

1 **Sense-antisense gene overlap causes evolutionary retention of the few**
2 **introns in *Giardia* genome and the implications**

3

4

5 Min Xue^{1,2}, Bing Chen¹, Qingqing Ye¹, Jingru Shao¹, Zhangxia Lyu¹, Jianfan Wen^{1,*}

6

7 ¹State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of
8 Zoology, Chinese Academy of Sciences, Kunming, Yunnan 650223, China.

9 ²Kunming College of Life Science, University of Chinese Academy of Sciences,
10 Kunming, Yunnan 650204, China.

11

12

13 ***Corresponding author:** Jianfan Wen, E-mail: wenjf@mail.kiz.ac.cn

14 E-mail for other authors:

15 Min Xue: xuemin@mail.kiz.ac.cn

16 Bing Chen: chenbing@mail.kiz.ac.cn

17 Qingqing Ye: yeqingqing@mail.kiz.ac.cn

18 Jingru Shao: shaojr@mail.kiz.ac.cn

19 Zhangxia Lyu: zhangxialv@outlook.com

20

21

22

23

24 **Abstract**

25 **Background:** It is widely accepted that the last eukaryotic common ancestor (LECA)
26 and early eukaryotes were intron-rich and intron loss dominated subsequent evolution,
27 thus the presence of only very few introns in some modern eukaryotes must be the
28 consequence of massive loss. But it is striking that few eukaryotes were found to have
29 completely lost introns. Despite extensive research, the causes of massive intron
30 losses remain elusive, and actually the reverse question – how the few introns are
31 retained under the pressure of loss is equally significant but was rarely studied, except
32 that it was conjectured that the essential functions of some introns prevent their loss.
33 The extremely few (eight) spliceosome-mediated cis-spliced introns in the relatively
34 simple genome of *Giardia lamblia* provide an excellent opportunity to explore this
35 question.

36 **Results:** Our investigation of the intron-containing genes and introns in *Giardia*
37 found three types of intron distribution patterns: ancient intron in ancient gene,
38 relatively new intron in ancient gene, and relatively new intron in relatively new gene,
39 which can reflect to some extent the dynamic evolution of introns in *Giardia*. Not
40 finding any special features or functional importance of these introns responsible for
41 the retention, we noticed and experimentally verified that some intron-containing
42 genes form sense-antisense gene pairs with functional genes on their complementary
43 strands, and that the introns just reside in the overlapping regions.

44 **Conclusions:** In *Giardia*'s evolution, despite constant pressure of intron loss, intron
45 gain can still occur in both ancient and newly-evolved genes, but only a few introns
46 have been retained; the evolutionary retention of introns is most likely not due to the
47 functional constraint of the introns themselves but the causes outside of introns, such
48 as the constraints imposed by other genomic functional elements overlapping with the

49 introns. These findings can not only provide some clues to find new genomic
50 functional elements -- in the areas overlapping with introns, but suggest that
51 “functional constraint” of introns may not be necessarily directly associated with
52 intron loss and gain, or that the real functions or the way of functioning of introns are
53 probably still outside of our current knowledge.

54

55 **Keywords:** Evolutionary retention of introns, Gene overlap, Intron evolution,
56 Genome evolution, *Giardia lamblia*.

57

58 **Background**

59 Spliceosomal introns are a common feature of all eukaryotic nuclear genomes, but
60 their number and density in a genome vary dramatically among different species [1, 2],
61 ranging from less than 0.5 intron/gene in some protists such as the Microsporidian
62 *Encephalitozoon* species [3] and *Cyanidioschyzon merolae* [4] to over 18 per gene in
63 *Symbiodinium minutum* [5] (even larger than those of most mammals, which is
64 generally over eight/gene [6]). Accumulating evidence suggests that the LECA and
65 early eukaryotes were relatively intron rich, with subsequent genome evolution
66 dominated by intron loss, and thus those contemporary eukaryotes with remarkably
67 few introns must have experienced massive intron loss secondarily [7-9]. But,
68 interestingly, no eukaryotes with sequenced genomes so far have been found to have
69 completely lost their introns except the two Microsporidia species, *Nematocida parisii*
70 and *Nematocida* sp1[10].

71 Unfortunately why introns were lost, especially massively lost in some eukaryotes,

72 remains obscure despite extensive research [11, 12]. Obviously, the reverse question
73 -- how introns, especially the few ones in intron-poor eukaryotes, can be retained
74 under the pressure of loss is equally important, but it was rarely carefully studied.
75 Although some authors thought that the reason for the retention of introns in genomes
76 is likely due to the essential functions of these introns [13-15], this ‘functional
77 constraint’ scenario – “only the introns with important functions can get rid of the fate
78 of being lost” lacks evidence that the lost introns are all useless or less useful than the
79 retained ones in any eukaryotes, and moreover, actually the functions of introns are
80 still far from being well understood [16]. Therefore, the investigation of the
81 evolutionary retention of intron might be helpful not only to answering the question
82 about intron loss but also to understanding the function and evolution of introns.

83 *Giardia lamblia* is a parasitic protozoan belonging to Diplomonadida (Excavata). It
84 has a very minimalistic genome, compact in structure and content [17], and only eight
85 spliceosomal introns were found in its genome [17-22]. Thus it can be speculated that
86 this organism must have undergone massive intron loss, with very few left in the
87 genome. Therefore, this organism may provide an excellent opportunity for exploring
88 how the few introns were retained. In the present work, by investigating the
89 intron-containing genes and the few introns of *Giardia*, besides finding the
90 distribution patterns that can reflect the dynamic evolution of intron in *Giardia*, we
91 observed and experimentally confirmed an interesting phenomenon that
92 sense-antisense (SAS) gene overlaps appear in the areas of some introns, and thus
93 “overlap constraint” might be at least one of the causes for preventing introns from

94 being lost, though it is uncertain whether the other retained introns also overlap with
95 any unknown genomic functional elements yet. The implications of these findings for
96 intron evolution and function are discussed.

97

98 **Results**

99 **Characteristics of the intron-containing genes and their introns in *Giardia*** 100 **genomes**

101 In the genome database of *Giardia*, GiardiaDB, four of the eight *Giardia*
102 intron-containing genes are annotated to code proteins with sequence similarity to
103 known proteins, and the other four to code hypothetical proteins. Our investigation
104 (mainly by sequence comparative analysis) indicated that: 1) the former four
105 intron-containing genes are all common eukaryotic-conserved genes, which are most
106 likely vertically inherited from the LECA and thus are very ‘ancient’, while the latter
107 four are all *Giardia*-specific genes (not found in other eukaryotes including other
108 Excavata species), which thus most likely emerged after the divergence of *Giardia*
109 from other Excavata and thus are ‘relatively new’ genes compared with the ‘ancient’
110 ones above; 2) the introns in the three (GL50803-15604, GL50803-15124,
111 GL50803-17244) of the four ancient genes are eukaryotic-conserved (Additional file
112 1), and thus they are ‘ancient’ introns in ancient genes, while the intron in the other
113 one ancient gene (GL50803-27266) is a *Giardia*-specific intron (not found in other
114 eukaryotes including other Excavata species), and thus this intron most likely
115 emerged after the divergence of *Giardia* from other Excavata and is a ‘relatively new’

116 intron in an ancient gene; 3) all four *Giardia*-specific ('relatively new')
117 intron-containing genes (GL50803-37070, GL50803-35332, GL50803-15525 and
118 GL50803- 86945), which account for only about 0.6 percent of all the ~700
119 *Giardia*-specific protein-coding genes in the genome[17], each contain an
120 *Giardia*-specific intron (not found in other eukaryotes including other Excavata
121 species), and thus the four introns all are 'relatively new' introns in 'relatively new'
122 genes (Table 1).

123 These observations suggest that: 1) while *Giardia* massively lost its introns, new
124 introns also arose in both the ancient and relatively newly evolved genes; 2) the
125 pressure of intron loss seems to constantly exist in the whole evolutionary process of
126 *Giardia*, but only a few of both the ancient and newly-emerged introns have been
127 retained in the genome.

128 To find the reason why these few introns can be retained in *Giardia* genome under
129 strong pressure of loss, we investigated the characteristics of these introns in many
130 aspects. It had been shown that seven of the eight *Giardia* introns are bounded by
131 canonical GT-AG splice signals, only one, the [2Fe-2S] ferredoxin (GL50803-27266)
132 intron, has a noncanonical splice signal CT-AG [19]. The sizes of the eight introns are
133 all small (most of them are less than 40 bp long and are not the multiple of three) and
134 do not have any conserved sequence motifs. Our further analysis predicted no special
135 secondary structures that would be able to form in these introns. Besides, our survey
136 also showed that there were not any reports about alternative splicing of the two
137 introns in genes GL50803-15525 and GL50803-86945 [18, 22]. Taken together, these

138 results suggest that the retention of the few introns seems to be neither due to the
139 structural features nor necessarily due to the functional importance of these introns.

140 Interestingly, on the complementary strands, we found that two intron-containing
141 genes, GL50803-17244 (ribosomal protein L7a gene) and GL50803-37070 (a
142 “hypothetical protein” gene), each have an antisense gene, GL50803-20429 and
143 GL50803-28204, which are just annotated as “hypothetical protein” and “unspecified
144 product” in the genome database, respectively. That is, the two intron-containing
145 genes and their antisense gene form SAS gene pairs. We thought this phenomenon
146 might be related to the intron retention. Nevertheless, the two anti-sense genes need to
147 be further verified, and the details of the overlaps with their sense genes also need to
148 be analyzed in detail.

149 **Verification of the antisense genes**

150 The strand-specific RT-PCR of the two antisense genes, GL50803-28204 and
151 GL50803-20429, generated two products with the expected lengths of 172 and 288 bp,
152 respectively. The sequencing further confirmed that the two products just seem to be
153 transcribed from the opposite direction of the two sense (intron-containing) genes,
154 GL50803-37070 and GL50803-17244, respectively, and the two introns are just
155 located within the two overlapping regions of the two SAS gene pairs, respectively
156 (Figure 1).

157 The RACE of the two antisense genes generated a 1232 bp product for
158 GL50803-28204 and a 1177 bp product for GL50803-20429. After sequencing and
159 comparing with their corresponding genomic DNA sequences in the GiardiaDB

160 database, we found the two antisense genes contain no introns, especially in the
161 regions corresponding to the two introns of the sense genes. But the software ORF
162 Finder predicted that the largest ORFs of the two antisense genes are only 264 bp for
163 GL50803-28204 and 363 bp for GL50803-20429, and moreover, no proteins
164 homologous to the putative proteins coded by the two largest ORFs could be found in
165 other organisms including *Giardia's* close relative *Spiroucleus* in GenBank.
166 Therefore, the two antisense genes have transcriptional activity and are most likely
167 *Giardia*-specific non-coding genes.

168 **Discussion**

169 To investigate the reasons for the evolutionary retention of the few introns in the
170 eukaryotes having undergone massive intron loss, we chose the extremely intron-poor
171 eukaryote *Giardia* as the study object. When investigating the characteristics
172 (including distribution patterns) of the eight intron-containing genes and their
173 corresponding introns in *Giardia* genome, we found that in spite of the massive loss
174 of introns, intron gain also occurred in both *Giardia's* ancient and relatively
175 newly-evolved genes; it turns out that the selective pressure of intron loss seems to
176 constantly exist in the whole evolutionary course of *Giardia*, but a few of both the
177 ancient and the newly-emerged introns have been retained in modern *Giardia*.

178 To explore how these few introns can be retained under the constant strong pressure
179 of loss. First, we investigated whether there exist some features in these few
180 intron-containing genes and the introns probably responsible for the retention, but
181 failed to find any special regularities or unusualnesses in many aspects such as

182 splicing signal, intron size and secondary structure, and alterative splicing. This
183 suggests that there are neither any special structural features nor necessarily any
184 functional importance of introns responsible for the intron retention. This is consistent
185 with the fact that so many introns, at least part of which definitely possesses important
186 functions, have been lost in intron-poor eukaryotes like *Giardia*. Thus, the reasons for
187 the retention might lie outside the intron-containing genes and the introns themselves.

188 Interestingly, we noticed that on the complementary strands of two of the eight
189 intron-containing genes, GL50803-37070 and GL50803-17244, there exist
190 correspondingly two anti-sense genes, GL50803-28204 and GL50803-20429, though
191 they are just annotated as “product unknown” in GiardiaDB. By strand-specific
192 RT-PCR, RACE and sequencing, we got the transcripts and sequences of the two
193 genes and found they both have no introns. Thus the two anti-sense genes have been
194 verified to be really genes that are actively transcribed. And actually the anti-sense
195 gene GL50803-20429 has been reported to be a mRNA gene being expressed during
196 excystation and encystation, and in trophozoites but not cysts [23]. As for the other
197 anti-sense gene GL50803-28204, it has a quite short putative ORF but has no
198 homologs in other organisms including *Spironucleus*. Although the corresponding
199 DNA sequence regions in the four other *G. lamblia* isolates (DH, P15, GS and GS-B)
200 with genomic data exhibit significant similarities ($\geq 83.9\%$ similarity) to those of the
201 two anti-sense genes, there are not any annotations and transcriptome information
202 about those regions. Besides, the total RNA were processed using Poly(A)
203 Polymerase to add a poly(A) tail at the 3'ends before we performed rapid

204 amplification of their cDNA 3'ends, thus from the experiment we still did not know
205 whether the transcripts of GL50803-28204 are polyadenylated or not, namely, mRNA
206 or not. Therefore, we can only conjecture that the two anti-sense genes are either
207 non-protein-coding genes or *Giardia*-specific protein-coding genes. Considering that
208 many identified non-coding RNAs in *Giardia* overlap with protein-coding genes on
209 the antisense strands [24], the antisense gene might also be noncoding RNA gene. But
210 there is still no tangible evidence for what the two genes are despite our lots of
211 experimental efforts (not shown) to identify them. Whatever the antisense genes code
212 for, our work showed that they are functional genes and form SAS gene overlap with
213 their intron-containing sense genes, and that the introns just reside in the overlap
214 regions. Considering that the gene sequence mutation (especially deletion) cannot
215 occur randomly, the antisense genes must have imposed the restriction of variation
216 (especially deletion) on the introns of the sense genes in the overlapping areas, and
217 thus such a kind of SAS gene overlap must have prevented the introns from being
218 lost.

219 As for the other six introns, we did not find any ORFs on their corresponding
220 complementary strands. Although we also experimentally examined whether their
221 complementary strands (especially the areas overlapping with these introns) are
222 transcribed, no transcripts were found (Additional file 2). Nevertheless, it is uncertain
223 whether the corresponding complementary strands of these introns are resided by
224 some unknown genomic functional elements which are not transcribed at all. If this is
225 true, these introns are also retained by the same cause as the former two ones.

226 Certainly, it is also possible that the six introns are retained by other unknown reasons.
227 We have also analyzed the intron regions of many intron-poor eukaryotes including
228 *Microsporidia*, *Trichomonas*, *Spiroplasma*, but unfortunately did not find such
229 sense-antisense gene overlaps as in *Giardia* (data not shown). But considering some
230 genes overlap with the UTRs of the adjacent genes (as found in *Antonospora locustae*
231 and *Encephalitozoon cuniculi* [25]), we can not obtain the UTRs information of those
232 genes in current database, thus many of the overlaps might be able to be found out.
233 More importantly, some introns might overlap with unknown genomic functional
234 elements including non-coding and non-transcribed ones, since there are so huge
235 remaining component of eukaryotic genomes, much of which was traditionally
236 regarded as "junk" and is still undetermined. This might be the important cause for
237 that few SAS gene overlaps in intron regions that can be identified at present.

238 Theoretically, overlapping with any genomic functional elements on either the
239 same strand or the complementary one (namely, either same-strand overlapping or
240 different-strand overlapping) can result in intron retention, as long as the introns are
241 just in the overlapping areas. Therefore, since such overlapping structures are widely
242 distributed in eukaryotes [26], it can be expected that quite a number of introns in
243 diverse eukaryotes may also be retained due to this kind of "overlap constraint". We
244 believe that more and more examples might be able to be found in diverse eukaryotes
245 in the future. Conversely, such an intron retention phenomenon probably can provide
246 a valuable clue to find new genomic functional elements – in the overlapping area
247 with introns.

248 **Conclusions**

249 In conclusion, by investigating the extremely intron-poor eukaryote *Giardia*, we
250 have revealed some interesting findings about the dynamic evolution of introns in the
251 intron-poor eukaryotes: the pressure of intron loss may constantly exist in these
252 eukaryotes, but new introns can still arise either in ancient genes or new-evolved
253 genes, but only a few introns can be retained in the genome; the retention of the few
254 introns is not caused by special features or functional constraint of the introns
255 themselves but due to the reasons outside of the introns, and “overlap constraint”
256 imposed by other genomic functional elements overlapping with the introns is at least
257 an important one of the causes. First, our findings not only support the “intron-rich
258 ancestor” theory, but also can explain why few eukaryotes were found to be
259 completely intronless. Second, our finding may conversely provide a clue to find new
260 genomic functional elements (which was probably traditionally regarded as “junk”
261 and is still undetermined) in such kinds of overlap regions. Most importantly, our
262 work implicates that “functional constraint” of introns is not necessarily directly
263 associated with intron loss and gain, or that the real functions or the way of
264 functioning of introns are probably still outside of our current knowledge. Therefore
265 our work may be able to shed some new lights on the research of evolution and
266 function of introns and genomes.

267 **Methods**

268 **Database and bioinformatics methods**

269 The template sequences for designing primers of *Giardia* genes were downloaded

270 from GiardiaDB (December 1, 2017 release) [27]. The software ORF Finder were
271 used to predict the ORFs of the RACE products (see blow), then the predicted ORFs
272 were used as queries to search their homologs with Blastp against the NCBI
273 non-redundant protein sequences (nr) database. The program RNAfold web server
274 was used to analyze the secondary structure of introns. The sequences of the four
275 genes and their coding proteins, 2Fe-2S ferredoxin, 26S proteasome non-ATPase
276 regulatory subunit 4, Dynein light chain, and Ribosomal Protein L7A from other
277 organisms were identified and collected by Blastp searching against GenBank with
278 *Giardia*'s corresponding sequences as queries. Protein alignments were generated
279 with ClustalX 2 applying default alignment parameters. The introns in the genes were
280 determined by comparing cDNA and gene sequences. The other four
281 intron-containing genes with annotated as hypothetical protein also identified by
282 sequence comparative analysis to determine whether they are *Giardia*-specific or not.

283 ***Giardia* cultures**

284 The cell line of WB isolate (assemblage A), namely WB clone C6 (ATCC 50803),
285 was used in the study. Its cultures were routinely grown in filter-sterilized TYI-S-33
286 medium supplemented with bovine bile in glass screw cap tubes at 37 °C and were
287 sub-cultured every 3 to 4 days.

288 **RNA extraction**

289 *Giardia* total RNA was extracted and treated to remove any contaminated genomic
290 DNA by RNAPrep Pure Cell/Bacteria Kit (TIANGEN) using about 5×10^6 *Giardia*
291 trophozoites that were harvested by ice-slush incubation and centrifuged at 6000g for

292 5 min according to the manufacturer's instructions. The quality and quantity of the
293 RNA preparation were assessed with agarose gel electrophoresis and the absorption at
294 260 and 280nm.

295 **Strand-specific RT-PCR**

296 First-strand cDNAs of the two antisense genes, GL50803-28204 and
297 GL50803-20429, were synthesized from 500ng DNase-I treated total RNA per
298 reaction at 54°C 30 min, 99°C 5 min, and 5°C 5 min with the specific antisense
299 primer 4C and 9C respectively instead of the reverse primers in the kit, and then
300 amplified with specific primers pairs of 4A/4S and 9A/9S by using RNA PCR Kit
301 (AMV) Ver.3.0 (Takara, Japan). The PCR conditions were as follows: 94°C for 30 s,
302 followed by 30 cycles of 94°C for 30 s, 55°C for 30 s and 72°C for 45s. The PCR
303 products were purified using the Wizard SV Gel and PCR Clean-Up System kit
304 (Qiagen, Germany), and cloned into pMD-19T vectors from 35 ng purified PCR
305 products using TaKaRa pMD-19T Vector Cloning Kit (TaKaRa, Japan) according to
306 the manufacturer's instructions. Then, the ligation products were transformed into
307 DH5 α Chemically Competent *E. coli*. Colony PCR with vector-specific primers
308 provided in the kit was adopted to select colonies. These selected colonies were
309 sequenced using vector-specific forward and reverse primers by Sangon Biology
310 Company (Shanghai, China).

311 **Rapid amplification of cDNA ends**

312 The total RNA were processed using Poly(A) Polymerase(TaKaRa, Japan) to add a
313 poly(A) tail at the 3'ends of the RNA before performing rapid amplification of their

314 cDNA 3'ends. We experimentally determined the 3'ends by using nested PCR primer
315 (3R4O/3R4I and 3R9O/3R9I) according to the RNA PCR Kit (AMV) Ver.3.0 (Takara,
316 Japan). 5'-RACE was performed by using a SMARTer RACE 5'/3'Kit (TaKaRa, Japan)
317 with 500ng total RNA as the template and the gene-specific 5'-RACE primers 5R4
318 and 5R9 for the two antisense genes, GL50803-28204 and GL50803-20429,
319 according to the manufacturer's instructions. Both the 3'-RACE and 5'-RACE primers
320 (Additional file 3) were designed based on the transcripts from the Strand-specific
321 RT-PCR. The RACE-PCR products were analyzed by agarose gel electrophoresis and
322 sequenced as described above.

323 **List of abbreviations**

324 LECA: last eukaryotic common ancestor.

325 SAS: sense-antisense.

326 RACE: rapid amplification of cDNA ends.

327 **Declarations**

328 **Ethics approval and consent to participate**

329 Not applicable

330 **Consent for publication**

331 Not applicable

332 **Availability of data and materials**

333 All data generated or analysed during this study are included in this published
334 article and its supplementary information files.

335 **Competing interests**

336 The authors declare that they have no competing interests.

337 **Funding**

338 This work was supported by the National Natural Science Foundation of China
339 (NSFC) (grant numbers 31572256, 31772452, 31401972 and 31401973) and the
340 Natural Science Foundation of Yunnan Province (grant number 2015FB181)

341 **Authors' contributions**

342 J.W. designed and supervised this study. M.X performed genetic characterization
343 work. M.X., B.C., Q.Y., J.S., and Z.L. analyzed the data. M.X. and J.W. wrote the
344 manuscript. All authors read and approved the final manuscript.

345 **Acknowledgements**

346 Not applicable

347

348 **References**

- 349 1. Csuros M, Rogozin IB, Koonin EV: A detailed history of intron-rich
350 eukaryotic ancestors inferred from a global survey of 100 complete genomes.
351 *PLoS Comput Biol* 2011, 7(9):e1002150.
- 352 2. Roy SW: Intron-rich ancestors. *Trends Genet* 2006, 22(9):468-471.
- 353 3. Pombert JF, Selman M, Burki F, Bardell FT, Farinelli L, Solter LF, Whitman
354 DW, Weiss LM, Corradi N, Keeling PJ: Gain and loss of multiple functionally
355 related, horizontally transferred genes in the reduced genomes of two
356 microsporidian parasites. *Proceedings of the National Academy of Sciences of*

- 357 *the United States of America* 2012, 109(31):12638-12643.
- 358 4. Stark MR, Dunn EA, Dunn WS, Grisdale CJ, Daniele AR, Halstead MR, Fast
359 NM, Rader SD: Dramatically reduced spliceosome in *Cyanidioschyzon*
360 *merolae*. *Proc Natl Acad Sci U S A* 2015, 112(11):E1191-1200.
- 361 5. Shoguchi E, Shinzato C, Kawashima T, Gyoja F, Mungpakdee S, Koyanagi R,
362 Takeuchi T, Hisata K, Tanaka M, Fujiwara M *et al*: Draft assembly of the
363 *Symbiodinium minutum* nuclear genome reveals dinoflagellate gene structure.
364 *Current biology : CB* 2013, 23(15):1399-1408.
- 365 6. Roy SW, Gilbert W: The evolution of spliceosomal introns: patterns, puzzles
366 and progress. *Nat Rev Genet* 2006, 7(3):211-221.
- 367 7. Rodriguez-Trelles F, Tarrío R, Ayala FJ: Origins and evolution of
368 spliceosomal introns. *Annu Rev Genet* 2006, 40:47-76.
- 369 8. Irimia M, Roy SW: Origin of spliceosomal introns and alternative splicing.
370 *Cold Spring Harb Perspect Biol* 2014, 6(6).
- 371 9. Koonin EV: Intron-dominated genomes of early ancestors of eukaryotes. *J*
372 *Hered* 2009, 100(5):618-623.
- 373 10. Cuomo CA, Desjardins CA, Bakowski MA, Goldberg J, Ma AT, Becnel JJ,
374 Didier ES, Fan L, Heiman DI, Levin JZ *et al*: Microsporidian genome analysis
375 reveals evolutionary strategies for obligate intracellular growth. *Genome Res*
376 2012, 22(12):2478-2488.
- 377 11. Cohen NE, Shen R, Carmel L: The role of reverse transcriptase in intron gain
378 and loss mechanisms. *Molecular biology and evolution* 2012, 29(1):179-186.

- 379 12. Mitrovich QM, Tuch BB, De La Vega FM, Guthrie C, Johnson AD: Evolution
380 of yeast noncoding RNAs reveals an alternative mechanism for widespread
381 intron loss. *Science* 2010, 330(6005):838-841.
- 382 13. Oswald A, Oates AC: Control of endogenous gene expression timing by
383 introns. *Genome biology* 2011, 12(3):107.
- 384 14. Rearick D, Prakash A, McSweeney A, Shepard SS, Fedorova L, Fedorov A:
385 Critical association of ncRNA with introns. *Nucleic Acids Res* 2011,
386 39(6):2357-2366.
- 387 15. Chorev M, Carmel L: The function of introns. *Front Genet* 2012, 3:55.
- 388 16. Carmel L, Wolf YI, Rogozin IB, Koonin EV: Three distinct modes of intron
389 dynamics in the evolution of eukaryotes. *Genome Res* 2007, 17(7):1034-1044.
- 390 17. Morrison HG, McArthur AG, Gillin FD, Aley SB, Adam RD, Olsen GJ, Best
391 AA, Cande WZ, Chen F, Cipriano MJ *et al*: Genomic minimalism in the early
392 diverging intestinal parasite *Giardia lamblia*. *Science* 2007,
393 317(5846):1921-1926.
- 394 18. Kamikawa R, Inagaki Y, Hashimoto T: Secondary loss of a cis-spliced intron
395 during the divergence of *Giardia intestinalis* assemblages. *BMC Res Notes*
396 2014, 7:413.
- 397 19. Nixon JE, Wang A, Morrison HG, McArthur AG, Sogin ML, Loftus BJ,
398 Samuelson J: A spliceosomal intron in *Giardia lamblia*. *Proc Natl Acad Sci U*
399 *S A* 2002, 99(6):3701-3705.
- 400 20. Russell AG, Shutt TE, Watkins RF, Gray MW: An ancient spliceosomal intron

- 401 in the ribosomal protein L7a gene (Rpl7a) of *Giardia lamblia*. *BMC Evol Biol*
402 2005, 5:45.
- 403 21. Roy SW, Hudson AJ, Joseph J, Yee J, Russell AG: Numerous fragmented
404 spliceosomal introns, AT-AC splicing, and an unusual dynein gene expression
405 pathway in *Giardia lamblia*. *Mol Biol Evol* 2012, 29(1):43-49.
- 406 22. Franzen O, Jerlstrom-Hultqvist J, Einarsson E, Ankarklev J, Ferella M,
407 Andersson B, Svard SG: Transcriptome profiling of *Giardia intestinalis* using
408 strand-specific RNA-seq. *PLoS Comput Biol* 2013, 9(3):e1003000.
- 409 23. Birkeland SR, Preheim SP, Davids BJ, Cipriano MJ, Palm D, Reiner DS,
410 Svard SG, Gillin FD, McArthur AG: Transcriptome analyses of the *Giardia*
411 *lamblia* life cycle. *Molecular and biochemical parasitology* 2010,
412 174(1):62-65.
- 413 24. Chen XS, White WT, Collins LJ, Penny D: Computational identification of
414 four spliceosomal snRNAs from the deep-branching eukaryote *Giardia*
415 *intestinalis*. *PLoS One* 2008, 3(8):e3106.
- 416 25. Corradi N, Gangaeva A, Keeling PJ: Comparative profiling of overlapping
417 transcription in the compacted genomes of microsporidia *Antonospora*
418 *locustae* and *Encephalitozoon cuniculi*. *Genomics* 2008, 91(4):388-393.
- 419 26. Kumar A: An overview of nested genes in eukaryotic genomes. *Eukaryot Cell*
420 2009, 8(9):1321-1329.
- 421 27. Shabalina SA, Ogurtsov AY, Kondrashov VA, Kondrashov AS: Selective
422 constraint in intergenic regions of human and mouse genomes. *Trends In*

423 *Genetics* 2001, 17(7):373-376.

424

425 **Table 1. The integrated information of the eight spliceosome-mediated cis-spliced**
 426 **introns and their host genes and complementary strands.**

Intron-containing gene			Intron		Complementary strand
Gene ID	Product	Age	Age	Size (bp)	
GL50803-27266	2Fe-2S ferredoxin		relatively new	35	No ORF, no transcripts detected
GL50803-15604	26S proteasome non-ATPase regulatory subunit 4	ancient	ancient	29	No ORF, no transcripts detected
GL50803-15124	Dynein light chain			32	No ORF, no transcripts detected
GL50803-17244	Ribosomal Protein L7A			109	Antisense gene with transcripts
GL50803-37070	Hypothetical protein			41	Antisense gene with transcripts
GL50803-35332	Hypothetical protein	relatively new	relatively new	220	No ORF, no transcripts detected
GL50803-15525	Hypothetical protein			33	No ORF, no transcripts detected
GL50803-86945	Hypothetical protein			36	No ORF, no transcripts detected

427

428 Figure Captions

429 **Fig 1. Results of strand-specific RT-PCR and sequencing of the two antisense**
 430 **genes, and the schematic diagram of this two SAS gene pairs.**

431 **A.** Lane 1, Strand-specific RT-PCR product of the GL50803-20429; Lane 4, the Strand-specific
 432 RT-PCR product of the GL50803-28204; Lane 2 and lane 3, negative controls (with no RTase)
 433 corresponding to lane 1 and lane 4, respectively; M, molecular markers. **B.** Nucleotide sequence of
 434 GL50803-28204 gene acquired by Strand-specific RT-PCR and sequencing. The locations of the
 435 primers are underlined. **C.** Nucleotide sequence of GL50803-20429 gene acquired by Strand-
 436 specific RT-PCR and sequencing. **D.** Schematic diagram of the two SAS gene pairs. The sequence
 437 lengths of GL50803-28204 and GL50803-20429 are according to the Strand-specific RT-PCR
 438 products, and the lengths of GL50803-37070 and GL50803-17244 are based on the GiardiaDB

439 database. Arrow represents the orientation of transcription; and the dashed box and solid lines

440 represent introns and exons, respectively.

441

442 **Additional files**

443 **Additional file 1:** The conservative analysis of Giardia's introns among diverse

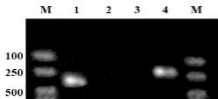
444 eukaryotes

445 **Additional file 2:** Results of strand-specific RT-PCR of the complementary areas of

446 the other six introns of *G. lamblia*.

447 **Additional file 3:** The primers designed for the strand-specific RT-PCR and RACE of

448 the complementary areas of the eight introns of *G. lamblia*.

A**B****GL50803-28204:**

ATGCCGATAAAGATAAAGGACGCGCGAAAG
 ACGCTCTTGGGGGAGATGCGTGCAACGGGA
 GCGTGCTTCTCTGTCTTTCTGCGTGTTAGACG
 TCGGGTGGAACGAGCAATGAAACATACGCCA
 TCTCTGGATTTTTTTCTTTGGCGGCTTTATGG
CAGGCTACCGATTACC

C**GL50803-20429:**

AAAGGTGGGCTTGTCCTCTGGGTTACAGTCCGGTGAAGCAGACGCTGGTTCGTCTTCTTC
 AGATGGACGAGCTTGCCAGATCGCCCTTA GTGCGAACGATGGCGTACGGGACGCC
 ATCTTGTGACAGAGTGTGGGAAGCCAAAGTACGAGCTGTGGGTGAGTTGTCAGGTGA
 ACAGCGAAGTCCACCCGCTGACAACACACAACCCGCAATCAGAGGTGTGTGCGGTGTCAG
 CGGACGGCTCCTCGCGCATAAGAACATACTTCAAGGGGGTTCGACATCATTGCAATCA

D**GL50803-37070:****GL50803-28204:****GL50803-17244:****GL50803-20429:**