasker (v.3.2.8,
426 <http://www.repeatmasker.org>) to calculate the distribution of each repeat and its abundance in the
427 genome.

428

429 Gene prediction and functional annotation

430 To predict protein-coding genes, Augustus (v. 3.0.1) [57] was trained for *S. proliferum* and *S.*
431 *erinaceieuropaei*, individually, based on a training set of 500 non-overlapping, manually curated
432 genes. To obtain high-confidence curated genes, a selection of gene models from gene predictions
433 based on Augustus *S. mansoni* parameters, were manually curated in Artemis using aligned
434 RNA-seq data and BLAST matches against the NCBI database. RNA-seq reads were mapped to
435 the genomes using Hisat2 (parameters: --rna-strandness RF --min-intronlen 20 --max-intronlen
436 10000) [58]. Based on the Hisat2 alignments, the bam2hints program (part of the Augustus
437 package) was used to create the intron hints, with minimum length set to 20 bp. Augustus were
438 run with trained parameters using all the hints for that species as input. Introns starting with 'AT'
439 and ending with 'AC' were allowed (--allow_hinted_splicesites=atac). A weight of 10^5 was given to
440 intron and exonpart hints from RNA-seq. If Augustus predicted multiple, alternatively spliced
441 transcripts for a gene, we only kept the transcript corresponding to the longest predicted protein
442 for further analyses.

443 Functional annotations were performed on the gene models based on multiple pieces of evidence
444 including BLASTP search against NCBI nr database and the latest version Pfam search (ver. 30.0)
445 with HMMER3 [59]. Gene ontology (GO) terms were assigned to genes using Blast2Go (v2) [60]
446 with BLAST search against NCBI nr database and the InterProScan results.

447

448 Species tree reconstruction

449 Amino acid sequences in each single-copy gene family were aligned using MAFFT version
450 v7.22152 [61], poorly aligned regions were trimmed using GBLOCKS v0.91b53 [62], and then the
451 trimmed alignments were concatenated. A maximum-likelihood phylogenetic tree was produced
452 based on the concatenated alignment using RAxML v8.2.754 [63] with 500 bootstrap replicates.
453 The best-fitting substitution model for each protein alignment was identified using the RAxML
454 option (-m PROTGAMMAAUTO). Mitochondrial genome phylogeny was also constructed by the
455 same method using 12 protein coding genes on mitogenomes.

456

457 Gene family analysis

458 To estimate branch or lineage specific gain and loss of orthologous gene families, OrthoFinder
459 [64] and CAFÉ (v3) [65] under parameters “-p 0.01, -r 1000” were used.

460

461 Screening for horizontally transferred genes

462 To screen potential horizontal gene transfers (HGTs) into the *S. proliferum* and *S. erinaceieuropaei*
463 lineages, we used DarkHorse v2, which detects phylogenetically atypical proteins based on
464 phylogenetic relatedness of blastp hits against a taxonomically diverse reference database using
465 a taxonomically-weighted distance algorithm [66]. Options (-n 1 -b 0.5 -f 0.1) were used in
466 DarkHorse HGT screening.

467

468 Positive Selection Scans (dNdS)

469 To analyse selection pressures in *S. proliferum* genes, the ETE3 Python package [67] for CODEML
470 [68] was employed to calculate the non-synonymous (dN) and synonymous (dS) substitutions
471 rates, and the ratio (dN/dS or ω). Nucleotide sequences of single copy orthologue genes from 12
472 cestode species (*S. proliferum*, *S. erinaceieuropaei*, *Diphyllobothrium latum*, *Schistocephalus*
473 *solidus*, *Hymenolepis diminuta*, *Hymenolepis nana*, *Hydatigera taeniaeformis*, *Taenia solium*,
474 *Taenia asiatica*, *Echinococcus multilocularis*, *Echinococcus granulosus*, *Mesocestoides corti*)
475 were aligned based on amino acid alignment using Pal2aln v14 [69] with the parameters (-
476 nomismatch and -nogap). dN/dS were estimated using branch-site models with *S. proliferum* as
477 the foreground and other branches in the tree as the background. The non-null model (bsA) were
478 compared with the null model (bsA1) for each tree using a likelihood ratio test (LRT), where log-
479 likelihood ratios were compared to a chi-square distribution with 1 degree of freedom. False
480 discovery rate (FDR) correction were performed over all the P-values and genes showing FDR
481 <0.05 were manually curated before obtaining final dN/dS values.

482 Test for relaxed selection was performed using the RELAX tool [70] with aforementioned single
483 copy orthologue gene sets. The relaxation parameter k was calculated for each branch and tested
484 by LRT with *S. proliferum* as foreground and the others as background.

485

486 RNAseq analysis

487 For gene expression analyses, *S. proliferum* plerocercoid worms were grouped into two types
488 based on the morphology and proliferation activity; worms vigorously branching to form structure
489 like “Medusa head” and worms under static form to form like “Wasabi root” (Figure 4a). Worms
490 were collected from infected mice on ~50 weeks post inoculation. RNA was extracted from the
491 individual worms and sequenced as described above. RNAseq reads were mapped to the *S.*
492 *proliferum* reference genomes (v2.2) using Hisat2 [58] (parameters: --rna-strandness RF --min-
493 intronlen 20 --max-intronlen 10000). Mapped read count of each gene was calculated using HTSeq

494 [71] with options (-s no, -a 10, -m union) and differential expression analyses were performed
495 using EdgeR v3.2.4 [72]. A transcript was identified as differentially expressed in a pairwise
496 comparison if the following criteria were met: false discovery rate (FDR) \leq 0.001 and fold change
497 \geq 2.0. FPKM values were calculated using Cufflinks packages v2.2.1 [73] and used to generate for
498 multidimensional scaling (MDS) plots and gene expression heatmaps.

499

500

501 References

- 502 1. Ijima I: *On a New Cestode Larva Parasitic in Man (Plerocercoides Prolifer)*. 1905.
- 503 2. Stiles CW: **The occurrence of a proliferating cestode larva (Sparganum proliferum) in**
504 **man in Florida.** *Hyg Lab Bull* 1908, **40**:7-18.
- 505 3. Meric R, Ilie MI, Hofman V, Rioux-Leclercq N, Michot L, Haffaf Y, Nelson AM, Neafie
506 RC, Hofman P: **Disseminated infection caused by Sparganum proliferum in an AIDS**
507 **patient.** *Histopathology* 2010, **56**:824-828.
- 508 4. Nakamura T, Hara M, Matsuoka M, Kawabata M, Tsuji M: **Human proliferative**
509 **sparganosis: a new Japanese case.** *American journal of clinical pathology* 1990,
510 **94**:224-228.
- 511 5. Kikuchi T, Maruyama H: **Human proliferative sparganosis update.** *Parasitology*
512 *International* 2019:102036.
- 513 6. Iwata S: **On the branched plerocercoid (Sparganum proliferum) from Japanese snake.**
514 *Prog Med Parasitol Jpn* 1972, **4**:587-590.
- 515 7. Mueller JF, Strano AJ: **Sparganum proliferum, a sparganum infected with a virus?** *The*
516 *Journal of parasitology* 1974:15-19.
- 517 8. Kokaze A, Miyadera H, Kita K, Machinami R, Noya O, de Noya BA, Okamoto M, Horii T,
518 Kojima S: **Phylogenetic identification of Sparganum proliferum as a**
519 **pseudophyllidean cestode.** *Parasitology International* 1997, **46**:271-279.
- 520 9. Miyadera H, Kokaze A, Kuramochi T, Kita K, Machinami R, Noya O, de Noya BA,
521 Okamoto M, Kojima S: **Phylogenetic identification of Sparganum proliferum as a**
522 **pseudophyllidean cestode by the sequence analyses on mitochondrial COI and**
523 **nuclear sdhB genes.** *Parasitology international* 2001, **50**:93-104.
- 524 10. Reuter M, Kreshchenko N: **Flatworm asexual multiplication implicates stem cells and**
525 **regeneration.** *Canadian Journal of Zoology* 2004, **82**:334-356.
- 526 11. Specht D, Vogt M: **ASEXUAL MULTIPLICATION OF MESOCESTOIDES TETRATHYRIDIA**
527 **IN LABORATORY ANIMALS.** *J Parasitol* 1965, **51**:268-272.
- 528 12. Bennett HM, Mok HP, Gkrania-Klotsas E, Tsai IJ, Stanley EJ, Antoun NM, Coghlan A,
529 Harsha B, Traini A, Ribeiro DM: **The genome of the sparganosis tapeworm Spirometra**
530 **erinaceiuropeaei isolated from the biopsy of a migrating brain lesion.** *Genome*
531 *biology* 2014, **15**:510.
- 532 13. Almeida GG, Coscarelli D, Melo MN, Melo AL, Pinto HA: **Molecular identification of**
533 **Spirometra spp. (Cestoda: Diphylobothriidae) in some wild animals from Brazil.**
534 *Parasitol Int* 2016, **65**:428-431.
- 535 14. Jeon HK, Park H, Lee D, Choe S, Kim KH, Sohn WM, Eom KS: **Genetic Identification of**
536 **Spirometra decipiens Plerocercoids in Terrestrial Snakes from Korea and China.**
537 *Korean J Parasitol* 2016, **54**:181-185.

- 538 15. Bernt M, Bleidorn C, Braband A, Dambach J, Donath A, Fritzsich G, Golombek A, Hadryš
539 H, Juhling F, Meusemann K, et al: **A comprehensive analysis of bilaterian
540 mitochondrial genomes and phylogeny.** *Mol Phylogenet Evol* 2013, **69**:352-364.
- 541 16. Tsai IJ, Zarowiecki M, Holroyd N, Garcarrubio A, Sanchez-Flores A, Brooks KL, Tracey
542 A, Bobes RJ, Fragoso G, Scitutto E, et al: **The genomes of four tapeworm species reveal
543 adaptations to parasitism.** *Nature* 2013, **496**:57-63.
- 544 17. Lahti DC, Johnson NA, Ajie BC, Otto SP, Hendry AP, Blumstein DT, Coss RG, Donohue
545 K, Foster SA: **Relaxed selection in the wild.** *Trends Ecol Evol* 2009, **24**:487-496.
- 546 18. Cebria F, Newmark PA: **Morphogenesis defects are associated with abnormal
547 nervous system regeneration following roboA RNAi in planarians.** *Development* 2007,
548 **134**:833-837.
- 549 19. Basseres DS, Tizzei EV, Duarte AA, Costa FF, Saad ST: **ARHGAP10, a novel human gene
550 coding for a potentially cytoskeletal Rho-GTPase activating protein.** *Biochem Biophys
551 Res Commun* 2002, **294**:579-585.
- 552 20. Wenemoser D, Lapan SW, Wilkinson AW, Bell GW, Reddien PW: **A molecular wound
553 response program associated with regeneration initiation in planarians.** *Genes Dev*
554 2012, **26**:988-1002.
- 555 21. Lei K, Thi-Kim Vu H, Mohan RD, McKinney SA, Seidel CW, Alexander R, Gotting K,
556 Workman JL, Sanchez Alvarado A: **Egf Signaling Directs Neoblast Repopulation by
557 Regulating Asymmetric Cell Division in Planarians.** *Dev Cell* 2016, **38**:413-429.
- 558 22. Klemba M, Goldberg DE: **Biological roles of proteases in parasitic protozoa.** *Annu Rev
559 Biochem* 2002, **71**:275-305.
- 560 23. Frees D, Brondsted L, Ingmer H: **Bacterial proteases and virulence.** *Subcell Biochem*
561 2013, **66**:161-192.
- 562 24. Yike I: **Fungal proteases and their pathophysiological effects.** *Mycopathologia* 2011,
563 **171**:299-323.
- 564 25. Coleman JJ, Mylonakis E: **Efflux in Fungi: La Pièce de Résistance.** *PLOS Pathogens* 2009,
565 **5**:e1000486.
- 566 26. Florin L, Lang T: **Tetraspanin Assemblies in Virus Infection.** *Front Immunol* 2018,
567 **9**:1140.
- 568 27. Hunt VL, Hino A, Yoshida A, Kikuchi T: **Comparative transcriptomics gives insights into
569 the evolution of parasitism in Strongyloides nematodes at the genus, subclade and
570 species level.** *Sci Rep* 2018, **8**:5192.
- 571 28. Hunt VL, Tsai IJ, Coghlan A, Reid AJ, Holroyd N, Foth BJ, Tracey A, Cotton JA, Stanley
572 EJ, Beasley H, et al: **The genomic basis of parasitism in the Strongyloides clade of
573 nematodes.** *Nat Genet* 2016, **48**:299-307.
- 574 29. Zheng H, Zhang W, Zhang L, Zhang Z, Li J, Lu G, Zhu Y, Wang Y, Huang Y, Liu J, et al: **The
575 genome of the hydatid tapeworm Echinococcus granulosus.** *Nat Genet* 2013,
576 **45**:1168-1175.
- 577 30. Li W, Liu B, Yang Y, Ren Y, Wang S, Liu C, Zhang N, Qu Z, Yang W, Zhang Y, et al: **The
578 genome of tapeworm Taenia multiceps sheds light on understanding parasitic
579 mechanism and control of coenurosis disease.** *DNA Res* 2018, **25**:499-510.
- 580 31. International Helminth Genomes C: **Comparative genomics of the major parasitic
581 worms.** *Nature genetics* 2019, **51**:163-174.
- 582 32. Kondo N, Nikoh N, Ijichi N, Shimada M, Fukatsu T: **Genome fragment of Wolbachia
583 endosymbiont transferred to X chromosome of host insect.** *Proc Natl Acad Sci U S A*
584 2002, **99**:14280-14285.

- 585 33. Aikawa T, Anbutsu H, Nikoh N, Kikuchi T, Shibata F, Fukatsu T: **Longicorn beetle that**
586 **vectors pinewood nematode carries many Wolbachia genes on an autosome.** *Proc*
587 *Biol Sci* 2009, **276**:3791-3798.
- 588 34. Walsh AM, Kortschak RD, Gardner MG, Bertozzi T, Adelson DL: **Widespread horizontal**
589 **transfer of retrotransposons.** *Proc Natl Acad Sci U S A* 2013, **110**:1012-1016.
- 590 35. Ivancevic AM, Kortschak RD, Bertozzi T, Adelson DL: **Horizontal transfer of BovB and**
591 **L1 retrotransposons in eukaryotes.** *Genome Biol* 2018, **19**:85.
- 592 36. Carroll SB: **Homeotic genes and the evolution of arthropods and chordates.** *Nature*
593 1995, **376**:479-485.
- 594 37. Moens CB, Selleri L: **Hox cofactors in vertebrate development.** *Dev Biol* 2006,
595 **291**:193-206.
- 596 38. Olson PD, Zarowiecki M, James K, Baillie A, Bartl G, Burchell P, Chellappoo A, Jarero F,
597 Tan LY, Holroyd N, Berriman M: **Genome-wide transcriptome profiling and spatial**
598 **expression analyses identify signals and switches of development in tapeworms.**
599 *Evodevo* 2018, **9**:21.
- 600 39. Gold DA, Gates RD, Jacobs DK: **The early expansion and evolutionary dynamics of**
601 **POU class genes.** *Mol Biol Evol* 2014, **31**:3136-3147.
- 602 40. Dirgahayu P, Fukumoto S, Tademoto S, Kina Y, Hirai K: **Excretory/secretory products**
603 **from plerocercoids of Spirometra erinaceieuropaei suppress interleukin-1beta gene**
604 **expression in murine macrophages.** *Int J Parasitol* 2004, **34**:577-584.
- 605 41. Kina Y, Fukumoto S, Miura K, Tademoto S, Nunomura K, Dirgahayu P, Hirai K: **A**
606 **glycoprotein from Spirometra erinaceieuropaei plerocercoids suppresses**
607 **osteoclastogenesis and proinflammatory cytokine gene expression.** *Int J Parasitol*
608 2005, **35**:1399-1406.
- 609 42. Sakai T, Larsen M, Yamada KM: **Fibronectin requirement in branching morphogenesis.**
610 *Nature* 2003, **423**:876-881.
- 611 43. Wang S, Sekiguchi R, Daley WP, Yamada KM: **Patterned cell and matrix dynamics in**
612 **branching morphogenesis.** *The Journal of Cell Biology* 2017, **216**:559-570.
- 613 44. Moulinier R, Martinez E, Torres J, Noya O, de Noya BA, Reyes O: **Human proliferative**
614 **sparganosis in Venezuela: report of a case.** *The American journal of tropical medicine*
615 *and hygiene* 1982, **31**:358-363.
- 616 45. Alarcon de Noya B, Torres JR, Noya O: **Maintenance of Sparganum proliferum in vitro**
617 **and in experimental animals.** *Int J Parasitol* 1992, **22**:835-838.
- 618 46. Kokot M, Dlugosz M, Deorowicz S: **KMC 3: counting and manipulating k-mer statistics.**
619 *Bioinformatics* 2017, **33**:2759-2761.
- 620 47. Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz
621 MC: **GenomeScope: fast reference-free genome profiling from short reads.**
622 *Bioinformatics* 2017, **33**:2202-2204.
- 623 48. Jaron KS, Bast J, Ranallo-Benavidez TR, Robinson-Rechavi M, Schwander T: **Genomic**
624 **features of asexual animals.** *BioRxiv* 2018:497495.
- 625 49. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada
626 M, Nagayasu E, Maruyama H: **Efficient de novo assembly of highly heterozygous**
627 **genomes from whole-genome shotgun short reads.** *Genome research* 2014, **24**:1384-
628 1395.
- 629 50. Huang S, Kang M, Xu A: **HaploMerger2: rebuilding both haploid sub-assemblies from**
630 **high-heterozygosity diploid genome assembly.** *Bioinformatics* 2017, **33**:2577-2579.

- 631 51. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W: **Scaffolding pre-assembled**
632 **contigs using SSPACE**. *Bioinformatics* 2010, **27**:578-579.
- 633 52. Parra G, Bradnam K, Korf I: **CEGMA: a pipeline to accurately annotate core genes in**
634 **eukaryotic genomes**. *Bioinformatics* 2007, **23**:1061-1067.
- 635 53. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO:**
636 **assessing genome assembly and annotation completeness with single-copy**
637 **orthologs**. *Bioinformatics* 2015, **31**:3210-3212.
- 638 54. Hahn C, Bachmann L, Chevreur B: **Reconstructing mitochondrial genomes directly**
639 **from genomic next-generation sequencing reads--a baiting and iterative mapping**
640 **approach**. *Nucleic Acids Res* 2013, **41**:e129.
- 641 55. Bernt M, Donath A, Juhling F, Externbrink F, Florentz C, Fritzsche G, Putz J, Middendorf
642 M, Stadler PF: **MITOS: improved de novo metazoan mitochondrial genome**
643 **annotation**. *Mol Phylogenet Evol* 2013, **69**:313-319.
- 644 56. Carver T, Harris SR, Berriman M, Parkhill J, McQuillan JA: **Artemis: an integrated**
645 **platform for visualization and analysis of high-throughput sequence-based**
646 **experimental data**. *Bioinformatics* 2012, **28**:464-469.
- 647 57. Stanke M, Waack S: **Gene prediction with a hidden Markov model and a new intron**
648 **submodel**. *Bioinformatics* 2003, **19**:ii215-ii225.
- 649 58. Kim D, Langmead B, Salzberg SL: **HISAT: a fast spliced aligner with low memory**
650 **requirements**. *Nat Methods* 2015, **12**:357-360.
- 651 59. El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, Qureshi M, Richardson
652 LJ, Salazar GA, Smart A: **The Pfam protein families database in 2019**. *Nucleic acids*
653 *research* 2018, **47**:D427-D432.
- 654 60. Gotz S, Garcia-Gomez JM, Terol J, Williams TD, Nagaraj SH, Nueda MJ, Robles M, Talon
655 M, Dopazo J, Conesa A: **High-throughput functional annotation and data mining with**
656 **the Blast2GO suite**. *Nucleic Acids Res* 2008, **36**:3420-3435.
- 657 61. Katoh K, Standley DM: **MAFFT multiple sequence alignment software version 7:**
658 **improvements in performance and usability**. *Molecular biology and evolution* 2013,
659 **30**:772-780.
- 660 62. Castresana J: **Selection of conserved blocks from multiple alignments for their use in**
661 **phylogenetic analysis**. *Mol Biol Evol* 2000, **17**:540-552.
- 662 63. Stamatakis A: **RAXML version 8: a tool for phylogenetic analysis and post-analysis of**
663 **large phylogenies**. *Bioinformatics* 2014, **30**:1312-1313.
- 664 64. Emms DM, Kelly S: **OrthoFinder: solving fundamental biases in whole genome**
665 **comparisons dramatically improves orthogroup inference accuracy**. *Genome biology*
666 2015, **16**:157.
- 667 65. De Bie T, Cristianini N, Demuth JP, Hahn MW: **CAFE: a computational tool for the**
668 **study of gene family evolution**. *Bioinformatics* 2006, **22**:1269-1271.
- 669 66. Podell S, Gaasterland T: **DarkHorse: a method for genome-wide prediction of**
670 **horizontal gene transfer**. *Genome Biol* 2007, **8**:R16.
- 671 67. Huerta-Cepas J, Serra F, Bork P: **ETE 3: reconstruction, analysis, and visualization of**
672 **phylogenomic data**. *Molecular biology and evolution* 2016, **33**:1635-1638.
- 673 68. Yang Z: **PAML: a program package for phylogenetic analysis by maximum likelihood**.
674 *Bioinformatics* 1997, **13**:555-556.
- 675 69. Suyama M, Torrents D, Bork P: **PAL2NAL: robust conversion of protein sequence**
676 **alignments into the corresponding codon alignments**. *Nucleic acids research* 2006,
677 **34**:W609-W612.

- 678 70. Wertheim JO, Murrell B, Smith MD, Kosakovsky Pond SL, Scheffler K: **RELAX: detecting**
679 **relaxed selection in a phylogenetic framework.** *Mol Biol Evol* 2015, **32**:820-832.
- 680 71. Anders S, Pyl PT, Huber W: **HTSeq—a Python framework to work with high-**
681 **throughput sequencing data.** *Bioinformatics* 2015, **31**:166-169.
- 682 72. Robinson MD, McCarthy DJ, Smyth GK: **edgeR: a Bioconductor package for**
683 **differential expression analysis of digital gene expression data.** *Bioinformatics* 2010,
684 **26**:139-140.
- 685 73. Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn
686 JL, Pachter L: **Differential gene and transcript expression analysis of RNA-seq**
687 **experiments with TopHat and Cufflinks.** *Nature protocols* 2012, **7**:562.
- 688

689 **Declarations**

690 **Ethics approval and consent to participate:** All animal experiments in this study were performed
691 under the applicable laws and guidelines for the care and use of laboratory animals, as specified
692 in the Fundamental Guidelines for Proper Conduct of Animal Experiment and Related Activities
693 in Academic Research Institutions under the jurisdiction of the Ministry of Education, Culture,
694 Sports, Science and Technology, Japan, 2006.

695 **Consent for publication:** Not applicable.

696 **Availability of data and materials:** All sequence data from the genome projects have been deposited
697 at DDBJ/ENA/GenBank under BioProject accession PRJEB35374 and PRJEB35375. All
698 relevant data are available from the authors.

699 **Competing interests:** The authors declare that they have no competing interests

700 **Funding:** This work was supported by Japan Society for the Promotion of Science (JSPS) KAKENHI
701 Grant Numbers 26460510 and 19H03212, AMED 18fk0108009h0003 and JST CREST Grant
702 Number JPMJCR18S7.

703 **Authors' Contributions:** T.Ki., T.Ku. and H.M. conceived the study. T.Ki. contributed to study design.
704 V.L.H., H.M. and T.Ki. wrote the manuscript with inputs from others. BAN, ON, SK prepared
705 biological samples. Ki and T.Ku. conducted experiments. V.L.H., M.D., Y.M., A.T. and T.Ki.
706 completed genome assembly and analysed genome data.

707 **Acknowledgements:** Genome data analyses were partly performed using the DDBJ supercomputer
708 system. We thank Ryusei Tanaka, Asuka Kounosu, Akemi Yoshida for assistance and comments.

709 **Author Information** Correspondence and requests for materials should be addressed to T.K.
710 (taisei_kikuchi@med.miyazaki-u.ac.jp).

711

712

713 Table 1. Statistics of genome assemblies

	<i>Sparganum proliferum</i> (v2.2)	<i>Spirometra erinaceieur opaei</i> (v2.0)	<i>Spirometra erinaceieur opaei</i> (UK) (WBPS13)	<i>Diphyllobothrium latum</i> (WBPS13)	<i>Schistocephalus solidus</i> (WBPS13)	<i>Hymenolepis microstoma</i> (WBPS13)	<i>Taenia solium</i> (WBPS13)	<i>Echinococcus multilocularis</i> (WBPS13)
Assembly size (Mb)	653.4	796.0	1258.7	531.4	539.4	182.1	122.4	114.5
Num. scaffolds	7,388	5,723	482,608	140,336	56,778	3,643	11,237	1,288
Average (kb)	88.4	139.0	2.6	3.8	10.0	50.0	11.2	889.3
Largest scaff (kb)	8,099	5,490	90	80	595	2,234	740	15,981
N50 (kb)	1,242	821	5	7	32	767	68	5,229
N90 (kb)	110	167	1	2	5	41	5	213
Gaps (kb)	51,020	77,788	128,163	38,407	22,091	259	164	336
CEGMA completeness complete/partial (%)	61.7/81.5	58.5/80.2	29.4/45.9	49.6/65.3	76.6/87.9	91.9/92.7	87.1/90.7	93.2/93.2
Average CEG number complete/partial	1.1/1.2	1.1/1.3	1.8/2.2	1.5/1.6	1.2/1.3	1.1/1.1	1.2/1.2	1.1/1.1
BUSCO completeness (Metazoa dataset/Eukaryota dataset)	72.0/88.1	71.9/88.5	33.6/37.3	38.1/53.8	70.3/86.2	78.6/90.4	72.6/85.5	72.2/88.1
Num. coding genes	25,627	30,751	39,557	19,966	20,228	12,373	12,481	10,273
Coding gene size (median; aa)	665.0	627.0	200.0	216.0	455.0	709.0	609.0	596.0

714

715

716

717 Table 2. Statistics of repeats in genomes

	<i>Sparganum proliferum</i> (v2.2)		<i>Spirometra erinaceieuropaei</i> (v2)	
	num element	% in bp	num element	% in bp
SINEs:	49435	1.59	45184	1.09
LINEs:	390951	26.32	519275	31.90
LINE/Penelope	139623	7.95	214116	10.41
LINE/RTE-BovB	162469	10.24	188656	11.18
LINE/CR1	75037	7.59	101503	9.37
LTR element:1	18276	1.79	25374	1.88
LTR/Gypsy	16179	1.56	24544	1.81
DNA element:	22386	1.38	54802	2.48
DNA/CMC-EnSpm	5795	0.35	16161	0.69
DNA/TcMar-Tc1	8162	0.59	7416	0.53
Small RNA:	2906	0.15	2955	0.08
Satellites:	10962	0.39	5823	0.16
Simple repeat:	79986	1.21	68909	0.76
Low complexity:	3799	0.04	5690	0.06
Unclassified:	378820	20.68	425608	15.92
TOTAL	1004498 (55.01%)		1185136 (55.14%)	

718

719

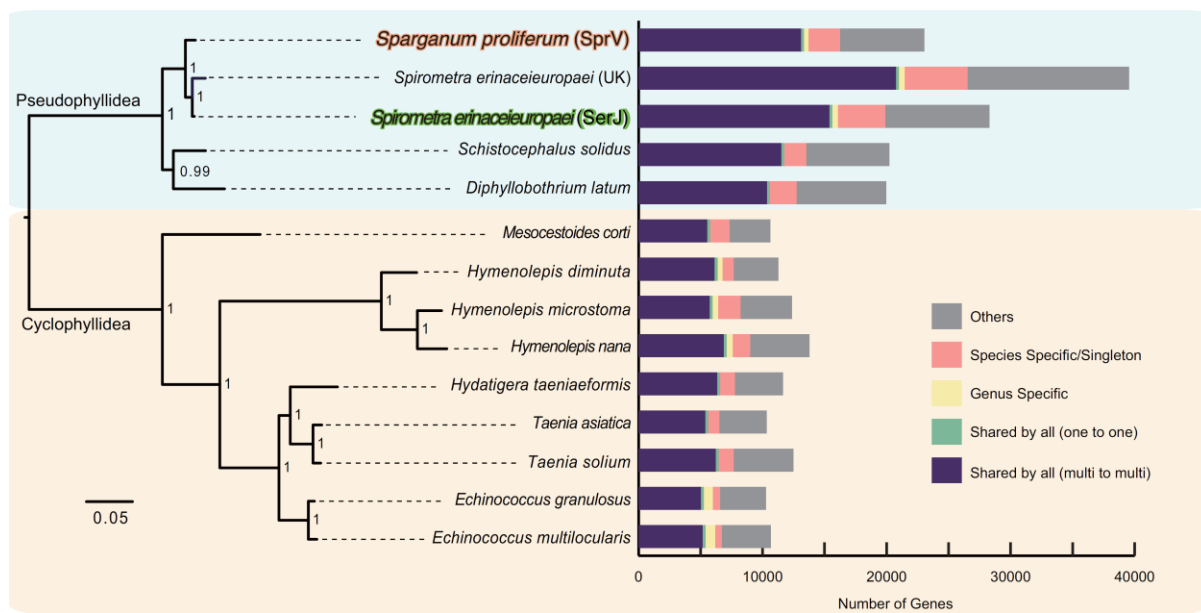
720

721 Table 3. Homeobox complement in *S. proliferum* and *S. erinaceiropaei* compared with other tapeworms and bilaterians

homeobox class	<i>Sparganum proliferum</i> (v2.2)	<i>Spirometra erinaceiropaei</i> (v2.0)	<i>Taenia solium</i> (WBPS13)	<i>Echinococcus multilocularis</i> (WBPS13)	<i>Branchiostoma floridae</i>
ANTP	25	30	36	25	58
PRD	10	8	11	15(18)	21
CUT	3	4	3	3	4
SINE	3	4	2	3	3
TALE	8	10	11	12	10
CERS	2	1	2	2	1
POU	3	4	4	5	6
LIM	6	6	7	7	8
ZF	4	4	3(2)	3(2)	4
Total	64	71	79	75	115

722

723

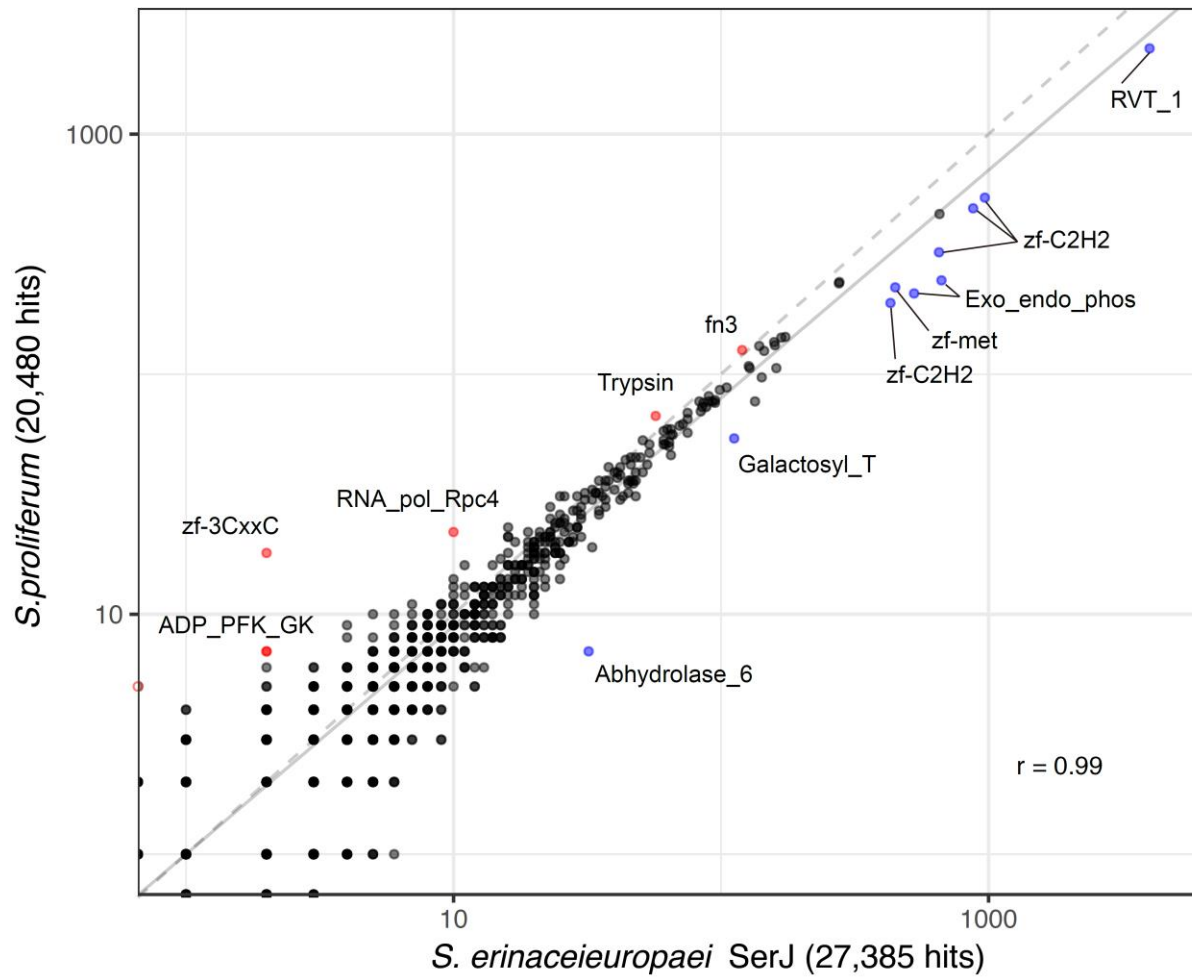


724

725

726 Figure 1. Phylogeny and gene contents; genes are categorized in a stack bar, and the length of
727 stack bar is proportional to number of genes.

728



729

730

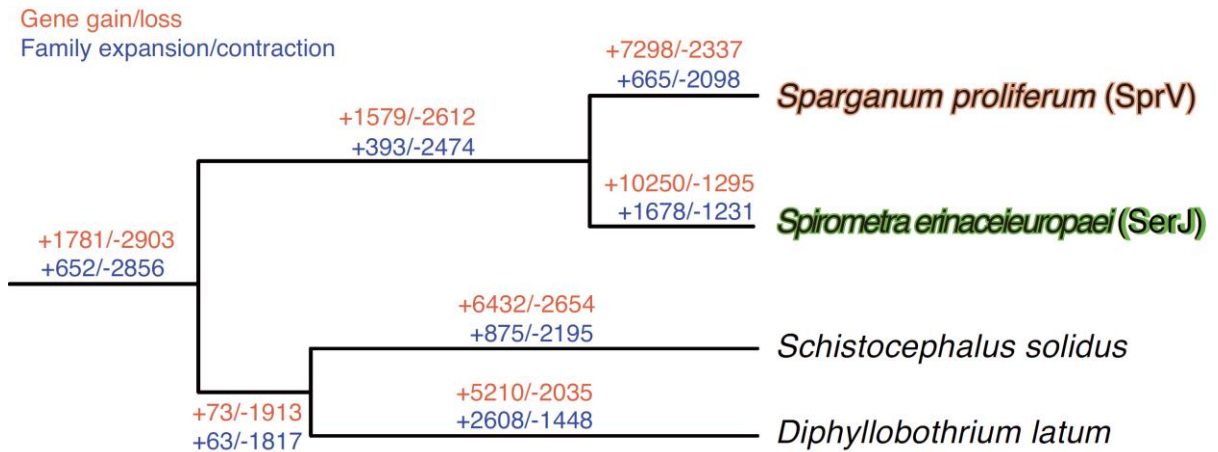
731 Figure 2. A scatterplot showing the abundance of Pfam domains in *S. proliferum* and *S.*
732 *erinaceiuropeaei* genomes; Pfam domains that are more enriched in *S. proliferum* than in *S.*
733 *erinaceiuropeaei* are highlighted in red. Those enriched in *S. erinaceiuropeaei* relative to *S.*
734 *proliferum* are highlighted in blue.

735

736

737

738



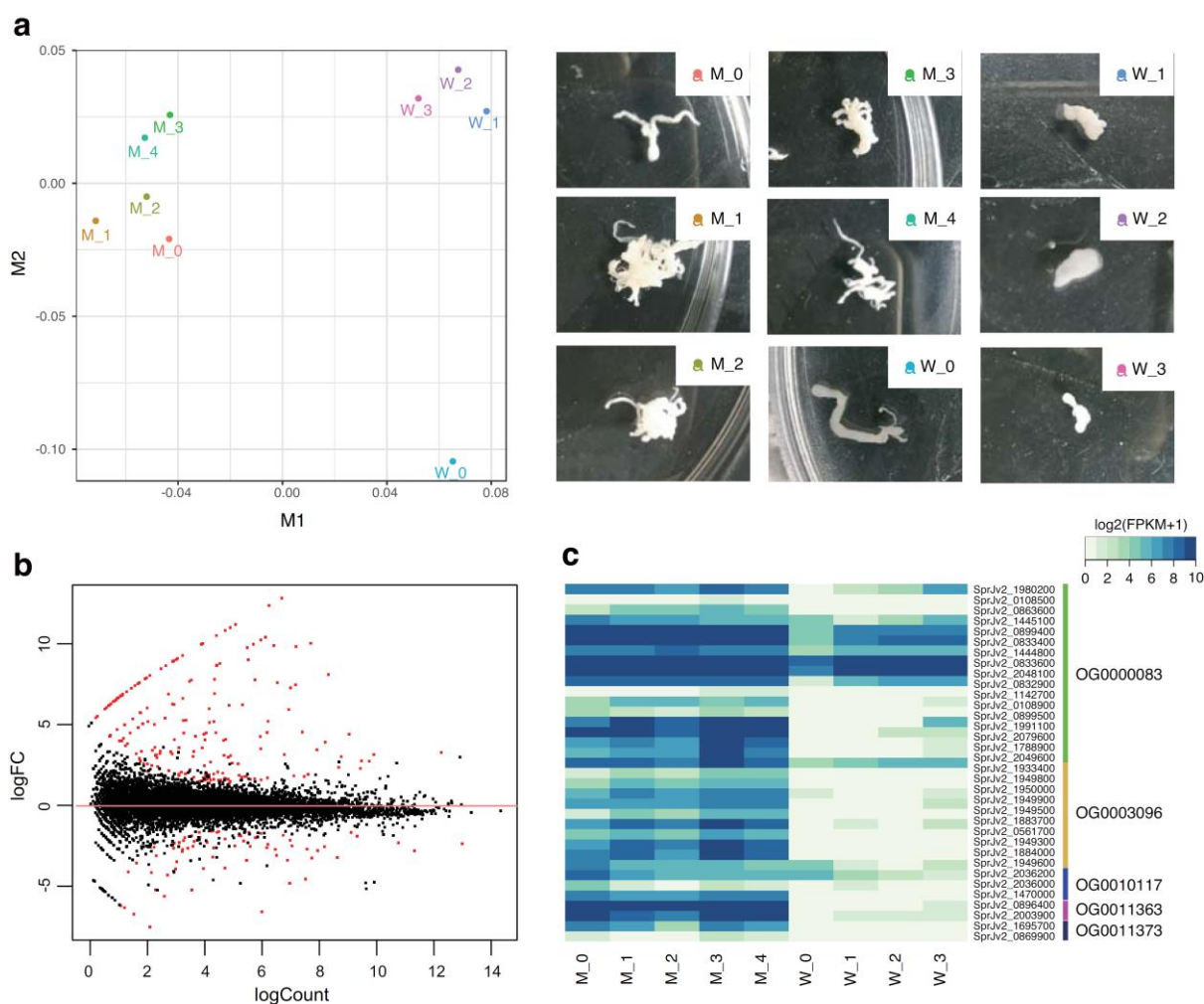
739

740

741 Figure 3. Gene family evolution of selected cestode species was inferred using computational
742 analysis of gene family evolution (CAFE). Numbers on each branch (or lineage) indicate specific
743 gains/losses of that branch (or lineage).

744

745



746

747 Figure 4. Comparison of gene expression in highly branching worms (medusa-head form) relative
 748 to static worms (wasabi-root form) of *S. proliferum*: A) multidimensional scaling (MDS) analyses of
 749 RNA-seq samples clearly separate the two forms by dimension 1. Pictures of used samples are
 750 shown on the right. B) Bland-Altman (MA) plot of the two-form comparison; dots represent
 751 transcripts and log₂ fold changes (medusa-head/wasabi-root) plotted against average abundance
 752 in counts per million. Red dots indicate differentially expressed transcripts with false discovery
 753 rates (FDR) of < 0.05 and fold changes of > 2. C) Heatmap of gene families encoding novel
 754 secreted proteins; the heat map shows log₂ fragments per kilobase per million reads mapped
 755 (FPKM) values for 5 gene families.

756