

1 **The compact genome of *Giardia muris* reveals important**
2 **steps in the evolution of intestinal protozoan parasites**

3

4 **Feifei Xu¹, Alejandro Jiménez-González¹, Elin Einarsson^{1,&}, Ásgeir Ástvaldsson¹,**
5 **Dimitra Peirasmaki¹, Lars Eckmann², Jan O. Andersson¹, Staffan G. Svärd¹, Jon**
6 **Jerlström-Hultqvist^{3,#}**

7 ¹Department of Cell and Molecular Biology, BMC, Box 596, Uppsala University,
8 SE-751 24 Uppsala, Sweden

9 ²Department of Medicine, University of California, San Diego, La Jolla, California,
10 USA

11 ³Department of Biochemistry and Molecular Biology, Dalhousie University, Halifax,
12 Canada

13 [&]Current address: Department of Biochemistry, University of Cambridge, Cambridge,
14 England

15

16 # Corresponding author: Dr. Jon Jerlström-Hultqvist

17 Department of Biochemistry and Molecular Biology

18 Sir Charles Tupper Medical Building

19 5850 College Street

20 Halifax NS Canada

21 B3H 4R2

22

23

24 Tel.: +1 902-494-2881; Fax: +1 902-494-1355

25 *E-mail adress:* jon.jerlstromhultqvist@dal.ca

26

27 **Running title:** Streamlined *Giardia* genome

28 **Key words:** parasite, diplomonad, *Giardia*, streamlined, intestinal colonization,

29 evolutionary biology, horizontal gene transfer, antigenic variation

30

31

32 **Abstract**

33 Diplomonad parasites of the genus *Giardia* have adapted to colonizing different hosts,
34 most notably the intestinal tract of mammals. The human-pathogenic *Giardia* species,
35 *Giardia intestinalis*, has been extensively studied at the genome and gene expression
36 level, but no such information is available for other *Giardia* species. Comparative data
37 would be particularly valuable for *Giardia muris*, which colonizes mice and is commonly
38 used as a prototypic *in vivo* model for investigating host responses to intestinal parasitic
39 infection. Here we report the draft-genome of *G. muris*. We discovered a highly
40 streamlined genome, amongst the most densely encoded ever described for a nuclear
41 eukaryotic genome. *G. muris* and *G. intestinalis* share many known or predicted virulence
42 factors, including cysteine proteases and a large repertoire of cysteine-rich surface
43 proteins involved in antigenic variation. Different to *G. intestinalis*, *G. muris* maintains
44 tandem arrays of pseudogenized surface antigens at the telomeres, whereas intact surface
45 antigens are present centrally in the chromosomes. The two classes of surface antigens
46 engage in genetic exchange. Reconstruction of metabolic pathways from the *G. muris*
47 genome suggest significant metabolic differences to *G. intestinalis*. Additionally, *G.*
48 *muris* encodes proteins that might be used to modulate the prokaryotic microbiota. The
49 responsible genes have been introduced in the *Giardia* genus via lateral gene transfer
50 from prokaryotic sources. Our findings point to important evolutionary steps in the
51 *Giardia* genus as it adapted to different hosts and it provides a powerful foundation for
52 mechanistic exploration of host-pathogen interaction in the *G. muris* – mouse
53 pathosystem.

54 **Importance**

55 The *Giardia* genus comprises eukaryotic single-celled parasites that infect many
56 animals. The *Giardia intestinalis* species complex, which can colonize and cause

57 diarrheal disease in humans and different animal hosts has been extensively explored
58 at the genomic and cell biologic levels. Other *Giardia* species, such as the mouse
59 parasite *Giardia muris*, have remained uncharacterized at the genomic level,
60 hampering our understanding of *in vivo* host-pathogen interactions and the impact of
61 host dependence on the evolution of the *Giardia* genus. We discovered that the *G.*
62 *muris* genome encodes many of the same virulence factors as *G. intestinalis*. The *G.*
63 *muris* genome has undergone genome contraction, potentially in response to a more
64 defined infective niche in the murine host. We describe differences in metabolic and
65 microbiome modulatory gene repertoire, mediated mainly by lateral gene transfer, that
66 could be important for understanding infective success across the *Giardia* genus. Our
67 findings provide new insights for the use of *G. muris* as a powerful model for
68 exploring host-pathogen interactions in giardiasis.
69

70 **Background**

71 Many eukaryotes have evolved from free-living to parasitic lifestyles over
72 evolutionary time, yet parasitism has developed independently in different taxonomic
73 groups and is therefore characterized by many unique features^{1,2}. Comparative
74 genomics provides an opportunity to investigate the factors of parasitism such as loss
75 of morphological, metabolic and genomic complexity, and consequently reduced
76 evolutionary potential for a free-living lifestyle². It can also identify the drivers and
77 consequences of a parasitic lifestyle and generate new testable ecological and
78 evolutionary hypotheses³.

79 *Giardia* is a protozoan parasite that non-invasively colonizes the intestinal
80 tract of many vertebrates. The human pathogen, *Giardia intestinalis*, is estimated to
81 cause 300 million cases of giardiasis in the world each year, being a major cause of
82 diarrheal disease⁴. Giardiasis is also a problem in domestic animals, and the zoonotic
83 potential of *Giardia* has been highlighted in recent years⁵. *In vitro* models of the
84 interaction of *G. intestinalis* with human cells have helped to unravel clues to how
85 *Giardia* causes disease⁶⁻⁸, such as the importance of, the adhesive disc for
86 attachment⁹, flagella for motility^{10,11}, secreted cysteine proteases for interference with
87 host defenses¹²⁻¹⁶, interactions with the intestinal microbiota^{6,17}, differentiation into
88 cysts for transmission^{4,18} and interference with nitric oxide (NO) production^{19,20}.
89 Despite this progress it remains uncertain whether these *in vitro* models are
90 representative of the natural infection, particularly because animal models of *G.*
91 *intestinalis* infection have significant limitations. For example, infection of mice, the
92 most commonly used laboratory animals, with human *G. intestinalis* isolates is
93 unreliable and requires manipulations such as antibiotic conditioning⁶. *Giardia muris*,
94 one of six recognized species of *Giardia*²¹, has been used as a mouse model since the

1960s for exploring the pathogenesis and immunological responses of the mammalian host to infection²². The availability of knock-out mice and other host-related resources makes *G. muris* a powerful model to investigate host-pathogen interactions⁶. The life cycle and infective process of *G. muris* is closely related to infection by *G. intestinalis*²³. Major findings in *Giardia* biology such as flagellar and disc function, cellular differentiation^{22,24,25}, and immunity^{23,26–29} have been pioneered with *G. muris*, and later been shown to be transferable to human *G. intestinalis* infections^{29,30}. Unfortunately, research on *G. muris* has been hampered by the lack of genome information and gene expression data⁵.

Here we describe the draft genome of *G. muris*, representing the first genome of any *Giardia* outside of the *G. intestinalis* species complex. We performed comparative genomics with free-living (*Kipferlia bialata*³¹) and parasitic (*G. intestinalis*^{32–34} and *S. salmonicida*³⁵) relatives to *G. muris* to determine how *G. muris* may have evolved into an intestinal pathogen of rodents.

110 **Results**

111 **Genome assembly**

We extracted DNA from freshly excysted *G. muris* cysts purified from the feces of infected mice (Fig. S1A) and assembled a high-quality draft genome using sequences obtained by PacBio and Illumina technologies. In addition, we generated RNA-Seq data for gene prediction and gene expression analyses with total RNA extracted from cysts, recently excysted cells (excystozoites), and trophozoites isolated from the small intestine of infected mice (Fig. S1A).

The *G. muris* draft genome consists of 59 contigs spanning 9.8 Mbp, which is notably smaller than the *G. intestinalis* WB genome (12.6 Mbp, Table 1). Most of the

120 genome (9.0 Mbp, 92%) is found on five contigs (>1 Mbp). Of the remaining short
121 contigs (<30 kbp), ten are terminated in telomeric repeats (TAGGG), suggesting they
122 are the terminal points of five chromosomes. The karyotype of *G. muris* was
123 previously shown to consist of four separable chromosomes³⁶. We hypothesize that
124 our five major contigs represent a total of five chromosomes in *G. muris*, two of
125 which are so close in size (1.290 and 1.297 Mbp) that they were not readily resolved
126 using pulsed-field gels, and are named accordingly from 1 to 5 from largest to
127 smallest in size (Fig. 1). 42 of the 44 small contigs contain ribosomal DNA (rDNA)
128 clusters that encode 28S, 18S, and 5.8S rRNAs. In fact, rDNA clusters make up 2.0%
129 of the total genome, and account for 91.6% of the identified repeats (Supplementary
130 Methods). Half of the contigs terminated by telomeric repeats have adjacent rDNA
131 clusters (Fig. S1B), suggesting that multiple of the *G. muris* chromosomes³⁶, like
132 those in *G. intestinalis*³⁷, have long repeats of rDNAs close to the telomeres. In
133 contrast to *G. intestinalis* chromosomes³⁷, no retrotransposon sequences were found in
134 the telomeric regions and overall very few retrotransposon sequences were detected in
135 the *G. muris* genome.

136 Allelic sequence heterozygosity (ASH) in the assembly was estimated to be
137 0.016% (Table 1), equivalent to the low level found in the *G. intestinalis* WB genome
138 (0.026%, Table 1). Distribution of ASH along chromosomes showed only weak
139 clustering in certain areas, particularly at the ends of chromosomes (Fig. 1).

140

141 **Genome streamlining and synteny**

142 Gene prediction and manually curated annotation identified 4653 protein coding
143 genes in *G. muris* (Table 1). This makes 84.5% of the genome coding, counting also the
144 tRNAs and rRNAs (Table 1). Thus, the *G. muris* genome is an example of a very compact

145 eukaryotic genome. Consistent with that, the average intergenic size is 264 bp (Table 1),
146 with a prominent skew towards shorter intergenic regions for a high proportion of genes
147 (median size at 37 bp, Fig. 2C). The compactness of the genome is also illustrated by
148 multiple instances of overlapping genes, with 441 genes (9.5% of all genes) showing an
149 average overlapping size of 21 bp (spanning 1-327 bp) with neighboring genes.

150 The *G. muris* and the new improved *G. intestinalis* WB genome
151 (AACB03000000, F. Xu and S.G. Svärd) do not maintain clear chromosomal synteny
152 even though both are assembled as five near-complete chromosomes (Fig. S2).
153 However, local synteny (Fig. 2AB) was obtained among 3,043 one-to-one orthologs
154 (Table S1A) (an average amino acid similarity of 44.7%) shared by the two genomes.
155 Comparing local synteny, it becomes obvious that *G. muris* keeps shorter orthologous
156 gene and intergenic region sizes (Fig. 2BCD).

157 **Gene regulation**

158 We could not identify any universal, conserved promoter motifs shared by all
159 *G. muris* genes except for an enrichment of A residues around the start codon (Fig.
160 S3A), which resembles observations in *G. intestinalis*³⁸. The streamlining of the *G.*
161 *muris* genome was also apparent at the 3' end of genes where the putative
162 polyadenylation signal, which is similar to the one described in *G. intestinalis*³⁸, is
163 overlapping with the stop codon for most genes (Fig. S3B). Genes up-regulated early
164 during encystation in *G. intestinalis* have specific promoter elements^{39,40}. Most of
165 these genes were also identified in *G. muris*, with one notable exception of cyst-wall
166 protein 3 (Fig. 3C). Encystation-related genes in *G. muris* share promoter motifs (Fig.
167 3A), similar to the Myb binding sites found in *G. intestinalis*⁴⁰, suggesting a similar
168 type of regulation. We also noted that the encystation-related genes are among the

169 most highly expressed genes in *G. muris* trophozoites *in vivo* (Fig. S3C, Table S1B),
170 similar to *G. intestinalis* infection in mice^{41,42}.

171 Very few genes in *G. intestinalis* contain introns, with only eight known *cis*-
172 spliced and four *trans*-spliced genes (five *trans*-introns)⁴³⁻⁴⁵. Similarly, only three *cis*-
173 and no *trans*-introns were identified in the parasitic diplomonad *S. salmonicida*³⁵,
174 whereas the free-living fornicate *K. bialata* has on average seven *cis*-introns per protein
175 encoding gene³¹. *G. muris* maintains homologs to the eight *cis*-spliced *G. intestinalis*
176 genes, but has only three retained introns (Fig. S4A). All four *trans*-spliced genes in *G.*
177 *intestinalis* have homologs in *G. muris* with conserved splicing motifs (Fig. S4). Mining
178 genes with similar motifs did not reveal additional intron-containing genes in *G. muris*.
179 Similar to *G. intestinalis*⁴³ all the *trans*-spliced genes in *G. muris* preserve a similar
180 cleavage motif TCCTTTACTCAA (Fig. S4C) as the RNA processing sequence motif⁴³.
181 Thus, we observe a reduction of introns in *G. muris*, and *cis*-introns seem to be easier to
182 lose than *trans*-introns.

183 **VSPs and antigenic variation in *G. muris***

184 Variant specific-surface proteins (VSPs) in *G. intestinalis* are characterized as
185 cysteine-rich proteins with frequent CXXC motifs and a conserved C-terminal
186 transmembrane (TM) domain followed by a cytoplasmic pentapeptide (CRGKA, Fig.
187 S5A). We identified 265 VSP homologs in *G. muris*. Their C-terminal pentapeptide
188 (GCRGK, Fig. S5A, Table S1C) differed slightly from that in *G. intestinalis*. However,
189 the cysteine and arginine residues in the pentapeptide, which are known to be post
190 translationally modified in *G. intestinalis*, are conserved^{46,47}. In addition, the preceding 24
191 aa of the *G. muris* VSP TM domain show conservation to the TM domain of *G.*
192 *intestinalis* VSPs (Fig. S5A). Most *G. muris* VSPs contain the conserved GGCY motif
193 present in most *G. intestinalis* VSPs (Table S1C). Since *bona fide* VSPs need signal

194 peptides (SPs) at the N-terminus to guide VSPs to the parasite surface, we divided this
195 group into two subgroups; proteins with predicted SPs are called VSPs, whereas VSP
196 proteins without SPs are referred to as pseudogenized VSPs (ψ VSPs). The 26 complete
197 VSP genes (16 unique at 98% identity to each other) are mostly located chromosome-
198 centrally (Fig. 1, Table 3). Seven pairs of VSP genes were identified with identical
199 sequences arranged either as head-to-head (2 pairs) or tail-to-tail (5 pairs) (Table 4).
200 Sequences in between the different VSP pairs resemble NimA (never in mitosis gene a)-
201 related kinase (NEKs), ankyrin repeat proteins (ARPs) and zinc-finger domains. There
202 are also 4 copies of identical VSPs clustering close to the 3' end of chromosome 2 (Fig.
203 2) with tandem repeats and sequences resembling NEKs and zinc-finger domains
204 interspersed (Table 4).

205 In contrast to complete VSPs, most ψ VSPs (183 / 239) are found in linear
206 arrays (n=17) in *G. muris*, herein defined as having >3 ψ VSPs genes (Fig. S5BC).
207 Strikingly, nine out of the ten ends of the main contigs have a ψ VSP array at or close
208 to the ends of the chromosomes (telomere-adjacent) containing a total of 131 genes
209 (Table S1C). The only main contig that ends without an array has a cluster of two
210 ψ VSPs close to the chromosome terminus. We found a single ψ VSP array consisting
211 of 12 genes in a chromosome position that was non-telomere adjacent (on
212 chromosome 5) (Fig. 1). The ψ VSP arrays vary in copy numbers (5-23 genes) and the
213 terminal part of the ψ VSP array is always arranged with the tail-end towards the
214 chromosome terminus. The tandem arrangement of the gene arrays suggested that
215 they were generated by gene duplication. Two of the terminal arrays are scrambled
216 and have a shift in the ψ VSP array directionality at the site of an intact VSP (Fig.
217 S5B).

218 We constructed a phylogeny of all the VSPs and ψ VSPs in *G. muris* to
219 investigate their evolutionary dynamics. The phylogeny revealed relaxed clustering of
220 ψ VSP genes originating in each linear array (Fig. S5C). However, the internal linear
221 array on chromosome 5 represents a noteworthy exception showcasing a very recent
222 gene duplication event. Interestingly, the great majority of full-length VSPs (23/28)
223 are clustered in the phylogeny, including the VSPs in the scrambled ψ VSP arrays,
224 despite these genes being distributed in physically separate chromosomal locations
225 across all five primary scaffolds (Fig. S5C). The few non-clustered VSP genes in the
226 phylogeny that are not part of pairs or clusters are found directly adjacent to ψ VSPs
227 genes. The relaxed clustering of VSPs and ψ VSPs suggested that these genes might
228 be undergoing periodical recombination or gene conversion.

229 The VSPs and ψ VSPs showed distinct expression patterns. Essentially all the
230 ψ VSP genes were non-transcribed in the three surveyed life-stages (Fig. S5D). VSPs, on
231 the other hand, showed on average higher expression with one or a few loci displaying
232 dominant expression in the different life-stages (Fig. S5D).

233 There is also another, less characterized VSP-related cysteine-rich protein family
234 in *G. intestinalis*, High Cysteine Membrane Proteins (HCMPs)⁴⁸, with 62 members³².
235 Many are highly up-regulated during interaction with intestinal epithelial cells⁴⁹. The
236 HCMPs have several CXXC and CXC motifs, one VSP-like transmembrane domain but
237 with longer C-terminals than in the VSPs⁴⁹. The 34 genes matching these criteria in the *G.*
238 *muris* genome were named HCMPs after the corresponding gene family in *G. intestinalis*.
239 They are found spread-out on the five chromosomes (Fig. 1).

240 **Multigene families in *G. muris***

241 The largest multigene families in *G. intestinalis* outside the VSPs and HCMPs are
242 the NEKs⁵⁰ and ankyrin repeat containing proteins (Protein 21.1)³². There are 230 NEKs

243 in *G. muris* (Fig. 1, Table 4), making up 71% of its kinome, slightly more than what was
244 found in *G. intestinalis*⁵⁰. We classified ankyrin repeat containing proteins further into
245 three groups. The ankyrin repeat protein-1 (ARP-1) with only ankyrin repeats, ARP-2
246 with ankyrin repeats plus zinc finger domains, and ARP-3 with both ankyrin repeats plus
247 domains other than zinc finger domains. The NEKs and the different classes of ARPs are
248 scattered throughout the chromosomes without obvious clustering (Fig. 1). A
249 phylogenetic analysis revealed that 79 of the NEKs are conserved as 1:1 orthologs
250 between *G. muris* and *G. intestinalis* (Fig. S6A). Each species has one massively
251 expanded cluster of NEKs with *G. muris* having the largest with 104 members and the
252 one in *G. intestinalis* having 79 members. ARPs show a similar evolutionary stability
253 with 132 conserved 1:1 orthologs between species (Fig. S6B). *G. muris* shows a major
254 species-specific expansion of 91 genes whereas the largest expanded clusters in *G.*
255 *intestinalis* amounts to two groups of 15 genes each. The partly shared domain-structure
256 of NEKs and ARPs prompted us to investigate their relationship by a network analysis
257 employing reciprocal blastp (1e-05 cutoff). The two groups are not recovered as clearly
258 separated clusters but form partially overlapping networks, indicating that there might be
259 recombination or gene conversion in between (Fig. S6C). The larger numbers of genes in
260 these multigene families in *G. muris* compared to *G. intestinalis* and their evolutionary
261 dynamics are intriguing given the otherwise streamlined features of the *G. muris* genome,
262 perhaps suggesting that they have unique roles in adaptation to their murine hosts.

263 **Virulence factors in *Giardia***

264 *G. intestinalis* is not known to possess classical virulence factors, such as
265 enterotoxins, but several genes are important for colonization of the host and thus for
266 pathogenesis. These include genes for motility¹⁰, the adhesive disc for attachment⁹,
267 secreted cysteine proteases that can degrade host defensive factors^{12,13}, and cysteine-rich

268 surface protein like the VSPs⁵¹ and the HCMPs⁴⁸ that undergo antigenic variation. The
269 cytoskeletal protein repertoire in *G. muris* is very similar to *G. intestinalis* apart from
270 several fragmented alpha-tubulins (3 complete genes with homologs in *G. intestinalis* and
271 9 incomplete gene fragments). The adhesive disc is a unique cytoskeletal structure of
272 *Giardia* parasites essential for attachment of the trophozoite in the small intestine, but is
273 missing in other fornicates⁹. The first detailed studies of the adhesive disc were performed
274 on *G. muris* trophozoites²⁴, but more recent work has mostly focused on *G. intestinalis*.
275 The vast majority (82 of 85) of *G. intestinalis* disc proteins⁹ were also identified in *G.*
276 *muris* (Table S1D); 12 were NEK kinases and 27 were ARP-1 proteins. Two of the three
277 *G. intestinalis* disc proteins not found in *G. muris* (Table S1D) localize to a structure in
278 the *G. intestinalis* disc on top of the ventral groove, but this structure is missing in the *G.*
279 *muris* disc²², suggesting functional disc differences. Many of the disc proteins that are
280 immunodominant during *G. intestinalis* infections (e.g. alpha-giardins, beta-giardin,
281 SALP-1, alpha- and beta-tubulin⁵² are highly expressed (here defined as >500 FPKM) in
282 *G. muris* trophozoites in the small intestine (Table S1B).

283 Proteases are important virulence factors in many pathogens and cysteine-protease
284 activities have been suggested to play a role in *Giardia* virulence^{4,13}. We identified 81
285 proteins classified as proteases in the *G. muris* genome, compared to 96 proteins
286 identified in an identical search in *G. intestinalis* WB (Table S1E). The largest family of
287 proteases in *G. muris*, with 15 members, are papain-like cysteine proteases (C1A family).
288 This protein family is also the largest protease group in *G. intestinalis* with 21
289 members^{53,54}. Several conserved groups of proteases were found to have been present in
290 the ancestor to *Giardia* and *Spironucleus*, although we also found evidence for lineage-
291 specific gene loss and expansion in *G. muris* (Fig. S7). The most highly expressed

292 cysteine protease of *G. muris in vivo* is the closest homolog to the highest expressed
293 protease in *G. intestinalis* WB⁵⁴.

294 **Metabolic pathways in *G. muris***

295 Our metabolic reconstruction identified 95 metabolic pathways in *G. muris*
296 compared to 98 pathways detailed in *G. intestinalis* WB; four pathways were unique in *G.*
297 *muris* and eight in *G. intestinalis* WB. Even though the overall metabolism is highly
298 similar between *G. muris* and *G. intestinalis* WB, the genes for several specific enzymes
299 and their putative reactions show distinct differences. Thus, 18 unique reactions (14
300 enzymes) were predicted in *G. muris* and 25 in *G. intestinalis* WB (10 enzymes) (Table
301 2). Several of these unique proteins showed moderate-high identity to prokaryotic
302 proteins (Table S1F). Five of these prokaryote-like genes are present in the genome of
303 assemblage B strain *G. intestinalis* GS.

304 The potential utilization of carbohydrate sources for glycolysis is different in *G.*
305 *muris* compared to *G. intestinalis*. Fructokinase and mannose 6-phosphate isomerase
306 enable *G. muris* to use fructose and mannose 6-phosphate unlike *G. intestinalis*. In both
307 cases, the enzymes were acquired from bacteria via lateral gene transfer (Table S1F,
308 Figure S8AB). Curiously, the *G. muris* gene for mannose 6-phosphate isomerase is found
309 next to a bacterial transcriptional regulator/sugar kinase gene. Phylogenetic analyses
310 show that *G. muris* mannose 6-phosphate isomerase and the transcriptional
311 regulator/sugar kinase gene clusters deep in the Bacteroidetes group. This gene
312 arrangement is observed in bacteria of the genus *Alistipes*, whose genomes harbor the
313 most similar homologs, supporting the notion that a single event of lateral gene transfer
314 best explains the origin of these genes in *G. muris* (Figure S8BC). Both genes have the A-
315 rich initiator that precedes the start codon in most *G. muris* genes and are expressed in *G.*
316 *muris* trophozoites (Table S1B).

317 The utilization of glycerol for ATP synthesis via glycerol kinase has been
318 suggested in *G. intestinalis* upon depletion of primary carbon sources⁵⁵. This enzyme is
319 found in both *Spironucleus* and *Trichomonas* but has been lost in *G. muris*. Another
320 notable metabolic difference to *G. intestinalis* WB is the lack of pyrophosphate-
321 dependent pyruvate phosphate dikinase (PPDK) that leaves a single, less energy efficient,
322 route from phosphoenol pyruvate to pyruvate via pyruvate kinase in *G. muris*.

323 *G. muris* is predicted to synthesize coenzyme A from pantothenate, employing a
324 bifunctional phosphopantothenoyl decarboxylase-phosphopantothenate synthase. The
325 same pathway is described in *S. salmonicida*³⁵, but the complete pathway is missing in *G.*
326 *intestinalis* (Table 2).

327 As in all other studied metamonads, *G. muris* encodes the arginine dihydrolase
328 (ADH) pathway that enables the use of arginine as an energy source and at the same time,
329 reduces the available free arginine in the environment, preventing nitric oxide (NO)
330 production in host cells⁵⁶. NO efficiently kills *G. intestinalis* trophozoites and the main
331 scavenger enzyme for NO in *G. intestinalis* is Flavohemoprotein⁵⁷, which is lacking in *G.*
332 *muris* (Table 2). Arginine is an important modulator of virulence in many infectious
333 organisms since it interferes with NO production⁵⁸. *G. muris* encodes arginases which
334 converts arginine directly to ornithine and urea. Arginases are present in *G. intestinalis*
335 GS and *T. vaginalis*, representing an ancestral acquisition in Metamonada followed by
336 subsequent losses in *Spironucleus* and *G. intestinalis* assemblage A and E (Table 2 and
337 Fig. S8D).

338 We noticed that several of the bacterial derived genes that are shared with *G.*
339 *intestinalis* GS are clustered together on the chromosomes in *G. muris* in highly dynamic
340 genomic regions. For example, arginase genes are found in a four-gene genomic region
341 that is present in two adjacent copies in chromosome 2, two on chromosome 1 and one on

342 chromosome 5 (dark grey filled arrows, Fig. 4A). Intact arginase genes are only found on
343 chromosome 2 adjacent to the genes encoding 2,5-diketo-D-gluconic acid reductase. All
344 the duplication events have occurred between ARPs (red arrows) and NEKs (dark red
345 arrows). The homologous regions on the three chromosomes likely originated via two
346 duplication events.

347 Nucleotide substitutions have accumulated since the duplication events and in-
348 frame stop codons (marked by red asterisk in Fig. 4A) have rendered three of the arginase
349 genes pseudogenized. The small ORFs sitting at the other side of arginase are hypothetical
350 proteins and have similar homologs in other parts of the genome, but the sequences of their
351 duplicated homologs in those regions are pseudogenized (dotted light grey block in Fig.
352 4A). Both genes have their closest relatives in bacteria, even though the genes can be
353 found in other eukaryotes (Fig. S13E). Phylogenetic analyses indicate the genes might
354 have been transferred from different bacterial donors multiple times into different
355 eukaryotic lineages (Fig. S13E).

356 A second, distantly related 2,5-diketo-D-gluconic acid reductase gene copy (Fig.
357 4B, Fig. S8F) in the genome is present in three different genomic locations together with
358 two other enzymes, carboxymuconolactone decarboxylase (CMD) and ketosteroid
359 isomerase-like protein. All three genes, constitute putative lateral gene transfers (Fig.
360 S8GH). This three-gene region is close to ψ VSPs (orange arrows) and ARPs.

361 362 **Interaction with the intestinal microbiota**

363 *Giardia* trophozoites colonize the intestinal lumen where they can potentially
364 interact with other intestinal microbes. Although little is functionally known about such
365 interactions and their consequences for parasite survival, four proteins encoded in the *G.*
366 *muris* genome could play a role in these interactions.

367 Bactericidal/permeability-increasing (BPI) proteins are innate immune defense
368 proteins that bind to lipopolysaccharide and display potent killing activity against gram-
369 negative bacteria by increasing membrane permeability. Beyond this basic function BPI
370 proteins might also act as effectors in controlling mutualistic symbioses⁵⁹. Homologs of
371 BPI proteins are found in *G. muris* and *G. intestinalis*⁶⁰, but it remains to be determined if
372 they have anti-microbial activity.

373 *G. muris* encodes tryptophanase, an enzyme that metabolizes tryptophan to
374 pyruvate with concomitant release of indole and ammonia. While pyruvate can be utilized
375 in energy metabolism, indole and its metabolites have been shown to affect gut
376 microbiota composition, possibly by interfering with quorum-sensing systems, and might
377 be able to influence host health⁶¹. Phylogenetic analysis of this protein showed that this
378 enzyme represents an ancestral acquisition in diplomonads with subsequent loss in *G.*
379 *intestinalis* (Fig. S8I).

380 Two more proteins with potential importance for microbiota interactions are
381 encoded in *G. muris*: Tae4 and quorum-quenching N-acyl-homoserine lactonase. The
382 Tae4 proteins are wide-spread amidases that were first described in association with the
383 T6SS system effector Tae4 in *Salmonella* Typhimurium⁶². The Tae4 proteins degrade
384 bacterial peptidoglycan by hydrolyzing the amide bond, γ -D-glutamyl-mDAP (DL-bond)
385 of Gram-negative bacteria⁶², and is required for interbacterial antagonism and successful
386 gut colonization by *S. Typhimurium*⁶³. Quorum-quenching N-acyl-homoserine lactonase
387 degrades N-acyl-homoserine lactone, a molecule used by both Gram-positive and Gram-
388 negative bacteria, for quorum sensing⁶⁴. Our phylogenetic analysis supports the lateral
389 acquisition of both genes (Fig. S8JK). While Tae4 has been a recent acquisition in *G.*
390 *muris*, quorum-quenching N-acyl-homoserine lactonase was present in the common

391 ancestor of *G. muris* and *G. intestinalis* and lost in *G. intestinalis* WB and P15 (Fig.
392 S8K).

393 **Discussion**

394 Our data shows that *G. muris* has an even more compact genome than *G. intestinalis*,
395 whose genome is already known to be highly streamlined³². Genome compaction via
396 reduction of mobile or repetitive elements have been seen in other eukaryotic
397 parasites^{1,65}. *G. muris* appears to fall into this category as it encodes no known classes
398 of mobile elements and repetitive elements are mostly confined to telomeric contexts.
399 The shortness of intergenic regions in *G. muris* ranks among the most extreme
400 recorded for any eukaryote, even shorter than Microsporidia which are known as the
401 most compact and reduced eukaryotic genomes⁶⁶. The global synteny map of *G. muris*
402 to *G. intestinalis* indicates many frequent small-scale genome rearrangements that
403 often favors a more efficient gene packing in *G. muris* thus allowing shorter
404 intergenic regions. This evidence of gene shuffling and the fact that there is very little
405 evidence of genome degradation would argue for optimization of growth as the
406 driving force of *G. muris* genome streamlining.

407 *G. muris* trophozoites have not been grown axenically *in vitro*, which has
408 hampered exploration of its genome, gene regulation and metabolism^{5,21}, and has
409 limited the use of *G. muris* as an *in vitro* model system for the human parasite *G.*
410 *intestinalis* and other intestinal protozoan parasites⁵. We identified several metabolic
411 differences between *G. muris* and *G. intestinalis* that might indicate avenues to
412 successful strain axenization. Most of these differences are represented by instances
413 of lateral gene transfer of metabolic genes or losses thereof in either *G. intestinalis* or
414 *G. muris*. *G. intestinalis* isolates are typically poor at infecting mice. Despite this, the
415 assemblage B isolate GS, which shares more metabolic enzymes with *G. muris* than

416 the assemblage A isolate WB, is better able to establish infection in mice than WB.
417 This suggests that the shared metabolic capacity of *G. muris* and *G. intestinalis* GS
418 enables survival in the murine intestinal tract. Additionally, *G. muris* might be able to
419 interact or interfere with intestinal Gram-positive and Gram-negative bacteria and this
420 could be a key to establish successful infections. *G. intestinalis*, which lacks some of
421 the putative microbiome modulators, such as Tryptophanase and Tae4, is dependent
422 on reduction of the small intestinal microbiota in order to efficiently infect mice⁶⁷. *G.*
423 *muris* is cleared from the murine host by secretory IgA²⁶, whereas the role of IgA in
424 anti-giardial defense is less clear for *G. intestinalis* GS^{26,68}. We speculate that *G.*
425 *muris* is more resistant to elimination by innate factors such as competition with the
426 normal microbiota, or host production of reactive oxygen species and/or NO, whereas
427 GS is more sensitive to innate factors and eliminated much faster within 1-3 weeks
428 (while *G. muris* clearance requires 4-8 weeks). Future insights into the importance of
429 innate factors in *G. muris* infection should be facilitated by the availability of the
430 complete genome sequence.

431 Sub-telomeric regions in parasitic protozoa often contain arrays of expanded
432 gene families that are under positive selection by the immune system⁶⁵. The relaxed
433 evolutionary pressure offered by keeping pseudogenized copies of surface antigens
434 might be an advantage for *G. muris* that allows genetic drift and recombination to
435 drive rapid and stealthy diversification, thus avoiding elimination by adaptive immune
436 defenses. It was previously reported that *G. muris* is capable of antigenic variation
437 and encodes VSP genes with high similarity and conserved structural features
438 (CRGKA pentapeptide) to those *G. intestinalis*⁶⁹. We failed to identify close
439 homologs to the previously sequenced *G. muris* VSPs in our *G. muris* genome. The
440 linear ψ VSP arrays in *G. muris* have previously been described in *G. intestinalis*⁷⁰.

441 Our phylogenetic analyses of *G. muris* VSPs and ψ VSPs revealed evidence of
442 recombination or segmental gene conversion, as previously demonstrated in *G.*
443 *intestinalis*⁷⁰. However, we recognized two clear differences in the VSP repertoire in
444 *G. muris* and *G. intestinalis*. First, *G. muris* encodes a low number of intact VSP loci
445 that are located internally on the chromosomes. Second, the ψ VSP arrays are almost
446 exclusively telomere adjacent, as opposed to *G. intestinalis* where this tendency is not
447 apparent⁷⁰. These aspects of the *G. muris* VSP repertoire resemble the antigenic
448 variation systems of *Pneumocystis* spp. and *Trypanosoma brucei*^{71,72}. Despite clear
449 mechanistic differences, all these systems have converged on having large reservoirs
450 of mostly telomeric positioned, arrayed genes that are transcriptionally silent and are
451 sources for recombination and gene conversion into expression sites.

452 The function of ankyrin repeat proteins and NEK kinases remains mostly
453 unknown. They represent, together with VSPs, the most dynamic protein families in
454 the *Giardia* genomes⁵⁰. The *Giardia* NEK kinases lack transmembrane domains and
455 have been suggested to target and localize to different intracellular structures with
456 their ankyrin repeats⁵⁰ and many of the *G. intestinalis* NEK kinases localize to
457 cytoskeletal structures, including the flagella and adhesive disc⁹. Rearrangements and
458 duplications in the *G. muris* genome are frequently associated with these large gene
459 families (Fig. 2A and Fig. 4), indicating they might serve as anchoring-points for
460 recombination.

461 Lateral gene transfer is an important shaping factor in the evolution of
462 metabolism in protists⁷³. The origins of laterally transferred genes in *G. muris* are here
463 inferred to be by prokaryotic sources that are members of the gastrointestinal flora, in
464 agreement to previous observations⁷⁴. Most of the putative differences in metabolic
465 potential in the *Giardia* genomes are attributable to lateral gene transfers, either by

466 lineage specific gene gain or loss. For example, the ability to utilize mannose has been
467 introduced from bacteria of the genus *Alistipes* via lateral gene transfer. This event is
468 supported by phylogenetic reconstruction, shows a high degree of sequence
469 conservation (>70% at the amino acid level) and displays maintained gene order to
470 the one seen in the closest related bacterial lineages (Figure S8B). Interestingly,
471 several lateral gene transfers were found clustered in amplified areas of the *G. muris*
472 chromosome. Curiously, *G. intestinalis* assemblage B also maintains clustered copies
473 of arginase and 2,5-diketo-D-gluconic acid reductase, while these genes have both
474 been lost in the *G. intestinalis* assemblage A and E lineages.

475 Anti-microbial peptides of several classes, such as defensins and trefoil-factor 3,
476 are up-regulated in the small intestine of *G. muris* infected mice²⁹. Secreted cysteine
477 proteases from *G. intestinalis* have been shown to be able to degrade defensins¹². We
478 detected prominent expression of several cysteine proteases in *G. muris*. The protease
479 with the highest expression is suggested to have a role in encystation and excystation⁵⁴.
480 Its *G. intestinalis* homolog is up-regulated and secreted during interactions with human
481 intestinal epithelial cells⁸ and it cleaves chemokines, tight junction proteins and
482 defensins^{12,75}. Thus, this is most likely also an important virulence factor in *G. muris*.

483 Our results from this study are summarized in a model of the evolution of
484 *Giardia*'s virulence traits in Figure 5. A number of characters important for *Giardia*'s
485 ability to infect the intestine of mammals are pre-parasitic inventions (such as
486 modified mitochondria and differentiation into transmissive cysts, Fig. 5) and some
487 are found in all diplomonad parasites (e.g. loss of metabolic functions, streamlined
488 microtubular cytoskeleton, expansion of gene families like ankyrins and cysteine
489 proteases and loss of introns, Fig. 5). *Giardia* specific innovations include the
490 adhesive disc for attachment, VSPs and High Cysteine Membrane Proteins for

491 antigenic variation and mitosomes involved in Fe-S complex synthesis (Fig. 5).
492 Whereas some are only found in *G. muris* (e.g. metabolic genes involved in
493 microbiota interactions, Fig. 5), suggesting adaptation to the intestinal environment of
494 mice. Our data shows that the environment in the host's intestine, most of all the
495 immune system and the microbiota, apply selective pressure for changes in the
496 genome, metabolic potential and the parasite surface proteome.

497 **Methods**

498 **Cell preparation and nucleic acid extraction.**

499 4.5×10^7 muris trophozoites (day 7 post infection) were collected from small intestines
500 of three C57 mice, washed once in PBS and pellet frozen at -80°C (Biosample
501 SAMN11231832). Viable cysts of *G. muris* isolate Roberts-Thomson passaged
502 through mice were obtained from Waterborne Inc. These cysts had been purified from
503 fecal material using Percoll and sucrose gradients (Biosample SAMN11231833).
504 DNA and RNA were extracted from 1×10^7 cysts using standard methods.
505 1×10^8 cysts were excysted according to the procedure in Feely et al.⁷⁶ (Biosample
506 SAMN11231834). RNA was purified from cell material equivalent to 1×10^7 cysts.
507 DNA for long-read sequencing was prepared from the remaining cysts as described in
508 Supplementary Methods.

509 **Sequencing, assembly and annotation**

510 Total genomic DNA was sequenced using both Illumina MiSeq and PacBio RS II
511 sequencers. The stranded transcriptome mRNA and the miRNA libraries were
512 sequenced with Illumina HiSeq 2000 system. The RNA samples extracted from
513 excysted cells and cysts prior excystation were prepared using the TruSeq stranded
514 mRNA sample preparation kit and sequenced by HiSeq 2500.
515 PacBio long reads were assembled de novo using the SMRT Analysis (v2.3.0)
516 pipeline. A detailed description of genome sequencing, assembly, annotation, synteny
517 analyses and RNA-Seq is available in Supplementary Methods.

518 **Pathway analysis**

519 The metabolic pathways of *G. muris* and *G. intestinalis* WB were predicted with a
520 combination of BlastKOALA⁷⁷ implemented in KEGG⁷⁸, Pathway Tools v21.5⁷⁹ and
521 GiardiaDB⁸⁰. The different predictions were combined and manually curated under

522 Pathway Tools⁷⁹. Pathway Tools function pathway hole filler⁸¹ was used to further
523 complete the pathway, and transport inference parser⁸² was used to infer transport
524 reaction(s) for transporters which were then verified with Conserved Domain
525 databases⁸³.

526 **Phylogenetic analysis**

527 *G. muris* sequences were used as queries to retrieve at least 5000 hits with e-value
528 <0.001, using BLASTP against the nr database and the organism-specific proteomes.
529 The datasets were aligned in the forward and reverse orientation using MAFFT
530 v6.603b⁸⁴ and PROBCONS v1.12⁸⁵. The four resulting alignments were combined
531 with T-COFFEE⁸⁶ and trimmed by BMGE v1.12⁸⁷. Maximum Likelihood (ML) trees
532 were computed using IQtree v1.6.5⁸⁸ under LG4X substitution model⁸⁹. Branch
533 supports were assessed using ultrafast bootstrap approximation (UFboot)⁹⁰ with
534 1,000 bootstrap replicates and 1,000 replicates for SH-like approximate likelihood
535 ratio test (SH-aLRT)⁹¹. A detailed description of the phylogenetic analyses is
536 available in Supplementary Methods.

537

538

539 **Figures**

540 **Figure 1. Circular representation of the *G. muris* chromosomes.** From outside
541 inward: five chromosomes, GC percent, unique genes (grey) including unique
542 metabolic genes in Table 2 (red), ARPs (greenblue) / NEKs (pink), VSPs (orange) /
543 ψ VSPs (blue) / HCMPs (purple), Coding percent / 5 kbp (green if ≤ 0.5), # SNPs /

544 kbp \geq 5 (red), BLASTN matches with $>95\%$ identity and > 1000 bp in size. Circular
545 plot was drawn with circlize⁹².

546 **Figure 2. Examples of synteny between *G. muris* and *G. intestinalis*.** A) A 50 kbp
547 region on chromosome 3 which share synteny to a 58 kbp region on chromosome 5 in
548 WB. Synteny plot was plotted using genoplots⁹³. Shades of red and blue represent
549 forward and inverted matches between orthologs. Genes are drawn as arrows in blue.
550 ARPs in red, NEKs in dark red, and VSPs in orange. Dark grey filled genes are
551 unique genes to that genome in comparison to the other. B) A 14 kbp region on
552 chromosome 1 which shares synteny to a 16 kbp region on chromosome 5 in *G.*
553 *intestinalis*. It uses the same color scheme as in A. C) Violin plots of intergenic sizes
554 of neighboring positional orthologs of *G. muris* and *G. intestinalis*, and the grey
555 vertical line represents the median intergenic size of *G. muris*. D) Violin plots of
556 positional orthologs sizes of *G. muris* and *G. intestinalis*, and the grey vertical line
557 represents the median ortholog size of *G. muris*. E) Histogram of the positional
558 ortholog size difference between *G. muris* and *G. intestinalis*.

559 **Figure 3. Gene regulation and organization of VSPs in *G. muris*.** A) Promoter
560 motifs shared by encystation-related genes. Motif 1 (gold in B) represents the general
561 promoter motif positioned directly adjacent to the start codon. Motif 2 (teal in B)
562 resembles the encystation-regulated promoter previously identified in *G. intestinalis*
563 (52). B) The distribution and position of motif 1 (gold) and motif 2 (teal) in chosen
564 genes regulated during encystation. C) *Giardia* cyst wall proteins. Cyst wall protein 3
565 is missing in *G. muris*. Signal peptide (pink). Acidic LRR-domain (grey). C-terminal
566 basic extension (green).

567 **Figure 4. Synteny plot of two duplicated regions in the genome.** A, B) Shades of
568 red and blue represent forward and inverted matches between neighboring sequences.

569 ARPs are drawn in red, NEKs in dark red and VSPs in orange. Pseudogenized genes
570 are drawn in dashed lines. Dark grey filled genes are unique genes to *G. muris* in
571 comparison with *G. intestinalis*. Point mutation in arginase is marked in red asterisk.
572 Homologous sequences that are not annotated in the genome are drawn in dashed box
573 on sides of the backbone grey line. Genes discussed in the paper: 21985 and 21992
574 encode the intact arginase genes, 21986 and 21992 encode the enzyme 2,5-diketo-D-
575 gluconic acid reductase, whereas 21984 and 21994 are hypothetical proteins in A;
576 14480 encodes 2,5-diketo-D-gluconic acid reductase, 24746 encodes CMD and 24745
577 encodes ketosteroid isomerase-like protein in B.

578 **Figure 5. Model of the evolution of virulence traits in *Giardia* parasites.** A set of
579 important diplomonad evolutionary innovations and their chronology is depicted at relevant
580 phylogenetic nodes. A) Free-living fornicate ancestor. B) Diplomonad ancestor. C) *Giardia*
581 ancestor. D) *Giardia muris*.

582

583 **Figure S1. A)** *G. muris* cells propagate as trophozoites in the mice small intestine.
584 The cell trophozoites undergo encystation into cysts and are shed in the feces. The
585 cysts are prompted to excyst as they encounter the conditions of the digestive tract.
586 The excyzoites transition into the vegetative trophozoites. Cell material for genomic
587 DNA was harvested from *G. muris* cysts from infected mice. **B)** Sub-telomeric
588 regions in *G. muris*. Sequence motifs in the 10 unitigs terminating in (TAGGG)ⁿ
589 telomeric sequences (purple). Ribosomal DNA sequences (black). Satellite-like
590 repeated sequences (green).

591 **Figure S2. Chromosome-scale synteny between *G. muris* and *G. intestinalis* WB.**
592 *G. muris* chromosomes are color coded. The links between the two genomes

593 representing sequence similarity identified by MUMmer promoter⁹⁴, are colored after
594 the *G. muris* chromosomal color. Circular plot was drawn with circlize⁹².
595 **Figure S3. Sequence logo around start codon (A) and stop codon (B) in *G. muris*.**
596 **C) *G. muris* gene expression from *in vivo* trophozoites.** Genes showing top-ranked
597 FPKM ratios (\log_2 (Trophs/Cysts)) between *in vivo* derived trophozoites and cysts
598 (cut-off value ≥ 4). The value is indicated in the teal circle. An orange dot indicates
599 that the ortholog of this gene is upregulated (>2 fold) in *G. intestinalis* at 7 hrs, 12 hrs
600 or 22 hrs into encystation³⁹.

601 **Figure S4. *Trans*- and *cis*-introns in *G. muris*.** A) Aligned *trans*- and *cis*-splice
602 junctions in *G. muris*. B) Splice-site associated motifs in *trans*-spliced exons. Base-
603 pairing stems are underlined. Bases in red denote the predicted cleavage motif. C)
604 Consensus alignment of the cleavage motifs in *G. muris*. D-E) Predicted secondary
605 structures of D) *cis*-spliced and E) *trans*-spliced genes. The secondary structure was
606 predicted using the mfold webserver using default settings and manually curated. GU-rich
607 stretches and cleavage motifs are boxed.

608 **Figure S5. A) Conservation of the VSP C-terminal in *G. muris* and *G. intestinalis*.**
609 Sequence logos of conserved motifs in the C-terminal of VSPs in *G. muris* (n=25)
610 (upper panel) and *G. intestinalis* WB (n=183) (lower panel) were generated using
611 MEME v5.0.5. Transmembrane and cytosolic pentapeptides were predicted by the
612 Phobius webserver. B) **VSP and ψ VSP evolutionary dynamics.** A) The ψ VSP arrays
613 (defined as >3 ψ VSP genes arrayed) in the *G. muris* genome arranged by
614 chromosome, termini arranged to the right. Arrays of ψ VSPs without chromosomal
615 context are shown below as orphan arrays. C) VSP cluster mostly together in a
616 phylogeny whereas most ψ VSP gene arrays show relaxed clustering. Intact VSP
617 genes are colored red, while ψ VSP genes are colored according to their resident

618 ψ VSP array in panel A. Genes with black label are not found in arrays. **D)**

619 **Expression measurements of VSP and ψ VSP in *G. muris*.** VSP and ψ VSP RNA-
620 seq read counts (FPKM) from trophozoites, cysts and excyzoites. Size of the circle
621 corresponds to expression level.

622 **Figure S6. A) Phylogenetic tree of *G. muris* and *G. intestinalis* NEKs. B) Phylogenetic**
623 **tree of *G. muris* and *G. intestinalis* ARPs. C) Network analysis of *G. muris* and *G.***
624 ***intestinalis* NEKs and ARPs. *G. muris* genes are represented in circle whereas *G.***
625 ***intestinalis* genes are shown in triangle. Blue indicates NEK and pink as ARP. Edges are**
626 **weighted and scaled by reciprocal BLAST scores with evalue 1e-05 as the cutoff.**

627 **Figure S7. Phylogenetic tree of C1 family cysteine proteases in diplomonads. Black:**
628 ***G. muris*, Red: *S. salmonicida*, Cyan: *G. intestinalis* GS/M, Green: *G. intestinalis* P15,**
629 **Purple: *G. intestinalis* WB.**

630 **Figure S8. Phylogenetic analyses of HGT genes. A) Fructokinase. B) Mannose-6-**
631 **phosphate isomerase. C) Transcriptional regulator/sugar kinase. D) Arginase. E) 2,5-**
632 **diketo-D-gluconic acid reductase 1. F) 2,5-diketo-D-gluconic acid reductase 2. G)**
633 **Carboxymuconolactone decarboxylase. H) Ketosteroid isomerase. I) Tryptophanase.**
634 **J) Tae4. K) Quorum-quenching N-acyl-homoserine lactonase.**

635

636 **Supplementary Methods: Additional methods section**

637 **Declarations**

638 **Ethics approval and consent to participate**

639 All animal studies were approved by the Institutional Animal Care and Use

640 Committees of the University of California, San Diego, USA.

641 **Consent for publication**

642 No applicable.

643 **Availability of data and material**

644 Raw DNA and RNA sequence reads are archived at NCBI Sequence Read Archive

645 (SRA) under accession number SRR8858297 - SRR8858305.

646 This Whole Genome Shotgun project has been deposited at DDBJ/EMBL/GenBank

647 under the accession PRJNA524057. The version described in this paper is version

648 VDLU00000000.1. The datasets generated and analyzed in this project are also

649 available via GiardiaDB.

650 **Competing interests**

651 The authors declare that they have no competing interests.

652 **Funding**

653 The project was supported by funding from Vetenskapsrådet-M (2017-02918) to SGS.

654 LE was supported by NIH grant DK120515.

655 **Authors' contributions**

656 SGS, LE, JA and JJH conceived the study. JJH and EE performed RNA and DNA

657 extractions. FX assembled the genome. FX, AJG, JJH, AA, DP, JA, SGS and EE

658 annotated the genome. FX, JJH, AJG and SGS analyzed the data. All authors wrote

659 the paper.

660 **Acknowledgements**

661 Elaine Hanson is acknowledged for technical support performing trophozoite purifications

662 from infected mice.

663 Phylogenetic analyses were performed on resources provided by the Swedish National

664 Infrastructure for Computing (SNIC) at Uppsala Multidisciplinary Center for Advanced

665 Computational Science (UPPMAX).

666

667 **Tables and additional files**

668 **Tables**

669 **Table 1** - Comparison of genome content between *G. muris*, *G. intestinalis*, *S.*

670 *salmonicida* and *K. bialata*

671

Species	<i>G. muris</i>	<i>G. intestinalis</i>	<i>S. salmonicida</i>	<i>K. bialata</i>
Genome size (Mbp)	9.8	12.6	12.9	51.0
Chromosomes	5 (59)	5 (35)	9 (233)	ND
(scaffolds)				(11,564)
G+C %	54.7	46.3	33.4	49.4
Number of protein	4653	4963	8067	17,389
encoding genes				
Mean / Median	578 / 428	635 / 457	373	333
protein size (aa)				
*Mean / Median	264 / 37	470 / 81	421	597
intergenic size (bp)				
§Coding density %	84.5 / 88.6	81.5 / 84.7	72.1	ND
Number of introns	3 <i>cis</i> ,5 <i>trans</i>	8 <i>cis</i> ,5 <i>trans</i>	3 <i>cis</i>	124,912
tRNA genes	68	65	145	ND
ASH %	0.016	0.028	0.15	ND
References	This study	Ref ³²	Ref ⁹⁵	Ref ⁹¹

672 * Mean/Median intergenic distance is based on all RNAs (mRNAs, tRNAs, rRNAs),

673 but not pseudogenized genes.

674 § Coding density: First value is based on all RNAs (mRNAs, tRNAs, rRNAs), but not

675 pseudogenized genes; Second value is based on all RNAs including ψ VSP.

676

677 **Table 2** - Lateral gene transfers in *G. muris* (*Gm*), *G. intestinalis* (*Gi* WB, *Gi* GS, *Gi* P15),

678 *S. salmonicida* (*Ss*), *Treponomas spp.* (*Trep*), *K. bialata* (*Kb*) and *T. vaginalis* (*Tv*).

Species	<i>Gm</i>	<i>Gi</i> <i>WB</i>	<i>Gi</i> <i>GS</i>	<i>Gi</i> <i>P15</i>	<i>Ss</i>	<i>Trep</i>	<i>Kb</i>	<i>Tv</i>
2.5-diketo-D-gluconic acid reductase	X		X				X	X
Arginase	X		X					X
Carboxymuconolactone decarboxylase	X		X					
Ferritin-like	X							X
Fructokinase	X							
Ketosteroid isomerase	X							
L-ascorbate-6-phosphate lactonase	X		X					
Maltose-O-acetyltransferase	X				X	X		
Mannose-6-phosphate isomerase	X				X			
Phosphopantothenate-cysteine ligase	X				X			
Quorum-quenching N-acyl-homoserine lactonase	X		X					
Ribonuclease 3	X						X	X
Tae4	X							
Tryptophanase	X				X	X		
β-phosphoglucomutase		X	X	X	X	X		
Extracellular nuclease		X	X	X	X	X	X	X

Flavoheмоprotein	X	X	X				
Glycerol kinase	X	X	X	X	X		
Inositol-3-phosphate synthase	X	X	X				
Methyltransferase	X	X	X				
NADPH-ferrihemoprotein	X	X	X	X	X		
Purine nucleoside phosphorylase	X	X	X	X	X	X	X
Pyruvate phosphate dikinase	X	X	X		X		
Sugar/H⁺ symporter	X	X	X				
Threonine dehydratase	X	X	X	X	X		

679

680 **Table 3** – Summaries of gene families within *Giardia*

	<i>G. intestinalis</i> *	<i>G. muris</i> *
NEK	184 (26)	216 (23)
ARP-1	269 (5)	298 (6)
ARP-2	33	86 (16)
ARP-3	8	33
VSP	133	26
ψVSP	208	239

681 * Values in () indicate the number of pseudogenized copies.

682 **Table 4** – Arrangement of VSP genes in the *G. muris* genome.

Chr	Geneid 1	Geneid 2	Arrangement	Gene size (aa)	Distance (bp)	Genes in between
1	20512	13275	-> <-	596	2025	ψARP-2

1	21145	21149	-> <-	594	1658	ψARP-2
1	13124	21374	-> <-	620	8400	ψARP-2, alpha-tubulin
2	16008	21957	<- ->	624	8794	NEK, ARP-2
2	22301	22304	-> <-	523	1656	ψARP-2
2	12920	22758 (22764, 22769)	<- ->	515	10612 (6824, 8428)	NEK, TRP, Zinc
4	24220	24228	-> <-	623	8303	ψARP-2, alpha-tubulin
5	24787	24792	<- ->	619	7158	NEK, Zinc, ARP-2

683

684 **Table S1 A)** One-to-one ortholog list between *G. muris* and *G. intestinalis* WB. Orthologs
685 inferred by OrthoMCL. **B)** RNA-Seq expression values in excel file. Sorted according to
686 FPKM values. **C)** VSP and ψVSP gene metadata. '-' in the pentapeptide motif represents
687 lack of the 3' end sequence. **D)** List of ventral disc proteins in *G. muris*. **E)** Proteases
688 detected in *G. muris* and *G. intestinalis*. **F)** Detailed information about hits of unique genes
689 in *G. muris*.

References

1. Poulin, R. & Randhawa, H. S. Evolution of parasitism along convergent lines: from ecology to genomics. *Parasitology* **142**, S6–S15 (2015).
2. Hupalo, D. N., Bradic, M. & Carlton, J. M. The impact of genomics on population genetics of parasitic diseases. *Current Opinion in Microbiology* **23**, 49–54 (2015).
3. Jackson, A. P. *et al.* Kinetoplastid Phylogenomics Reveals the Evolutionary Innovations Associated with the Origins of Parasitism. *Current Biology* **26**, 161–172 (2016).
4. Einarsson, E., Ma'ayeh, S. & Svärd, S. G. An up-date on *Giardia* and giardiasis. *Current Opinion in Microbiology* **34**, 47–52 (2016).
5. Cacciò, S. M., Lalle, M. & Svärd, S. G. Host specificity in the *Giardia duodenalis* species complex. *Infect. Genet. Evol.* **66**, 335–345 (2018).
6. Fink, M. Y. & Singer, S. M. The Intersection of Immune Responses, Microbiota, and Pathogenesis in Giardiasis. *Trends in Parasitology* **33**, 901–913 (2017).
7. Ma'ayeh, S. Y. *et al.* Responses of the Differentiated Intestinal Epithelial Cell Line Caco-2 to Infection With the *Giardia intestinalis* GS Isolate. *Front Cell Infect Microbiol* **8**, 244 (2018).
8. Ma'ayeh, S. Y. *et al.* Characterization of the *Giardia intestinalis* secretome during interaction with human intestinal epithelial cells: The impact on host cells. *PLoS Negl Trop Dis* **11**, e0006120 (2017).
9. Nosala, C., Hagen, K. D. & Dawson, S. C. 'Disc-o-Fever': Getting Down with *Giardia*'s Groovy Microtubule Organelle. *Trends in Cell Biology* **28**, 99–112 (2018).

10. Nosala, C. & Dawson, S. C. The Critical Role of the Cytoskeleton in the Pathogenesis of *Giardia*. *Current Clinical Microbiology Reports* **2**, 155–162 (2015).
11. McNally, S. G. & Dawson, S. C. Eight unique basal bodies in the multi-flagellated diplomonad *Giardia lamblia*. *Cilia* **5**, (2016).
12. Liu, J. *et al.* Secreted *Giardia intestinalis* cysteine proteases disrupt intestinal epithelial cell junctional complexes and degrade chemokines. *Virulence* **9**, 879–894 (2018).
13. Ortega-Pierres, G. *et al.* Giardipain-1, a protease secreted by *Giardia duodenalis* trophozoites, causes junctional, barrier and apoptotic damage in epithelial cell monolayers. *International Journal for Parasitology* **48**, 621–639 (2018).
14. Amat, C. B. *et al.* Cysteine Protease–Dependent Mucous Disruptions and Differential Mucin Gene Expression in *Giardia duodenalis* Infection. *The American Journal of Pathology* **187**, 2486–2498 (2017).
15. Cotton, J. A. *et al.* *Giardia duodenalis* Cathepsin B Proteases Degrade Intestinal Epithelial Interleukin-8 and Attenuate Interleukin-8-Induced Neutrophil Chemotaxis. *Infection and Immunity* **82**, 2772–2787 (2014).
16. Cotton, J. A. *et al.* *Giardia duodenalis* Infection Reduces Granulocyte Infiltration in an In Vivo Model of Bacterial Toxin-Induced Colitis and Attenuates Inflammation in Human Intestinal Tissue. *PLoS ONE* **9**, e109087 (2014).
17. Allain, T., Amat, C. B., Motta, J.-P., Manko, A. & Buret, A. G. Interactions of *Giardia* sp. with the intestinal barrier: Epithelium, mucus, and microbiota. *Tissue Barriers* **5**, e1274354 (2017).

18. Mendez, T. L. *et al.* Sphingolipids, Lipid Rafts, and Giardial Encystation: The Show Must Go On. *Current Tropical Medicine Reports* **2**, 136–143 (2015).
19. Stadelmann, B., Merino, M. C., Persson, L. & Svärd, S. G. Arginine Consumption by the Intestinal Parasite *Giardia intestinalis* Reduces Proliferation of Intestinal Epithelial Cells. *PLoS ONE* **7**, (2012).
20. Eckmann, L. *et al.* Nitric Oxide Production by Human Intestinal Epithelial Cells and Competition for Arginine as Potential Determinants of Host Defense Against the Lumen-Dwelling Pathogen *Giardia lamblia*. *The Journal of Immunology* **164**, 1478–1487 (2000).
21. Helmy, Y. A. *et al.* Occurrence and distribution of *Giardia* species in wild rodents in Germany. *Parasites & Vectors* **11**, (2018).
22. Friend, D. S. The fine structure of *Giardia muris*. *The Journal of Cell Biology* **29**, 317–332 (1966).
23. Dann, S. M., Le, C. H. Y., Hanson, E. M., Ross, M. C. & Eckmann, L. *Giardia* Infection of the Small Intestine Induces Chronic Colitis in Genetically Susceptible Hosts. *The Journal of Immunology* **201**, 548–559 (2018).
24. Holberton, D. V. Fine structure of the ventral disk apparatus and the mechanism of attachment in the flagellate *Giardia muris*. *J. Cell. Sci.* **13**, 11–41 (1973).
25. Schaefer, F. W., Rice, E. W. & Hoff, J. C. Factors promoting in vitro excystation of *Giardia muris* cysts. *Trans. R. Soc. Trop. Med. Hyg.* **78**, 795–800 (1984).
26. Langford, T. D. *et al.* Central importance of immunoglobulin A in host defense against *Giardia* spp. *Infect. Immun.* **70**, 11–18 (2002).

27. Davids, B. J. *et al.* Polymeric immunoglobulin receptor in intestinal immune defense against the lumen-dwelling protozoan parasite *Giardia*. *J. Immunol.* **177**, 6281–6290 (2006).
28. Dreesen, L. *et al.* *Giardia muris* infection in mice is associated with a protective interleukin 17A response and induction of peroxisome proliferator-activated receptor alpha. *Infect. Immun.* **82**, 3333–3340 (2014).
29. Manko, A. *et al.* *Giardia* co-infection promotes the secretion of antimicrobial peptides beta-defensin 2 and trefoil factor 3 and attenuates attaching and effacing bacteria-induced intestinal disease. *PLOS ONE* **12**, e0178647 (2017).
30. Saghaug, C. S. *et al.* Human Memory CD4+ T Cell Immune Responses against *Giardia lamblia*. *Clinical and Vaccine Immunology* **23**, 11–18 (2016).
31. Tanifuji, G. *et al.* The draft genome of *Kipferlia bialata* reveals reductive genome evolution in fornicate parasites. *PLOS ONE* **13**, e0194487 (2018).
32. Morrison, H. G. *et al.* Genomic minimalism in the early diverging intestinal parasite *Giardia lamblia*. *Science* **317**, 1921–1926 (2007).
33. Franzén, O. *et al.* Draft genome sequencing of giardia intestinalis assemblage B isolate GS: is human giardiasis caused by two different species? *PLoS Pathog.* **5**, e1000560 (2009).
34. Jerlström-Hultqvist, J. *et al.* Genome analysis and comparative genomics of a *Giardia intestinalis* assemblage E isolate. *BMC Genomics* **11**, 543 (2010).
35. Xu, F. *et al.* The Genome of *Spironucleus salmonicida* Highlights a Fish Pathogen Adapted to Fluctuating Environments. *PLoS Genetics* **10**, (2014).

36. Campbell, S. R., van Keulen, H., Erlandsen, S. L., Senturia, J. B. & Jarroll, E. L. *Giardia* sp.: Comparison of electrophoretic karyotypes. *Experimental Parasitology* **71**, 470–482 (1990).
37. Prabhu, A., Morrison, H. G., Martinez, C. R. & Adam, R. D. Characterisation of the subtelomeric regions of *Giardia lamblia* genome isolate WBC6. *Int. J. Parasitol.* **37**, 503–513 (2007).
38. Franzén, O. *et al.* Transcriptome Profiling of *Giardia intestinalis* Using Strand-specific RNA-Seq. *PLoS Computational Biology* **9**, e1003000 (2013).
39. Einarsson, E. *et al.* Coordinated Changes in Gene Expression Throughout Encystation of *Giardia intestinalis*. *PLOS Neglected Tropical Diseases* **10**, e0004571 (2016).
40. Morf, L. *et al.* The transcriptional response to encystation stimuli in *Giardia lamblia* is restricted to a small set of genes. *Eukaryotic Cell* **9**, 1566–1576 (2010).
41. Barash, N. R. *et al.* *Giardia* Colonizes and Encysts in High-Density Foci in the Murine Small Intestine. *mSphere* **2**, (2017).
42. Pham, J. K. *et al.* Transcriptomic Profiling of High-Density *Giardia* Foci Encysting in the Murine Proximal Intestine. *Frontiers in Cellular and Infection Microbiology* **7**, (2017).
43. Hudson, A. J. *et al.* Evolutionarily divergent spliceosomal snRNAs and a conserved non-coding RNA processing motif in *Giardia lamblia*. *Nucleic Acids Research* **40**, 10995–11008 (2012).

44. Kamikawa, R., Inagaki, Y., Tokoro, M., Roger, A. J. & Hashimoto, T. Split Introns in the Genome of *Giardia intestinalis* Are Excised by Spliceosome-Mediated trans-Splicing. *Current Biology* **21**, 311–315 (2011).
45. William Roy, S. Transcriptomic analysis of diplomonad parasites reveals a trans-spliced intron in a helicase gene in *Giardia*. *PeerJ* **5**, e2861 (2017).
46. Touz, M. C., Conrad, J. T. & Nash, T. E. A novel palmitoyl acyl transferase controls surface protein palmitoylation and cytotoxicity in *Giardia lamblia* VSP palmitoylation. *Molecular Microbiology* **58**, 999–1011 (2005).
47. Touz, M. C. *et al.* Arginine deiminase has multiple regulatory roles in the biology of *Giardia lamblia*. *Journal of Cell Science* **121**, 2930–2938 (2008).
48. Davids, B. J. *et al.* A new family of giardial cysteine-rich non-VSP protein genes and a novel cyst protein. *PLoS ONE* **1**, e44 (2006).
49. Ringqvist, E., Avesson, L., Söderbom, F. & Svärd, S. G. Transcriptional changes in *Giardia* during host-parasite interactions. *Int. J. Parasitol.* **41**, 277–285 (2011).
50. Manning, G. *et al.* The minimal kinome of *Giardia lamblia* illuminates early kinase evolution and unique parasite biology. *Genome Biology* **12**, R66 (2011).
51. Gargantini, P. R., Serradell, M. del C., Ríos, D. N., Tenaglia, A. H. & Luján, H. D. Antigenic variation in the intestinal parasite *Giardia lamblia*. *Current Opinion in Microbiology* **32**, 52–58 (2016).
52. Palm, J. E. D., Weiland, M. E.-L., Griffiths, W. J., Ljungström, I. & Svärd, S. G. Identification of immunoreactive proteins during acute human giardiasis. *J. Infect. Dis.* **187**, 1849–1859 (2003).

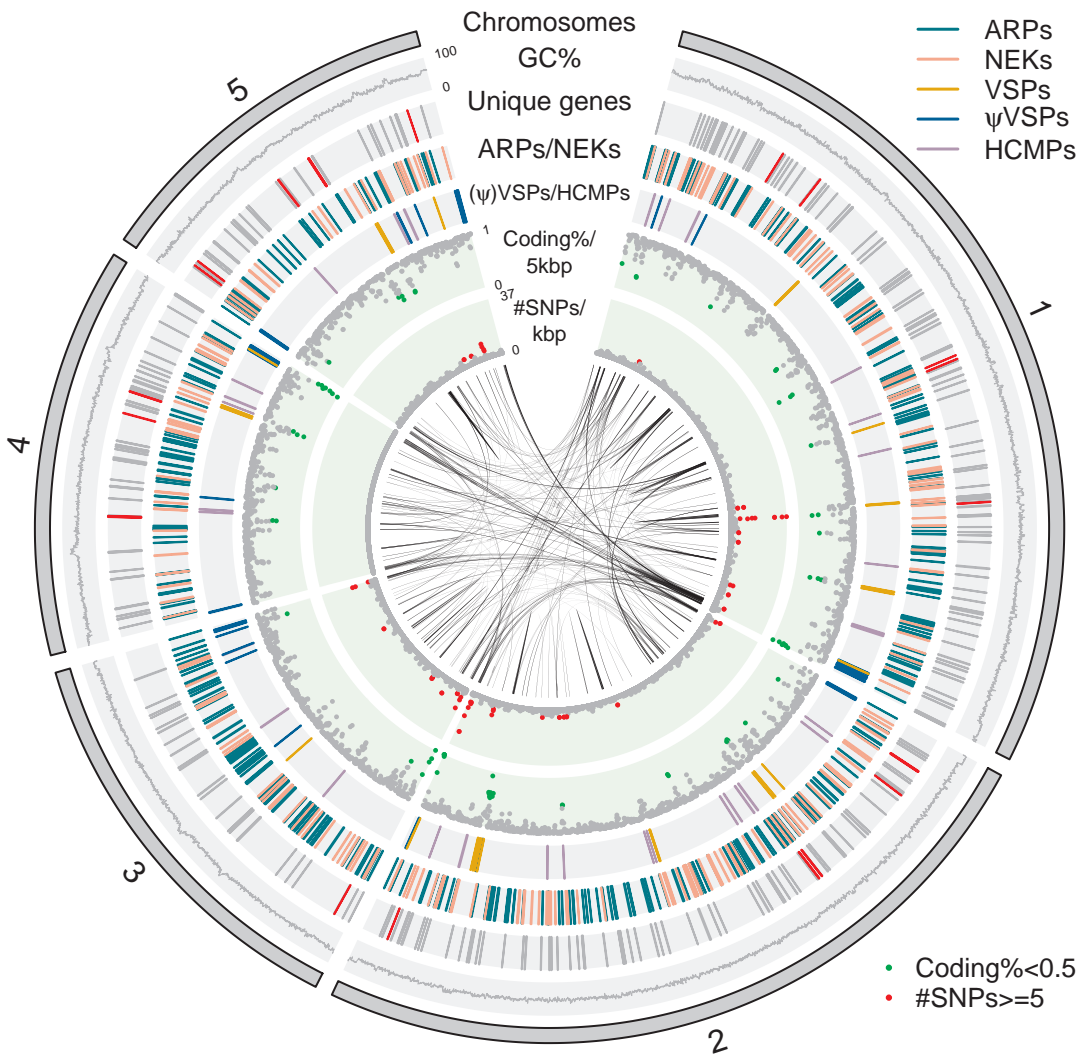
53. Liu, J., Svärd, S. G. & Klotz, C. *Giardia intestinalis* cystatin is a potent inhibitor of papain, parasite cysteine proteases and, to a lesser extent, human cathepsin B. *FEBS Letters* **593**, 1313–1325 (2019).
54. DuBois, K. N. *et al.* Identification of the major cysteine protease of *Giardia* and its role in encystation. *J. Biol. Chem.* **283**, 18024–18031 (2008).
55. Ansell, B. R. E. *et al.* Time-Dependent Transcriptional Changes in Axenic *Giardia duodenalis* Trophozoites. *PLoS Neglected Tropical Diseases* **9**, 1–24 (2015).
56. Stadelmann, B., Merino, M. C., Persson, L. & Svärd, S. G. Arginine Consumption by the Intestinal Parasite *Giardia intestinalis* Reduces Proliferation of Intestinal Epithelial Cells. *PLoS ONE* **7**, e45325 (2012).
57. Mastronicola, D. *et al.* Flavohemoglobin and nitric oxide detoxification in the human protozoan parasite *Giardia intestinalis*. *Biochem. Biophys. Res. Commun.* **399**, 654–658 (2010).
58. Das, P., Lahiri, A., Lahiri, A. & Chakravorty, D. Modulation of the arginase pathway in the context of microbial pathogenesis: A metabolic enzyme moonlighting as an immune modulator. *PLoS Pathogens* **6**, (2010).
59. Krasity, B. C., Troll, J. V., Weiss, J. P. & McFall-Ngai, M. J. LBP/BPI proteins and their relatives: conservation over evolution and roles in mutualism. *Biochemical Society Transactions* **39**, 1039–1044 (2011).
60. Ankarklev, J. *et al.* Comparative genomic analyses of freshly isolated *Giardia intestinalis* assemblage A isolates. *BMC Genomics* **16**, (2015).
61. Lee, J.-H., Wood, T. K. & Lee, J. Roles of Indole as an Interspecies and Interkingdom Signaling Molecule. *Trends in Microbiology* **23**, 707–718 (2015).

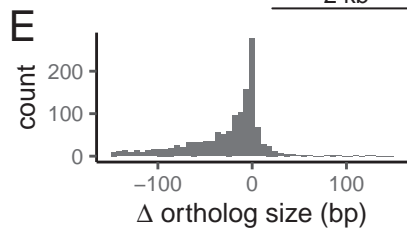
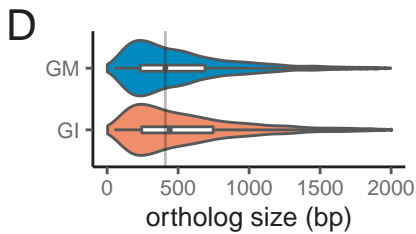
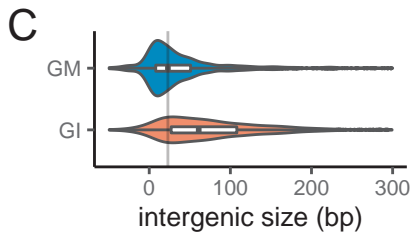
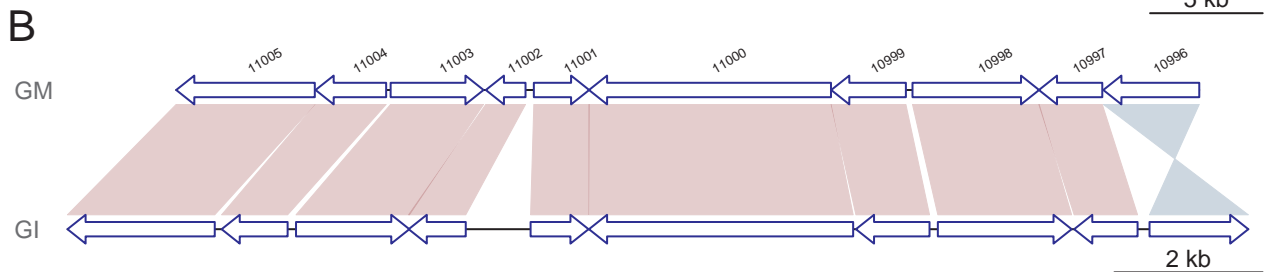
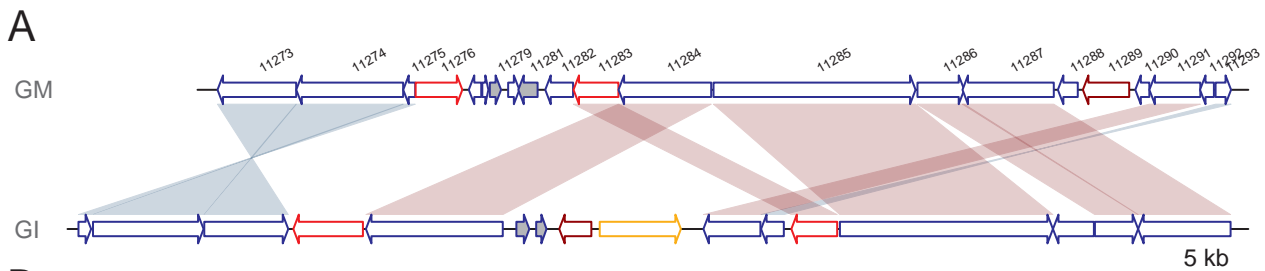
62. Russell, A. B. *et al.* A Widespread Bacterial Type VI Secretion Effector Superfamily Identified Using a Heuristic Approach. *Cell Host & Microbe* **11**, 538–549 (2012).
63. Sana, T. G. *et al.* *Salmonella* Typhimurium utilizes a T6SS-mediated antibacterial weapon to establish in the host gut. *Proceedings of the National Academy of Sciences* **113**, E5044–E5051 (2016).
64. Kim, M. H. *et al.* The molecular structure and catalytic mechanism of a quorum-quenching N-acyl-L-homoserine lactone hydrolase. *Proceedings of the National Academy of Sciences* **102**, 17606–17611 (2005).
65. Leckenby, A. & Hall, N. Genomic changes during evolution of animal parasitism in eukaryotes. *Current Opinion in Genetics & Development* **35**, 86–92 (2015).
66. Peyretaillade, E. *et al.* Exploiting the architecture and the features of the microsporidian genomes to investigate diversity and impact of these parasites on ecosystems. *Heredity (Edinb)* **114**, 441–449 (2015).
67. Singer, S. M. & Nash, T. E. The Role of Normal Flora in *Giardia lamblia* Infections in Mice. *The Journal of Infectious Diseases* **181**, 1510–1512 (2000).
68. Singer, S. M. & Nash, T. E. T-Cell-Dependent Control of Acute *Giardia lamblia* Infections in Mice. *Infection and Immunity* **68**, 170–175 (2000).
69. Ropolo, A. S., Saura, A., Carranza, P. G. & Lujan, H. D. Identification of Variant-Specific Surface Proteins in *Giardia muris* Trophozoites. *Infection and Immunity* **73**, 5208–5211 (2005).
70. Adam, R. D. *et al.* The *Giardia lamblia* vsp gene repertoire: characteristics, genomic organization, and evolution. *BMC Genomics* **11**, 424 (2010).

71. Schmid-Siegert, E. *et al.* Mechanisms of Surface Antigenic Variation in the Human Pathogenic Fungus *Pneumocystis jirovecii*. *mBio* **8**, (2017).
72. Mugnier, M. R., Stebbins, C. E. & Papavasiliou, F. N. Masters of Disguise: Antigenic Variation and the VSG Coat in *Trypanosoma brucei*. *PLoS Pathog.* **12**, e1005784 (2016).
73. Husnik, F. & McCutcheon, J. P. Functional horizontal gene transfer from bacteria to eukaryotes. *Nature Reviews Microbiology* **16**, 67–79 (2018).
74. Alsmark, C. *et al.* Patterns of prokaryotic lateral gene transfers affecting parasitic microbial eukaryotes. *Genome Biology* **14**, R19 (2013).
75. Liu, J., Fu, Z., Hellman, L. & Svärd, S. G. Cleavage specificity of recombinant *Giardia intestinalis* cysteine proteases: Degradation of immunoglobulins and defensins. *Molecular and Biochemical Parasitology* **227**, 29–38 (2019).
76. Feely, D. E., Gardner, M. D. & Hardin, E. L. Excystation of *Giardia muris* induced by a phosphate-bicarbonate medium: localization of acid phosphatase. *J. Parasitol.* **77**, 441–448 (1991).
77. Kanehisa, M., Sato, Y. & Morishima, K. BlastKOALA and GhostKOALA: KEGG Tools for Functional Characterization of Genome and Metagenome Sequences. *Journal of Molecular Biology* **428**, 726–731 (2016).
78. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Research* **45**, D353–D361 (2017).
79. Karp, P. D. *et al.* Pathway Tools version 19.0 update: software for pathway/genome informatics and systems biology. *Briefings in Bioinformatics* **17**, 877–890 (2016).

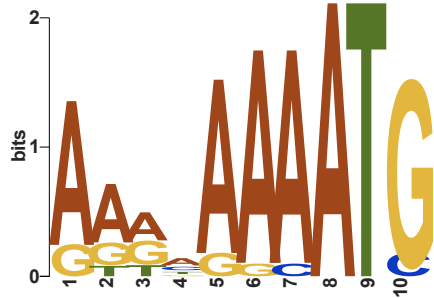
80. Aurrecochea, C. *et al.* GiardiaDB and TrichDB: integrated genomic resources for the eukaryotic protist pathogens *Giardia lamblia* and *Trichomonas vaginalis*. *Nucleic Acids Research* **37**, D526–D530 (2009).
81. Green, M. L. & Karp, P. D. A Bayesian method for identifying missing enzymes in predicted metabolic pathway databases. *BMC Bioinformatics* **5**, 76 (2004).
82. Lee, T. J., Paulsen, I. & Karp, P. Annotation-based inference of transporter function. *Bioinformatics* **24**, i259-267 (2008).
83. Marchler-Bauer, A. & Bryant, S. H. CD-Search: protein domain annotations on the fly. *Nucleic Acids Res.* **32**, W327-331 (2004).
84. Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic acids research* **30**, 3059–3066 (2002).
85. Do, C. B., Mahabhashyam, M. S. P., Brudno, M. & Batzoglou, S. ProbCons: Probabilistic consistency-based multiple sequence alignment. *Genome research* **15**, 330–40 (2005).
86. Notredame, C., Higgins, D. G. & Heringa, J. T-Coffee: A novel method for fast and accurate multiple sequence alignment. *J. Mol. Biol.* **302**, 205–217 (2000).
87. Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new software for selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evol. Biol.* **10**, 210 (2010).
88. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular Biology and Evolution* **32**, 268–274 (2015).

89. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermini, L. S. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature Methods* **14**, 587 (2017).
90. Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2: Improving the Ultrafast Bootstrap Approximation. *Molecular Biology and Evolution* **35**, 518–522 (2018).
91. Guindon, S. *et al.* New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance of PhyML 3.0. *Systematic Biology* **59**, 307–321 (2010).
92. Gu, Z., Gu, L., Eils, R., Schlesner, M. & Brors, B. circlize implements and enhances circular visualization in R. *Bioinformatics* **30**, 2811–2812 (2014).
93. Guy, L., Kultima, J. R. & Andersson, S. G. E. genoPlotR: comparative gene and genome visualization in R. *Bioinformatics* **26**, 2334–2335 (2010).
94. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol.* **5**, R12 (2004).
95. Xu, F. *et al.* The Genome of *Spironucleus salmonicida* Highlights a Fish Pathogen Adapted to Fluctuating Environments. *PLOS Genetics* **10**, e1004053 (2014).

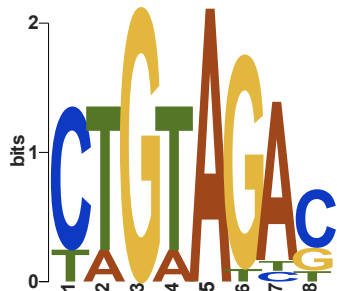




A Motif 1: General promoter motif



Motif 2: Encystation box



B

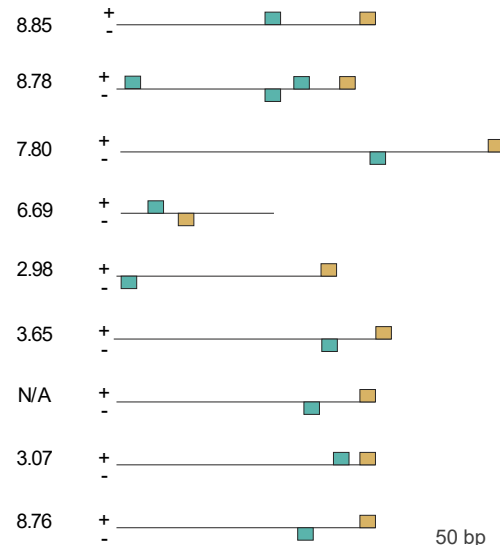
Gene

GMRT_13204	Cyst Wall Protein 2
GMRT_10774	Cyst Wall Protein 1
GMRT_13400	TMP-1 protein
GMRT_13283	SacB
GMRT_15542	Oxidoreductase
GMRT_10283	Synaptic glycoprotein SC2
GMRT_11019	MYB
GMRT_14284	Hypothetical protein
GMRT_15202	UDP-N-acetylglucosamine pyrophosphorylase

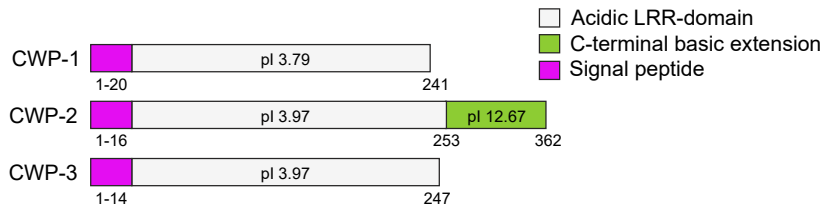
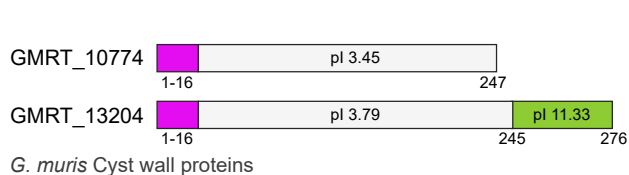
Up reg
(log₂)

Promoter motifs

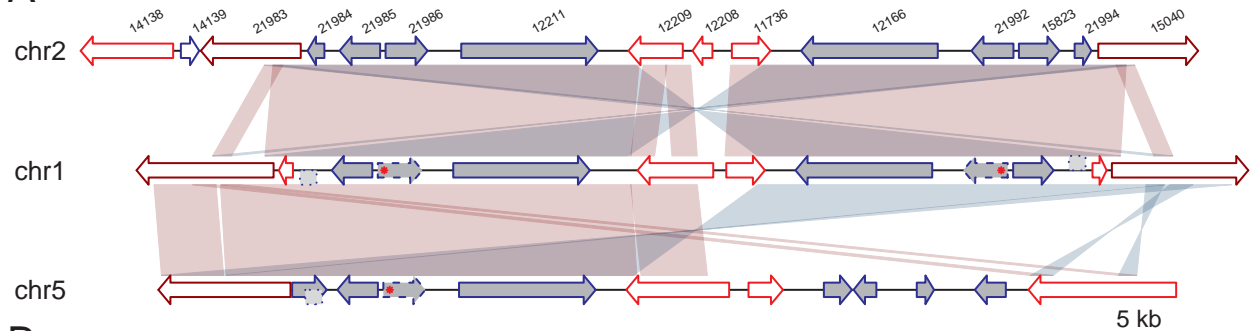
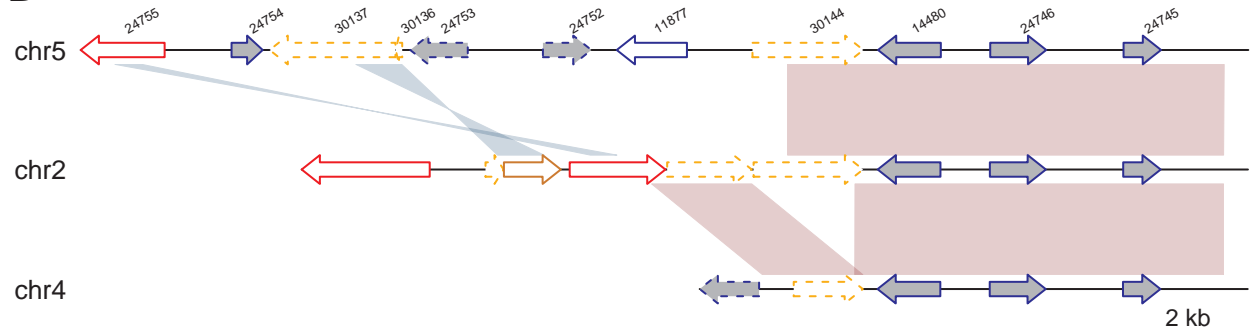
■ Motif 1
■ Motif 2



C



□ Acidic LRR-domain
■ C-terminal basic extension
■ Signal peptide

A**B**

Diplomonad adaptations

Metabolic adaptation : gene loss / gain via HGT

Cysteine-rich gene families

Heavy intron loss

