

1 **A free-living protist that lacks canonical eukaryotic DNA replication and segregation systems**

2 Dayana E. Salas-Leiva<sup>1</sup>, Eelco C. Tromer<sup>2,3</sup>, Bruce A. Curtis<sup>1</sup>, Jon Jerlström-Hultqvist<sup>1</sup>, Martin

3 Kolisko<sup>4</sup>, Zhenzhen Yi<sup>5</sup>, Joan S. Salas-Leiva<sup>6</sup>, Lucie Gallot-Lavallée<sup>1</sup>, Geert J. P. L. Kops<sup>3</sup>, John M.

4 Archibald<sup>1</sup>, Alastair G. B. Simpson<sup>7</sup> and Andrew J. Roger<sup>1\*</sup>

5 <sup>1</sup>Centre for Comparative Genomics and Evolutionary Bioinformatics (CGEB), Department of

6 Biochemistry and Molecular Biology, Dalhousie University, Halifax, NS, Canada, B3H 4R2

7 <sup>2</sup>Department of Biochemistry, University of Cambridge, Cambridge, United Kingdom

8 <sup>3</sup>Oncode Institute, Hubrecht Institute – KNAW (Royal Netherlands Academy of Arts and Sciences)

9 and University Medical Centre Utrecht, Utrecht, The Netherlands

10 <sup>4</sup>Institute of Parasitology Biology Centre, Czech Acad. Sci, České Budějovice, Czech Republic

11 <sup>5</sup>Guangzhou Key Laboratory of Subtropical Biodiversity and Biomonitoring, School of Life Science,

12 South China Normal University, Guangzhou 510631, China

13 <sup>6</sup>CONACyT-Centro de Investigación en Materiales Avanzados, Departamento de medio ambiente y

14 energía, Miguel de Cervantes 120, Complejo Industrial Chihuahua, 31136 Chihuahua, Chih., México

15 <sup>7</sup>Centre for Comparative Genomics and Evolutionary Bioinformatics (CGEB), Department of

16 Biology, Dalhousie University, Halifax, NS, Canada, B3H 4R2

17 \*corresponding author: Andrew.Roger@dal.ca

18 D.E.S-L ORCID iD: 0000-0003-2356-3351

19 E.C.T. ORCID iD: 0000-0003-3540-7727

20 B.A.C. ORCID iD: 0000-0001-6729-2173

21 J.J-H. ORCID iD: 0000-0002-7992-7970

22 M.K. ORCID iD: 0000-0003-0600-1867

- 23 Z.Y. ORCID iD: 0000-0002-0693-9989
- 24 J.S-L. ORCID iD: 0000-0002-4141-140X
- 25 L.G-L. ORCID iD:
- 26 G.J.P.L.K. ORCID iD: 0000-0003-3555-5295
- 27 J.M.A. ORCID iD: 0000-0001-7255-780X
- 28 A.G.B.S. ORCID iD:0000-0002-4133-1709
- 29 A.J.R. ORCID iD: 0000-0003-1370-9820
- 30

## 31 **Abstract**

32 Cells must replicate and segregate their DNA with precision. In eukaryotes, these processes are part  
33 of a regulated cell-cycle that begins at S-phase with the replication of DNA and ends after M-phase.  
34 Previous studies showed that these processes were present in the last eukaryotic common ancestor  
35 and the core parts of their molecular systems are conserved across eukaryotic diversity. However,  
36 some unicellular parasites, such as the metamonad *Giardia intestinalis*, have secondarily lost  
37 components of the DNA processing and segregation apparatuses. To clarify the evolutionary history  
38 of these systems in these unusual eukaryotes, we generated a high-quality draft genome assembly for  
39 the free-living metamonad *Carpediemonas membranifera* and carried out a comparative genomics  
40 analysis. We found that parasitic and free-living metamonads harbor a conspicuously incomplete set  
41 of canonical proteins for processing and segregating DNA. Unexpectedly, *Carpediemonas* species  
42 are further streamlined, completely lacking the origin recognition complex, Cdc6 and other replisome  
43 components, most structural kinetochore subunits including the Ndc80 complex, as well as several  
44 canonical cell-cycle checkpoint proteins. *Carpediemonas* is the first eukaryote known to have lost  
45 this large suite of conserved complexes, suggesting that it has a highly unusual cell cycle and that  
46 unlike any other known eukaryote, it must rely on a novel set of mechanisms to carry out these  
47 fundamental processes.

48

49 DNA replication, repair and segregation are critically important and conserved processes in  
50 eukaryotes that have been intensively studied in model organisms<sup>1</sup>. The initial step of DNA replication  
51 is accomplished by the replisome, a set of highly conserved proteins that is tightly regulated to  
52 minimize mutations<sup>2</sup>. The replisome relies on the interactions between cis-acting DNA sequences and  
53 trans-acting factors that serve to separate the template and promote RNA-primed DNA synthesis. This

54 occurs by the orderly assembly of the origin recognition (ORC), the pre-replicative (pre-RC), pre-  
55 initiation (pre-IC) and replication progression (RPC) complexes<sup>3-6</sup>. The synthesis of DNA usually  
56 encounters disruptive obstacles as replication proceeds and can be rescued either through template  
57 switching via trans-lesion or recombination-dependent synthesis. Trans-lesion synthesis uses  
58 replicative and non-replicative DNA polymerases to by-pass the lesion through multiple strategies that  
59 incorporate nucleotides opposite to it<sup>7</sup>, while recombination-dependent synthesis uses non-homologous  
60 or homologous templates for repair (reviewed in refs.<sup>8,9</sup>). Recombination-dependent synthesis occurs  
61 in response to single- or double-strand DNA breakage<sup>8,10,11</sup>. Other repair mechanisms occur throughout  
62 the cell cycle, fixing single-strand issues through base excision, nucleotide excision or mismatch  
63 repair, but they may also be employed during replication depending on the source of the damage. All  
64 of the repair processes are overseen by multiple regulation checkpoints that permit or stall DNA  
65 replication and the progression of the cell cycle. During M-phase the replicated DNA has to form  
66 attachments with the microtubule-based spindle apparatus via kinetochores, large multi-subunit  
67 complexes built upon centromeric chromatin<sup>12</sup>. Unattached kinetochores catalyse the formation of a  
68 soluble inhibitor of the cell cycle, preventing precocious chromosome segregation, a phenomenon  
69 known as the spindle assembly checkpoint (SAC)<sup>12</sup>. Failure to pass any of these checkpoints (*e.g.*,  
70 G1/S, S, G2/M and SAC checkpoints reviewed in refs.<sup>12-14</sup>) leads to genome instability and may result  
71 in cell death.

72 To investigate the diversity of DNA replication, repair, and segregation processes, we  
73 conducted a eukaryote-wide comparative genomics analysis with a special focus on metamonads, a  
74 major protist lineage comprised of parasitic and free-living anaerobes. Parasitic metamonads such as  
75 *Giardia intestinalis* and *Trichomonas vaginalis* are extremely divergent from model system  
76 eukaryotes, exhibit a diversity of cell division mechanisms (*e.g.*, closed/semi-open mitosis), possess

77 metabolically reduced mitosomes or hydrogenosomes instead of mitochondria, and lack several  
78 canonical eukaryotic features on the molecular and genomic-level<sup>15-17</sup>. Indeed, recent studies show  
79 that metamonad parasites have secondarily lost parts of the ancestral DNA replication and  
80 segregation apparatuses<sup>18,19</sup>. Furthermore, comparisons of metamonad proteins sequences to those of  
81 other eukaryotes reveal that they are often highly divergent compared to other eukaryotic homologs,  
82 indicating a high substitution rate in these organisms that is suggestive of error-prone replication  
83 and/or DNA repair<sup>20,21</sup>. Yet, it is unclear whether the divergent nature of proteins studied in  
84 metamonads is the result from the host-associated lifestyle or is a more ancient feature of  
85 Metamonada. To increase the representation of free-living metamonads in our analyses, we have  
86 generated a high-quality draft genome assembly of *Carpediemonas membranifera*, a flagellate  
87 isolated from hypoxic marine sediments<sup>22</sup>. Our analyses of genomes from across the tree of  
88 eukaryotes show that many systems for DNA replication, repair, segregation, and cell cycle control  
89 are ancestral to eukaryotes and highly conserved. However, metamonads have secondarily lost an  
90 extraordinarily large number of components. Most remarkably, the free-living *Carpediemonas*  
91 species have been drastically reduced further, having lost a large set of key proteins from the  
92 replisome and cell-cycle checkpoints (*i.e.*, including several from the kinetochore and repair  
93 pathways). We propose a hypothesis of how DNA replication may be achieved in this organism.

94

## 95 **Results**

### 96 **The *C. membranifera* genome assembly is complete.**

97 Our assembly for *C. membranifera* (a member of the Fornicata clade within metamonads, **Fig. 1**)  
98 is highly contiguous (**Table 1**) and has deep read coverage (*i.e.*, median coverage of 150× with short  
99 reads and 83× with long-reads), with an estimated genome completeness of 99.27% based on the

100 Mercury<sup>23</sup> method. 97.6% of transcripts mapped to the genome along their full length with an identity  
101 of  $\geq 95\%$  while a further 2.04% mapped with an identity between 90 - 95%. The high contiguity of the  
102 assembly is underscored by the large number of transcripts mapped to single contigs (90.2%), and  
103 since the proteins encoded by transcripts were consistently found in the predicted proteome, the latter  
104 is also considered to be of high quality. We also conducted BUSCO analyses, with the foreknowledge  
105 that genomic streamlining typical in Metamonada has led to the loss of many conserved proteins<sup>17,24,25</sup>.  
106 Our analyses show that previously completed metamonad genomes only encoded between 60% to 91%  
107 BUSCO proteins, while *C. membranifera* exhibits a relatively high 89% (**Table 1, Supplementary**  
108 **Information**). In any case, our coverage estimates for the *C. membranifera* genome for short and long  
109 read sequencing technologies are substantially greater than those found to be sufficient to capture  
110 genic regions that otherwise would have been missed (*i.e.*, coverage  $>52\times$  for long reads and  $>60\times$  for  
111 short paired-end reads, see ref.<sup>26</sup>). All these various data indicate that the draft genome of *C.*  
112 *membranifera* is nearly complete; if any genomic regions are missing, they are likely confined to  
113 difficult-to-sequence highly repetitive regions such as telomeres and centromeres.

114

### 115 **Extreme streamlining of the DNA replication apparatus in metamonads**

116 The first step in the replication of DNA is the assembly of ORC which serves to nucleate the pre-RC  
117 formation. The initiator protein Orc1 first binds an origin of replication, followed by the recruitment  
118 of Orc 2-6 proteins, which associate with chromatin<sup>27</sup>. As the cell transitions to G1 phase, the  
119 initiator Cdc6 binds to the ORC, forming a checkpoint control<sup>28</sup>. Cdt1 then joins Cdc6, promoting the  
120 loading of the replicative helicase MCM forming the pre-RC, a complex that remains inactive until  
121 the onset of S-phase when the ‘firing’ factors are recruited to convert the pre-RC into the pre-IC<sup>3-5</sup>.  
122 Additional factors join to form the RPC to stimulate replication elongation<sup>29</sup>. While the precise

123 replisome protein complement varies somewhat between different eukaryotes, metamonads show  
124 dramatic variation in ORC, pre-RC and replicative polymerases (**Fig. 1**). The presence-absence of  
125 ORC and Cdc6 proteins is notably patchy across Metamonada. Strikingly, whereas all most  
126 metamonads retain up to two paralogs of the core protein family Orc1/Cdc6 (here called Orc1 and  
127 Orc1/Cdc6-like, **Supplementary Figure 1**), plus some orthologs of Orc 2-6, all these proteins are  
128 absent in *C. membranifera* and *Carpediemonas frisia* (**Fig. 1, Supplementary Table 1**). The lack of  
129 these proteins in an eukaryote is unexpected and unprecedented, since their absence would be  
130 expected to make the genome prone to DSBs and impair DNA replication, as well as interfere with  
131 other non-replicative processes<sup>30</sup>. To rule out false negatives, we conducted further analyses using  
132 metamonad-specific HMMs (Hidden Markov Models), various other profile-based search strategies  
133 (**Supplementary Information**), tBLASTn<sup>31</sup> searches, and applied HMMER<sup>32</sup> on 6-frame assembly  
134 translations. These additional methods were sufficiently sensitive to identify these proteins in all  
135 nuclear genomes we examined, with the exception of the *Carpediemonas* species and the highly  
136 reduced, endosymbiotically-derived nucleomorphs of cryptophytes and chlorarachniophytes  
137 (**Supplementary Information, Supplementary Table 1, Supplementary Fig. 1 and 2**).  
138 *Carpediemonas* species are, therefore, the only known eukaryotes to completely lack ORC and Cdc6.

### 139 **DNA damage repair systems have undergone several modifications**

140 DNA repair occurs continuously during the cell cycle depending on the type or specificity of the  
141 lesion. Among the currently known mechanisms are base-excision repair (BER), nucleotide excision  
142 repair (NER), mismatch repair (MMR), and double strand break repair, with the latter conducted by  
143 either homologous recombination (HR), canonical non-homologous end joining (NHEJ) or alternative  
144 end joining (a-EJ)<sup>8,14</sup>. MMR can be coupled directly to replication or play a role in HR. MMR, BER  
145 and NER are present in all studied taxa (**Supplementary Table 1**), although our analyses indicate that

146 damage sensing and downstream functions in NER seem to be modified in the metamonad taxa  
147 Parabasalia and Fornicata due to the absence of the XPG and XPC sensor proteins.

148         Double strand breaks (DSBs) are extremely dangerous for cells and can occur as a result of  
149 damaging agents or from self-inflicted cuts during DNA repair and meiosis. NHEJ requires the  
150 heterodimer Ku70-Ku80 to recruit the catalytic kinase DNA-PKcs and accessory proteins.  
151 Metamonads lack all of these proteins, as do a number of other eukaryotes investigated here and in  
152 ref.<sup>33</sup>. The a-EJ system seems to be fully present in metamonads like *C. membranifera* and *T.*  
153 *vaginalis*, partial in others, and completely absent in parasitic diplomonads. NHEJ is thought to be the  
154 predominant mechanism for repairing DSBs in eukaryotes<sup>34</sup>, but since our analyses indicate this  
155 pathway is absent in metamonads and a-EJ is highly mutagenic<sup>8</sup>, the HR pathway is likely to be  
156 essential for DSB repair in most metamonads. Repair by the HR system occurs through multiple sub-  
157 pathways that are influenced by the extent of the similarity of the DNA template or its flanking  
158 sequences to the sequences near the break. HR complexes are recruited during DNA replication and  
159 transcription, and utilize DNA, transcript-RNA or newly synthesized transcript-cDNA as a  
160 homologous template<sup>11,35-40</sup>. These complexes are formed by recombinases from the RecA/Rad51  
161 family that interact with members of the Rad52 family and chromatin remodeling factors of the  
162 SNF2/SWI2 sub-family<sup>41,42</sup>. Although the recombinases Rad51A-D are all present in most eukaryotes,  
163 we found a patchy distribution in metamonads (**Supplementary Table 1, Supplementary Fig. 3**). All  
164 examined Fornicata have lost the major recombinase Rad51A and have two paralogs of the meiosis-  
165 specific recombinase Dmc1, as first noted in *Giardia intestinalis*<sup>43</sup>. Dmc1 has been reported to provide  
166 high stability to recombination due to strong D-loop resistance to strand dissociation<sup>44</sup>. The  
167 recombination mediator Rad52 is present in metamonads but Rad59 or Rad54 are not. Metamonads  
168 have no components of an ISWI remodeling complex yet retain a reduced INO80 complex. Therefore,



169 replication fork progression and HR are likely to occur under the assistance of INO80 alone. HR  
170 requires endonucleases and exonucleases, and our searches for proteins additional to those from the  
171 MMR pathway revealed a gene expansion of the Flap proteins from the Rad2/XPG family in some  
172 metamonads. We also found proteins of the PIF1 helicase family that encompasses homologs that  
173 resolve R-loop structures, unwind DNA–RNA hybrids and assists in fork progression in regular  
174 replication and HR<sup>45,46</sup>. Phylogenetic analysis reveals that although *Carpediemonas* species have  
175 orthologs that branch within a metamonad group in the main PIF1 clade (**Fig. 2**), they also possess a  
176 highly divergent clade of PIF1-like proteins. Each *Carpediemonas* species has multiple copies of PIF1-  
177 like proteins that have independently duplicated within each species; these may point to the *de novo*  
178 emergence of specialized functions in HR and DNA replication for these proteins. Metamonads appear  
179 capable of using all of the HR sub-pathways (*e.g.*, classical DSB repair, single strand annealing, break  
180 induced replication), but these are modified (**Supplementary Table 1, Supplementary Figure 3**).  
181 Overall, the presence-absence patterns of the orthologs involved in DSB repair in Fornicata point to  
182 the existence of a highly specialized HR pathway which is presumably not only essential for the cell  
183 cycle of metamonads but is also likely the major pathway for replication-related DNA repair and  
184 recombination.

185

### 186 **Modified DSB damage response checkpoints in metamonads**

187 Checkpoints constitute a cascade of signaling events that delay replication until DNA lesions are  
188 resolved<sup>13</sup>. The ATR-Chk1, ATM-Chk2 and DNA-PKcs pathways are activated by the interaction of  
189 TopBP1 and the 9-1-1 complex (Rad9-Hus1-Rad1) for DNA repair regulation during replication stress  
190 and response to DSBs<sup>47</sup>. The ATR-Chk1 signaling pathway is the initial response to ssDNA damage  
191 and is responsible for the coupling of DNA replication with mitosis, but when it is defective, the

192 ssDNA is converted into DSBs to activate the ATM-Chk2 pathway. The DNA-PKcs act as sensors of  
193 DSBs to promote NHEJ, but we found no homologs of DNA-PKcs in metamonads (**Supplementary**  
194 **Fig. 3**), which is consistent with the lack of a NHEJ repair pathway in the group. All the checkpoint  
195 pathways described are present in humans and yeasts, while the distribution of core checkpoint  
196 proteins in the remaining taxa is patchy. Notably, Fornicata lack several of the proteins thought to be  
197 needed to activate the signaling kinase cascades and, while orthologs of ATM or ATR kinases are  
198 present in some fornicates, there are no clear orthologs of Chk1 or Chk2 in metamonads except in  
199 *Monocercomonoides exilis* (**Supplementary Table 1, Supplementary Fig. 3**). *Carpediemonas* species  
200 and *Kipferlia bialata* contain ATM and ATR but lack Chk1, Chk2, Rad9 and Hus1. Diplomonads  
201 possess none of these proteins, except the free-living *Trepomonas* sp. PC1, which has only ATM. The  
202 depletion of Chk1 has been shown to increase the incidence of chromosomal breaks and mis-  
203 segregation<sup>48</sup> and the absence of Rad9 has been associated with changes in checkpoint responses in  
204 origin-deficient yeasts<sup>49</sup>. Together with the loss of sensors, these absences reinforce the idea that the  
205 checkpoint controls in Fornicata are non-canonical.

## 206 **Reduction of mitosis and meiosis machinery in metamonads**

207 Eukaryotes synchronize cell cycle progression with chromosome segregation by a kinetochore based  
208 signaling system called the spindle assembly checkpoint (SAC)<sup>50,51</sup> that is ancestral to all eukaryotes  
209 (**Fig. 3A, B**). Kinetochores primarily form microtubule attachments through the Ndc80 complex,  
210 which is connected through a large network of structural subunits to a histone H3-variant CenpA that  
211 is specifically deposited at centromeres<sup>12</sup>. To prevent premature chromosome segregation, unattached  
212 kinetochores catalyse the production of the Mitotic Checkpoint Complex (MCC)<sup>50</sup>, a cytosolic  
213 inhibitor of the Anaphase Promoting Complex/Cyclosome (APC/C), a large multi-subunit E3 ubiquitin  
214 ligase that drives progression into anaphase by promoting the proteolysis of its substrates such as

215 various Cyclins<sup>52</sup> (**Fig. 3A**). Our analysis indicates the reduction of ancestral complexity of these  
216 proteins in metamonads (**Fig. 3C, Supplementary Table 1, Supplementary Fig. 4**). Surprisingly,  
217 such reduction is most extensive in *Carpediemonas* species. We found that most structural kinetochore  
218 subunits, a microtubule plus-end tracking complex and all four subunits of the Ndc80 complex are  
219 absent (**Fig. 3C, Supplementary Fig. 4**). None of our additional search strategies led to the  
220 identification of Ndc80 complex members, making *Carpediemonas* the only known eukaryotic lineage  
221 without it, except for kinetoplastids, which appear to have lost the canonical kinetochore and replaced  
222 it by an analogous molecular system, although there is still some controversy about this loss<sup>53,54</sup>. With  
223 such widespread absence of kinetochore components it might be possible that *Carpediemonas*  
224 underwent a similar replacement process to that of kinetoplastids<sup>53</sup>. We did however find a potential  
225 candidate for the centromeric Histone H3-variant (CenpA) in *C. membranifera*. CenpA forms the basis  
226 of the canonical kinetochore in most eukaryotes<sup>55</sup> (**Supplementary Fig. 5**). On the other hand, the  
227 presence or absence of CenpA is often correlated with the presence/absence of its direct interactor  
228 CenpC<sup>19</sup>. Similar to diplomonads, *C. membranifera* lacks CenpC and therefore the molecular network  
229 associated with kinetochore assembly on CenpA chromatin may be very different.

230 Most metamonads encode all MCC components, but diplomonads lost the SAC response and  
231 the full APC/C complex<sup>56</sup>. In contrast, only *Carpediemonas* species and *K. bialata* have MCC subunits  
232 that contain the conserved short linear motifs to potentially elicit a canonical SAC signal<sup>52,57</sup>  
233 (**Supplementary Fig. 6**). Interestingly, not all of these motifs are present, and most are seemingly  
234 degenerate compared to their counterparts in other eukaryotic lineages (**Supplementary Fig. 6C**).  
235 Also, many other SAC-related genes are conserved, even in diplomonads (*e.g.*, Mad2, MadBub)<sup>56</sup>.  
236 Furthermore, the cyclins in *C. membranifera*, the main target of SAC signalling, have a diverged  
237 destruction motif (D-box) in their N-termini (**Supplementary Fig. 6C**). Collectively, our observations

238 indicate that *Carpediemonas* species could elicit a functional SAC response, but whether this would be  
239 kinetochore-based is unclear. Alternatively, SAC-related genes could have been repurposed for another  
240 cellular function(s) as in diplomonads<sup>56</sup>. Given that ORC has been observed to interact with the  
241 kinetochore (throughout chromosome condensation and segregation), centrioles and promotes  
242 cytokinesis<sup>30</sup>, the lack of Ncd80 and ORC complexes suggest that *Carpediemonas* species possess  
243 radically unconventional cell division systems.

244         Neither sexual nor parasexual processes have been directly observed in Metamonada<sup>43</sup>.  
245 Nonetheless, our surveys confirm the conservation of the key meiotic proteins in metamonads<sup>43</sup>,  
246 including Hap2 (for plasmogamy) and Gex1 (karyogamy). Unexpectedly, *Carpediemonas* species have  
247 homologs from the tmcB family that acts in the cAMP signaling pathway specific for sexual  
248 development in *Dictyostelium*<sup>58</sup>, and sperm-specific channel subunits (*i.e.*, CatSper  $\alpha$ ,  $\beta$ ,  $\delta$  and  $\gamma$ )  
249 reported previously only in Opisthokonta and three other protists<sup>59</sup>. In opisthokonts, the CatSper  
250 subunits enable the assembly of specialized Ca<sup>2+</sup> influx channels and are involved in the signaling for  
251 sperm maturation and motility<sup>59</sup>. In *Carpediemonas*, the tmcB family and CatSper subunits could  
252 similarly have a role in signaling and locomotion pathways required for a sexual cycle. As proteins in  
253 the cAMP pathway and Ca<sup>2+</sup> signaling cooperate to generate a variety of complex responses, the  
254 presence of these systems in *Carpediemonas* species but absence in all other sampled metamonads is  
255 intriguing and deserves further investigation. Even if these systems are not directly involved in a  
256 sexual cycle, the presence of Hap2 and Gex1 proteins is strong evidence that *C. membranifera* can  
257 reproduce sexually. Interestingly, based on the frequencies of single nucleotide polymorphisms, *C.*  
258 *membranifera* is predicted to be haploid (**Supplementary Fig. 7**). If this is correct, its sexual  
259 reproduction should include the formation of a zygote followed by a meiotic division to regain its  
260 haploid state<sup>60</sup>.

261 **Acquisition of DNA replication and repair proteins in *Carpediemonas* by lateral gene transfer**

262 The unprecedented absence of many components of canonical DNA replication, repair, and  
263 segregation systems in *Carpediemonas* species led us to investigate whether they had been replaced  
264 by analogous systems acquired by lateral gene transfer (LGT) from viruses or prokaryotes. We  
265 detected four Geminivirus-like replication initiation protein sequences in the *C. membranifera*  
266 genome but not in *C. frisia*, and helitron-related helicase endonucleases in both *Carpediemonas*  
267 genomes. All these genes were embedded in high-coverage eukaryotic scaffolds, yet all of them lack  
268 introns and show no evidence of gene expression in the RNA-Seq data. As RNA was harvested from  
269 log-phase actively replicating cell cultures, their lack of expression suggests it is unlikely that these  
270 acquired proteins were coopted to function in the replication of the *Carpediemonas* genomes.  
271 Nevertheless, the presence of Geminivirus protein-coding genes is intriguing as these viruses are  
272 known, in other systems (*e.g.*, plants, insects), to alter host transcriptional controls and reprogram the  
273 cell-cycle to induce the host DNA replication machinery<sup>61,62</sup>. We also detected putative LGTs of  
274 Endonuclease IV, RarA and RNase H1 from prokaryotes into a *Carpediemonas* ancestor  
275 **(Supplementary Information, Supplementary Fig. 8, 9 and 10)**. Of these, RarA is ubiquitous in  
276 bacteria and eukaryotes and acts during replication and recombination in the context of collapsed  
277 replication forks<sup>63,64</sup>. Interestingly, *Carpediemonas* appears to have lost the eukaryotic ortholog, and  
278 only retains the acquired prokaryotic-like RarA, a gene that is expressed (*i.e.*, transcripts are present  
279 in the RNA-Seq data). RNase Hs are involved in the cleavage of RNA from RNA:DNA hybrid  
280 structures that form during replication, transcription, and repair, and, while eukaryotes have a  
281 monomeric RNase H1 and a heterotrimeric RNase H2, prokaryotes have either one or both types.  
282 Eukaryotic RNase H1 removes RNA primers during replication and R-loops during transcription,  
283 and also participates in HR-mediated DSB repair<sup>65,66</sup>. The prokaryotic homologs have similar roles

284 during replication and transcription<sup>67</sup>. *C. membranifera* lacks a typical eukaryotic RNase H1 but has  
285 two copies of prokaryotic homologs. Both are located in scaffolds comprising intron-containing  
286 genes and have RNA-Seq coverage, clearly demonstrating that they are not from prokaryotic  
287 contaminants in the assembly.

288

## 289 **Discussion**

### 290 **Genome streamlining in metamonads**

291 The reductive evolution of the DNA replication and repair, and segregation systems and the low  
292 retention of proteins in the BUSCO dataset in metamonads demonstrate that substantial gene loss has  
293 occurred (**Supplementary information**), providing additional evidence for streamlining of gene  
294 content prior to the last common ancestor of Metamonada<sup>15-17</sup>. However, the patchy distribution of  
295 genes within the group suggests ongoing differential reduction in different metamonad groups. Such  
296 reduction – especially the unprecedented complete absence of systems such as the ORC, Cdc6 and  
297 kinetochore Ndc80 complexes in *Carpediemonas* species – demands an explanation. Whereas the loss  
298 of genes from varied metabolic pathways is well known in lineages with different lifestyles<sup>68-73</sup>, loss of  
299 cell cycle, DNA damage sensing and repair genes in eukaryotes is extremely rare. New evidence from  
300 yeasts of the genus *Hanseniaspora* suggests that the loss of proteins in these systems can lead to  
301 genome instability and long-term hypermutation leading to high rates of sequence substitution<sup>68</sup>. This  
302 could also apply to metamonads, especially fornicates, which are well known to have undergone rapid  
303 sequence evolution; these taxa form a highly divergent clade with very long branches in phylogenetic  
304 trees<sup>20,74</sup>. Most of the genes that were retained by Metamonada in the various pathways we examined  
305 were divergent in sequence relative to homologs in other eukaryotes and many of the gene losses  
306 correspond to proteins that are essential in model system eukaryotes. Gene essentiality appears to be

307 relative and context-dependent, and some studies have shown that the loss of ‘indispensable’ genes  
308 could be permitted by evolving divergent pathways that provide similar activities via chromosome  
309 stoichiometry changes and compensatory gene loss<sup>68-70,75</sup>.

310         The patchy distribution of genes from different ancestral eukaryotic pathways suggests that the  
311 last common ancestor of Metamonada had a broad gene repertoire for maintaining varied metabolic  
312 functions under fluctuating environmental conditions offered by diverse oxygen-depleted habitats.  
313 Although the loss of proteins and genomic streamlining are well known in parasitic diplomonads<sup>15,16</sup>,  
314 the Fornicata, as a whole, tend to have a reduced subset of the genes that are commonly found in core  
315 eukaryotic pathways. In general, such gene content reduction can partially be explained as the result of  
316 historical and niche-specific adaptations<sup>76</sup>. Yet, given that 1) genome maintenance mostly depends on  
317 the cell cycle checkpoints, DNA repair pathways, and their interactions<sup>14,77</sup>, 2) the lack of several  
318 proteins related to these pathways that were present in the last common ancestor of metamonads, 3)  
319 aneuploidy and high overall rates of sequence evolution have been observed in metamonads<sup>78,79</sup>, and,  
320 4) the loss of DNA repair genes can be associated with substantial gene loss and sequence instability  
321 that apparently boosts the rates of sequence evolution<sup>68</sup>, it is likely that genome evolution in the  
322 Fornicata clade has been heavily influenced by their error-prone DNA maintenance mechanisms.

323  
324 **Non-canonical replication initiation and replication licensing in *Carpodimonas*.**

325 Origin-independent replication has been observed in the context of DNA repair (reviewed in ref.<sup>10</sup>) and  
326 in origin-deficient or -depleted chromosomes in yeast<sup>80</sup>. These studies have highlighted the lack of (or  
327 reduction in) the recruitment of ORC and Cdc6 onto the DNA, but no study to date has documented  
328 regular eukaryotic DNA replication in the absence of genes encoding these proteins. While it is  
329 possible that extremely divergent versions of ORC and Cdc6 are governing the recognition of origins

330 of replication and replication licensing in *Carpodomonas* species, we have no evidence for this.  
331 Instead, our findings suggest the existence of an as-yet undiscovered underlying eukaryotic system that  
332 can accomplish eukaryotic DNA replication initiation and licensing. The existence of such a system  
333 has in fact already been suspected<sup>81</sup> given that: 1) Orc1- or Orc2-depleted human cells and mouse and  
334 fruit-fly ORC mutants are viable and capable of undergoing replication and endoreplication<sup>81-84</sup>, and 2)  
335 origin-independent replication at the chromosome level has been reported<sup>80,85,86</sup>. We propose that  
336 *Carpodomonas* species utilize an alternative DNA replication system based on a Dmc1-dependent HR  
337 mechanism that is origin-independent and mediated by RNA:DNA hybrids. Here we summarize  
338 evidence that such a mechanism is possible based on what is known in model systems and present a  
339 hypothetical model as to how it might occur in *Carpodomonas*.

340         During replication and transcription, the HR complexes, RNase H1 and RNA-interacting  
341 proteins are recruited onto the DNA to assist in its repair<sup>36,37,87</sup>. Remarkably, experiments show that  
342 HR is able to carry out full genome replication in archaea, bacteria, viruses, and linear mtDNA<sup>86,88-90</sup>,  
343 with replication fork progression rates that are comparable to those of regular replication<sup>27,91-94</sup>. A  
344 variety of *cis* and *trans* homologous sequences (*e.g.*, chromatids, transcript-RNA or -cDNA) can be  
345 used as templates<sup>27,36,40</sup>, and their length as well as the presence of one or two homologous ends likely  
346 influence a recombination execution checkpoint that decides which HR sub-pathway is utilized<sup>94</sup>. For  
347 example, in the absence of a second homologous end, HR by Rad51-dependent break-induced  
348 replication (BIR) can either use a newly synthesized DNA strand or independently invade donor  
349 sequences, such that the initial strand invasion intermediate creates a migrating D-loop and DNA is  
350 synthesized conservatively<sup>27,94,95</sup>. Studies have found that BIR does not require the assembly of an  
351 ORC complex and Cdc6 but the recruitment of the Cdc7, loading of MCM helicase, firing factors and  
352 replicative polymerases are needed for assembling the pre-RC complex<sup>27,94</sup>. The requirement of MCM



353 for BIR was questioned, as PIF1 helicase was found to be essential for long-range BIR<sup>96</sup>. However,  
354 recent evidence shows that MCM is typically recruited for unwinding DNA strands during HR and is  
355 likely needed together with PIF1 to enhance processivity<sup>97,98</sup>. All these proteins are also suspected to  
356 operate during origin-independent transcription-initiated replication (TIR), a still-enigmatic  
357 mechanism that is triggered by R-loops resulting from RNA:DNA hybrids during transcription<sup>10,11,99</sup>.

358         Considering the complement of proteins in *Carpediemonas* species discussed above, and that  
359 RNA:DNA hybrids are capable of promoting origin-independent replication in model systems<sup>11,39,100</sup>,  
360 we suggest that a Dmc1-dependent HR replication mechanism is enabled by excess of RNA:DNA  
361 hybrids in these organisms. In such a system, DSBs generated in stressed transcription-dependent R-  
362 loops could be repaired by HR with either transcript-RNA- or transcript-cDNA-templates and the *de*  
363 *novo* assembly of the replisome as in BIR (**Fig. 4**). The establishment of a replication fork could be  
364 favored by the presence of *Carpediemonas*-specific PIF1-like homologs, as these raise the possibility  
365 of the assembly of a multimeric PIF1 helicase with increased capability to bind multiple sites on the  
366 DNA, thereby facilitating DNA replication processivity and regulation<sup>45</sup>. The loss of Rad51A and the  
367 duplication of Dmc1 recombinases suggests that a Dmc1-dependent HR mechanism was likely  
368 enabled in the last common ancestor of Fornicata and this mechanism may have become the  
369 predominant replication pathway in the *Carpediemonas* lineage after its divergence from the other  
370 fornicates, ultimately leading to the loss of ORC and Cdc6 proteins.

371

### 372 **The impact of cell cycle dysregulation on genome evolution.**

373 DNA replication licensing and firing are temporally separated (*i.e.*, they occur at G1 and S phases  
374 respectively) and are the principal ways to counteract damaging over-replication<sup>6</sup>. As S-phase is  
375 particularly vulnerable to DNA errors and lesions, its checkpoints are likely more important for

376 preventing genome instability than those of G1, G2 or SAC<sup>101</sup>. Dysregulation is anticipated if no  
377 ORC/Cdc6 are present as licensing would not take place and replication would be blocked<sup>28</sup>. Yet this  
378 clearly does not happen in *Carpediemonas*. This implies that during late G1 phase, activation by  
379 loading the MCM helicase has to occur by an alternative mechanism that is still unknown but might  
380 already be in place in eukaryotes. Such a mechanism has long been suspected as it could explain the  
381 over-abundance and distribution patterns of MCM on the DNA (*i.e.*, the MCM paradox; reviewed in  
382 <sup>102</sup>).

383 In terms of the regulation of M-phase progression, the extremely divergent nature of the  
384 kinetochore in *C. membranifera* could suggest that it uses different mechanisms to execute mitosis  
385 and meiosis. It is known that in *Carpediemonas*-related fornicates such as retortamonads and in  
386 diplomonads, chromosome segregation proceeds inside a persisting nuclear envelope, with the aid of  
387 intranuclear microtubules, but with the mitotic spindle nucleated outside the nucleus (*i.e.*, semi-open  
388 mitosis)<sup>79</sup>. Although mitosis in *Carpediemonas* has not been directly observed, these organisms may  
389 also possess a semi-open mitotic system such as the ones found in other fornicates. Yet how the  
390 *Carpediemonas* kinetochore functions in the complete absence of the microtubule-binding Ndc80  
391 complex remains a mystery; it is possible that, like in kinetoplastids<sup>48</sup>, other molecular complexes have  
392 evolved in this lineage that fulfill the roles of Ndc80 and other kinetochore complexes.

393 Interestingly, a potential repurposing of SAC proteins seems to have occurred in the  
394 diplomonad *G. intestinalis*, as it does not arrest under treatment with microtubule-destabilizing drugs  
395 and Mad2 localizes to a region of the intracytoplasmic axonemes of the caudal flagella<sup>56</sup>. Other  
396 diplomonads have a similar SAC protein complement that may have a similar non-canonical function.  
397 In contrast to diplomonads, our investigations (**Fig. 3**) suggest that *Carpediemonas* species could elicit

398 a functional SAC response, although microtubule-disrupting experiments during mitosis will be  
399 needed to prove its existence.

400 In addition to the aforementioned apparent dysregulation of checkpoint controls in  
401 *Carpediemonas*, alternative mechanisms for chromosome condensation, spindle attachment, sister  
402 chromatid cohesion, cytokinesis, heterochromatin formation, and silencing and transcriptional  
403 regulation can also be expected in this organism due to the absence of ORC and Cdc6 (reviewed in  
404 refs<sup>30,103,104</sup>). All of the absences of canonical eukaryotic systems we have described for  
405 *Carpediemonas* suggest that a radically different cell cycle has evolved in this free-living protistan  
406 lineage. This underscores the fact that our concepts of universality and essentiality rely on studies of  
407 a very small subset of organisms. The development of *Carpediemonas* as a model system thus has  
408 great potential to enhance our understanding of fundamental DNA replication, repair and cell cycle  
409 processes. It could even reveal widely conserved alternative, but as-yet unknown, mechanisms  
410 underpinning the evolutionary plasticity of these systems across the eukaryote tree of life.

411

## 412 **Methods**

### 413 **Sequencing, assembly, and protein prediction for *C. membranifera***

414 DNA and RNA were isolated from log-phase cultures of *C. membranifera* BICM strain (see details in  
415 **Supplementary Information**). Sequencing employed Illumina short paired-end and long read  
416 (Oxford Nanopore MinION) technologies. For Illumina, extracted, purified DNA and RNA (*i.e.*,  
417 cDNA) were sequenced on the HiSeq 2000 (150 x 2 paired-end) at the Genome Québec facility.  
418 Illumina reads were quality trimmed (Q=30) and filtered for length (>40 bp) with Trimmomatic<sup>105</sup>.  
419 For MinION, the library was prepared using the 1D native barcoding genomic DNA (SQK-LSK108  
420 with EXP-NBD103) protocol (NBE\_9006\_v103\_revP\_21Dec2016). The final library (1070 ng) was

421 loaded on a R9.4 flow cell and sequenced for 48 h on the MinION Mk1B nanopore sequencer. The  
422 long reads were base-called and trimmed with Albacore v2.3.3 ([www.nanoporetech.com](http://www.nanoporetech.com)) and  
423 Porechop v0.2.3 ([www.github.com/rrwick/Porechop](http://www.github.com/rrwick/Porechop)), respectively. Canu v1.6<sup>106</sup> with default  
424 parameters and max genome size of 30Mb produced an assembly that was polished with Nanopolish  
425 v0.10.1<sup>107</sup>. The latter was iteratively error-corrected with the genomic paired-end Illumina reads  
426 using Unicycler<sup>108</sup>. The identification and removal of prokaryotic contigs was assisted by BLASTx  
427 and BLASTn searches against the nt database. Read-depth coverage at each position of the genomic  
428 scaffolds were obtained with samtools<sup>109</sup> and mosdepth v0.2.5<sup>110</sup>.

429 RNA-Seq reads were used for genome-independent assessments of the presence of the proteins  
430 of interest and to generate intron junction hints for gene prediction. For the independent assessments  
431 we obtained both a *de novo* and a genome-guided transcriptome assembly with Trinity v2.5.0<sup>111</sup>. Open  
432 reading frames were translated with TransDecoder v5.5.0 ([www.github.com/TransDecoder](http://www.github.com/TransDecoder)) and were  
433 included in all of our analyses. Gene predictions were carried out as follows: repeat libraries were  
434 obtained and masked with RepeatModeler 1.0 and RepeatMasker (<http://www.repeatmasker.org>).  
435 Then, RNA-Seq reads were mapped onto the assembly using Hisat2<sup>112</sup>, generating a bam file for  
436 GenMarkET<sup>113</sup>. This resulted in a list of intron hints used to train Augustus v3.2.3<sup>114</sup>. The genome-  
437 guided assembled transcriptome, genomic scaffolds and the newly predicted proteome were fed into  
438 the PASA pipeline<sup>115</sup> to yield a more accurate set of predicted proteins. Finally, the predicted proteome  
439 was manually curated for the proteins of interest.

#### 440 **Genome size, completeness, and ploidy assessments**

441 We estimated the completeness of the draft genome by 1) using the k-mer based and reference free  
442 method Merquy<sup>23</sup>, 2) calculating the percentage of transcripts that aligned to the genome, and 3)  
443 employing the BUSCO<sup>116</sup> framework. For method 1, all paired-end reads were used to estimate the

444 best k-mer and create ‘meryl’ databases necessary to apply Mercury<sup>23</sup>. For method 2, transcripts were  
445 mapped onto the genome using BLASTn and exonerate<sup>117</sup>. For method 3, the completeness of the  
446 draft genome was evaluated in a comparative setting by including the metamonads and using the  
447 universal single copy orthologs (BUSCO) from the Eukaryota (odb9) and protist databases  
448 (<https://busco.ezlab.org/>), which contain 303 and 215 proteins, respectively. Each search was run  
449 separately on the assembly and the predicted proteome for all these taxa. Unfortunately, both  
450 BUSCO database searches yielded false negatives in that several conserved proteins publicly  
451 reported for *T. vaginalis*, *G. intestinalis* and *Spironucleus salmonicida* were not detected due to the  
452 extreme divergence of metamonad homologs. Therefore, genome completeness was re-assessed with  
453 a phylogeny-guided search (**Supplementary Information**).

454 The ploidy of *C. membranifera* was inferred by *i*) counting k-mers with Mercury<sup>23</sup>, and *ii*)  
455 mapping 613,266,290 Illumina short reads to the assembly with Bowtie 2.3.1<sup>118</sup> and then using  
456 ploidyNGS<sup>119</sup> to calculate the distribution of allele frequencies across the genome. A site was deemed  
457 to be heterozygous if at least two different bases were present and there were at least two reads with  
458 the different bases. Positions with less than 10× coverage were ignored.

459

## 460 **Functional annotation of the predicted proteins**

461 Our analyses included the genomes and predicted proteomes of *C. membranifera* (reported here) as  
462 well as publicly available data for nine additional metamonads and eight other eukaryotes  
463 representing diverse groups across the eukaryotic tree of life (**Fig. 1, Supplementary Information**).  
464 Orthologs from each of these 18 predicted proteomes were retrieved for the assessment of core  
465 cellular pathways, such as DNA replication and repair, mitosis and meiosis and cell cycle  
466 checkpoints. For *C. membranifera*, we included the predicted proteomes derived from the assembly

467 plus the 6-frame translated transcriptomes. Positive hits were manually curated in the *C.*  
468 *membranifera* draft genome. A total of 367 protein queries were selected based on an extensive  
469 literature review and prioritizing queries from taxa in which they had been experimentally  
470 characterized. The identification of orthologs was as described for the BUSCO proteins but using  
471 these 367 queries for the initial BLASTp (**Supplementary Information**), except for kinetochore  
472 (KT), Spindle assembly check point (SAC) and anaphase-promoting complex-related genes (APC/C).  
473 For these, previously published refined HMMs with cut-offs specific to each orthologous group  
474 (see<sup>58</sup>) were used to query the proteomes with HMMER v3.1b2<sup>32</sup>. A multiple sequence alignment  
475 that included the newly-found hits was subsequently constructed with MAFFT v7.310<sup>120</sup> and was  
476 used in HMM searches for more divergent homologs. This process was iterated until no new  
477 significant hits could be found. As we were unable to retrieve orthologs of a number of essential  
478 proteins in the *C. membranifera* and *C. frisia* genomes, we embarked on additional more sensitive  
479 strategies to detect them using multiple different HMMs based on aligned homologs from archaea,  
480 metamonads, and broad samplings of taxa. Individual PFAM domains were searched for in the  
481 genomes, proteome and transcriptomes with e-value thresholds of  $10^{-3}$  (**Supplementary**  
482 **Information**). To rule out that failure to detect these proteins was due to insufficient sensitivity of  
483 our methods when applied them to highly divergent taxa, we queried 22 extra eukaryotic genomes  
484 with demonstrated high rates of sequence evolution, genome streamlining or unusual genomic  
485 features (**Supplementary Table 1, Supplementary Information**). Possible non-predicted or mis-  
486 predicted genes were investigated using tBLASTn and 6-frame translation HMMER searches of the  
487 genomic scaffolds. Also, as DNA replication and repair genes could have been acquired by lateral  
488 gene transfer into *Carpodimonas* species from prokaryotes or viruses, proteins from the DNA  
489 replication and repair categories whose best matches were to prokaryotic and viral homologs were

490 subjected to phylogenetic analysis using the methods described for the phylogeny-guided BUSCO  
491 analysis and using substitution models specified in the legend of each tree (**Supplementary**  
492 **Information**).

#### 493 **Data availability**

494 Genome assembly is available at NCBI under BioProject <XXXX>, accession number <XXXX>.  
495 DNA and RNA-Seq reads are available at SRA under accessions <XXXX> and <XXXX>,  
496 respectively.

#### 497 **References**

- 498 1 Yeeles, J. T., Deegan, T. D., Janska, A., Early, A. & Diffley, J. F. Regulated eukaryotic DNA  
499 replication origin firing with purified proteins. *Nature* **519**, 431-435 (2015).
- 500 2 Parker, M. W., Botchan, M. R. & Berger, J. M. Mechanisms and regulation of DNA  
501 replication initiation in eukaryotes. *Crit. Rev. Biochem. Mol. Biol.* **52**, 107-144 (2017).
- 502 3 Shen, Z. & Prasanth, S. G. Emerging players in the initiation of eukaryotic DNA replication.  
503 *Cell Div* **7**, 22 (2012).
- 504 4 Burgers, P. M. J. & Kunkel, T. A. Eukaryotic DNA replication fork. *Annu. Rev. Biochem.* **86**,  
505 417-438 (2017).
- 506 5 Riera, A. *et al.* From structure to mechanism-understanding initiation of DNA replication.  
507 *Genes Dev.* **31**, 1073-1088 (2017).
- 508 6 Reuswig, K. U. & Pfander, B. Control of eukaryotic DNA replication initiation-mechanisms  
509 to ensure smooth transitions. *Genes (Basel)* **10** (2019).
- 510 7 Waters, L. S. *et al.* Eukaryotic translesion polymerases and their roles and regulation in DNA  
511 damage tolerance. *Microbiol. Mol. Biol. Rev.* **73**, 134-154 (2009).

- 512 8 Chang, H. H. Y., Pannunzio, N. R., Adachi, N. & Lieber, M. R. Non-homologous DNA end  
513 joining and alternative pathways to double-strand break repair. *Nat. Rev. Mol. Cell Biol.* **18**,  
514 495-506 (2017).
- 515 9 Wright, W. D., Shah, S. S. & Heyer, W. D. Homologous recombination and the repair of  
516 DNA double-strand breaks. *J. Biol. Chem.* **293**, 10524-10535 (2018).
- 517 10 Ravoityte, B. & Wellinger, R. E. Non-canonical replication initiation: you're fired! *Genes*  
518 (*Basel*) **8** (2017).
- 519 11 Stuckey, R., Garcia-Rodriguez, N., Aguilera, A. & Wellinger, R. E. Role for RNA:DNA  
520 hybrids in origin-independent replication priming in a eukaryotic system. *Proc. Natl. Acad.*  
521 *Sci. U.S.A* **112**, 5779-5784 (2015).
- 522 12 Musacchio, A. & Desai, A. A molecular view of kinetochore assembly and function. *Biology*  
523 (*Basel*) **6** (2017).
- 524 13 Hustedt, N., Gasser, S. M. & Shimada, K. Replication checkpoint: tuning and coordination of  
525 replication forks in s phase. *Genes (Basel)* **4**, 388-434 (2013).
- 526 14 Hakem, R. DNA-damage repair; the good, the bad, and the ugly. *EMBO J.* **27**, 589-605  
527 (2008).
- 528 15 Adam, R. D. *et al.* Genome sequencing of *Giardia lamblia* genotypes A2 and B isolates (DH  
529 and GS) and comparative analysis with the genomes of genotypes A1 and E (WB and Pig).  
530 *Genome Biol. Evol.* **5**, 2498-2511 (2013).
- 531 16 Xu, F. *et al.* The genome of *Spironucleus salmonicida* highlights a fish pathogen adapted to  
532 fluctuating environments. *PLoS Genet.* **10**, e1004053 (2014).
- 533 17 Tanifuji, G. *et al.* The draft genome of *Kipferlia bialata* reveals reductive genome evolution  
534 in fornicate parasites. *PLoS One* **13**, e0194487 (2018).



- 535 18 Ocana-Pallares, E. *et al.* Origin recognition complex (ORC) evolution is influenced by global  
536 gene duplication/loss patterns in eukaryotic genomes. *Genome Biol. Evol.* **12**, 3878-3889  
537 (2020).
- 538 19 van Hooff, J. J., Tromer, E., van Wijk, L. M., Snel, B. & Kops, G. J. Evolutionary dynamics  
539 of the kinetochore network in eukaryotes as revealed by comparative genomics. *EMBO Rep.*  
540 **18**, 1559-1571 (2017).
- 541 20 Hampl, V. *et al.* Phylogenomic analyses support the monophyly of Excavata and resolve  
542 relationships among eukaryotic "supergroups". *Proc. Natl. Acad. Sci. U.S.A.* **106**, 3859-3864  
543 (2009).
- 544 21 Karnkowska, A. *et al.* A eukaryote without a mitochondrial organelle. *Curr. Biol.* **26**, 1274-  
545 1284 (2016).
- 546 22 Simpson, A. G. B. & Patterson, D. J. The ultrastructure of *Carpodiemonas membranifera*  
547 (Eukaryota) with reference to the "excavate hypothesis". *Eur. J. Protistol.* **35**, 353-370  
548 (1999).
- 549 23 Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality,  
550 completeness, and phasing assessment for genome assemblies. *Genome Biol.* **21**, 245 (2020).
- 551 24 Yubuki, N. & Leander, B. S. Evolution of microtubule organizing centers across the tree of  
552 eukaryotes. *Plant J.* **75**, 230-244 (2013).
- 553 25 Morrison, H. G. *et al.* Genomic Minimalism in the Early Diverging Intestinal Parasite *Giardia*  
554 *lamblia*. *Science* **317**, 1921-1926 (2007).
- 555 26 Ebbert, M. T. W. *et al.* Systematic analysis of dark and camouflaged genes reveals disease-  
556 relevant genes hiding in plain sight. *Genome Biol.* **20**, 97 (2019).

- 557 27 Lydeard, J. R. *et al.* Break-induced replication requires all essential DNA replication factors  
558 except those specific for pre-RC assembly. *Genes Dev.* **24**, 1133-1144 (2010).
- 559 28 Liu, J. *et al.* Structure and function of Cdc6/Cdc18: implications for origin recognition and  
560 checkpoint control. *Mol. Cell* **6**, 637-648 (2000).
- 561 29 Georgescu, R. E. *et al.* Reconstitution of a eukaryotic replisome reveals suppression  
562 mechanisms that define leading/lagging strand operation. *Elife* **4**, e04988 (2015).
- 563 30 Popova, V. V., Brechalov, A. V., Georgieva, S. G. & Kopytova, D. V. Nonreplicative  
564 functions of the origin recognition complex. *Nucleus* **9**, 460-473 (2018).
- 565 31 Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421  
566 (2009).
- 567 32 Eddy, S. R. Accelerated profile HMM searches. *PLoS Comp. Biol.* **7**, e1002195 (2011).
- 568 33 Nenarokova, A. *et al.* Causes and effects of loss of classical non-homologous end joining  
569 pathway in parasitic eukaryotes. *MBio* (2019).
- 570 34 van den Berg, J. *et al.* DNA end-resection in highly accessible chromatin produces a toxic  
571 break. *BioRxiv* (2019).
- 572 35 Wei, L., Levine, A. S. & Lan, L. Transcription-coupled homologous recombination after  
573 oxidative damage. *DNA Repair* **44**, 76-80 (2016).
- 574 36 Meers, C., Keskin, H. & Storici, F. DNA repair by RNA: templated, or not templated, that is  
575 the question. *DNA Repair (Amst)* **44**, 17-21 (2016).
- 576 37 Aymard, F. *et al.* Transcriptionally active chromatin recruits homologous recombination at  
577 DNA double-strand breaks. *Nat. Struct. Mol. Biol.* **21**, 366-374 (2014).

- 578 38 Wei, L. *et al.* DNA damage during the G0/G1 phase triggers RNA-templated, Cockayne  
579 syndrome B-dependent homologous recombination. *Proc. Natl. Acad. Sci. U.S.A* **112**, E3495-  
580 3504 (2015).
- 581 39 Keskin, H. *et al.* Transcript-RNA-templated DNA recombination and repair. *Nature* **515**, 436-  
582 439 (2014).
- 583 40 Storici, F., Bebenek, K., Kunkel, T. A., Gordenin, D. A. & Resnick, M. A. RNA-templated  
584 DNA repair. *Nature* **447**, 338-341 (2007).
- 585 41 Ceballos, S. J. & Heyer, W. D. Functions of the Snf2/Swi2 family Rad54 motor protein in  
586 homologous recombination. *Biochim. Biophys. Acta.* **1809**, 509-523 (2011).
- 587 42 Mazin, A. V. & Mazina, O. M. in *Molecular Life Sciences: An Encyclopedic Reference* (eds  
588 Robert D. Wells, Judith S. Bond, Judith Klinman, & Bettie Sue Siler Masters) 1009-1016  
589 (Springer New York, 2018).
- 590 43 Ramesh, M. A., Malik, S. B. & Logsdon, J. M., Jr. A phylogenomic inventory of meiotic  
591 genes; evidence for sex in *Giardia* and an early eukaryotic origin of meiosis. *Curr. Biol.* **15**,  
592 185-191 (2005).
- 593 44 Bugreev, D. V. *et al.* The resistance of DMC1 D-loops to dissociation may account for the  
594 DMC1 requirement in meiosis. *Nat. Struct. Mol. Biol.* **18**, 56-60 (2011).
- 595 45 Byrd, A. K. & Raney, K. D. Structure and function of Pif1 helicase. *Biochem. Soc. Trans.* **45**,  
596 1159-1171 (2017).
- 597 46 Hill, J., Eickhoff, P., Drury, L. S., Costa, A. & Diffley, J. F. X. The eukaryotic replisome  
598 requires an additional helicase to disarm dormant replication origins. *BioRxiv*,  
599 2020.2009.2017.301366 (2020).

- 600 47 Blackford, A. N. & Jackson, S. P. ATM, ATR, and DNA-PK: the trinity at the heart of the  
601 DNA damage response. *Mol. Cell.* **66**, 801-817 (2017).
- 602 48 Calzetta, N. L., Gonzalez Besteiro, M. A. & Gottifredi, V. Mus81-Eme1-dependent aberrant  
603 processing of DNA replication intermediates in mitosis impairs genome integrity. *Sci. Adv.* **6**  
604 (2020).
- 605 49 van Brabant, A. J., Buchanan, C. D., Charboneau, E., Fangman, W. L. & Brewer, B. J. An  
606 origin-deficient yeast artificial chromosome triggers a cell cycle checkpoint. *Mol. Cell* **7**, 705-  
607 713 (2001).
- 608 50 Sacristan, C. & Kops, G. J. Joined at the hip: kinetochores, microtubules, and spindle  
609 assembly checkpoint signaling. *Trends Cell Biol.* **25**, 21-28 (2015).
- 610 51 Kops, G. J. P. L., Snel, B. & Tromer, E. C. Evolutionary dynamics of the spindle assembly  
611 checkpoint in eukaryotes. *Curr. Biol.* **30**, R589-R602 (2020).
- 612 52 Alfieri, C., Zhang, S. & Barford, D. Visualizing the complex functions and mechanisms of the  
613 anaphase promoting complex/cyclosome (APC/C). *Open Biol.* **7** (2017).
- 614 53 Akiyoshi, B. & Gull, K. Discovery of unconventional kinetochores in kinetoplastids. *Cell*  
615 **156**, 1247-1258 (2014).
- 616 54 D'Archivio, S. & Wickstead, B. *Trypanosome* outer kinetochore proteins suggest conservation  
617 of chromosome segregation machinery across eukaryotes. *J. Cell Biol.* **216**, 379-391 (2017).
- 618 55 Drinnenberg, I. A., Henikoff, S. & Malik, H. S. Evolutionary turnover of kinetochore  
619 proteins: a ship of theseus? *Trends Cell Biol.* **26**, 498-510 (2016).
- 620 56 Markova, K. *et al.* Absence of a conventional spindle mitotic checkpoint in the binucleated  
621 single-celled parasite *Giardia intestinalis*. *Eur. J. Cell Biol.* **95**, 355-367 (2016).

- 622 57 Tromer, E., Bade, D., Snel, B. & Kops, G. J. Phylogenomics-guided discovery of a novel  
623 conserved cassette of short linear motifs in BubR1 essential for the spindle checkpoint. *Open*  
624 *Biol.* **6** (2016).
- 625 58 Muramoto, T., Takeda, S., Furuya, Y. & Urushihara, H. Reverse genetic analyses of gamete-  
626 enriched genes revealed a novel regulator of the cAMP signaling pathway in *Dictyostelium*  
627 *discoideum*. *Mech. Dev.* **122**, 733-743 (2005).
- 628 59 Cai, X., Wang, X. & Clapham, D. E. Early evolution of the eukaryotic Ca<sup>2+</sup> signaling  
629 machinery: conservation of the CatSper channel complex. *Mol. Biol. Evol.* **31**, 2735-2740  
630 (2014).
- 631 60 von Dassow, P. & Montresor, M. Unveiling the mysteries of phytoplankton life cycles:  
632 patterns and opportunities behind complexity. *J. Plankton Res.* **33**, 3-12 (2010).
- 633 61 Hanley-Bowdoin, L., Bejarano, E. R., Robertson, D. & Mansoor, S. Geminiviruses: masters at  
634 redirecting and reprogramming plant processes. *Nat. Rev. Microbiol.* **11**, 777-788 (2013).
- 635 62 He, Y.-Z. *et al.* A plant DNA virus replicates in the salivary glands of its insect vector via  
636 recruitment of host DNA synthesis machinery. *Proceedings of the National Academy of*  
637 *Sciences* **117**, 16928-16937 (2020).
- 638 63 Romero, H. *et al.* Single molecule tracking reveals functions for RarA at replication forks but  
639 also independently from replication during DNA repair in *Bacillus subtilis*. *Sci. Rep.* **9**, 1997  
640 (2019).
- 641 64 Yoshimura, A., Seki, M. & Enomoto, T. The role of WRNIP1 in genome maintenance. *Cell*  
642 *Cycle* **16**, 515-521 (2017).
- 643 65 Parajuli, S. *et al.* Human ribonuclease H1 resolves R-loops and thereby enables progression  
644 of the DNA replication fork. *J. Biol. Chem.* **292**, 15216-15224 (2017).

- 645 66 Posse, V. *et al.* RNase H1 directs origin-specific initiation of DNA replication in human  
646 mitochondria. *PLoS Genet.* **15**, e1007781 (2019).
- 647 67 Tadokoro, T. & Kanaya, S. Ribonuclease H: molecular diversities, substrate binding domains,  
648 and catalytic mechanism of the prokaryotic enzymes. *FEBS J.* **276**, 1482-1493 (2009).
- 649 68 Steenwyk, J. L. *et al.* Extensive loss of cell-cycle and DNA repair genes in an ancient lineage  
650 of bipolar budding yeasts. *PLoS Biol.* **17**, e3000255 (2019).
- 651 69 Grohme, M. A. *et al.* The genome of *Schmidtea mediterranea* and the evolution of core  
652 cellular mechanisms. *Nature* **554**, 56-61 (2018).
- 653 70 Sekelsky, J. DNA repair in *Drosophila*: mutagens, models, and missing genes. *Genetics* **205**,  
654 471-490 (2017).
- 655 71 Corradi, N. Microsporidia: eukaryotic intracellular parasites shaped by gene loss and  
656 horizontal gene transfers. *Annu. Rev. Microbiol.* **69**, 167-183 (2015).
- 657 72 Galindo, L. J. *et al.* Evolutionary genomics of *Metchnikovella incurvata* (Metchnikovellidae):  
658 an early branching microsporidium. *Genome Biol. Evol.* **10**, 2736-2748 (2018).
- 659 73 Albalat, R. & Canestro, C. Evolution by gene loss. *Nat. Rev. Genet.* **17**, 379-391 (2016).
- 660 74 Roger, A. J., Kolisko, M. & Simpson, A. G. B. in *Evolution of Virulence in Eukaryotic*  
661 *Microbes* 44-69 (2013).
- 662 75 Rancati, G. *et al.* Aneuploidy underlies rapid adaptive evolution of yeast cells deprived of a  
663 conserved cytokinesis motor. *Cell* **135**, 879-893 (2008).
- 664 76 Mendonca, A. G., Alves, R. J. & Pereira-Leal, J. B. Loss of genetic redundancy in reductive  
665 genome evolution. *PLoS Comput. Biol.* **7**, e1001082 (2011).
- 666 77 Lindahl, T. & Wood, R. D. Quality control by DNA repair. *Science* **286**, 1897-1905 (1999).

- 667 78 Tumova, P., Uzlikova, M., Jurczyk, T. & Nohynkova, E. Constitutive aneuploidy and  
668 genomic instability in the single-celled eukaryote *Giardia intestinalis*. *Microbiologyopen* **5**,  
669 560-574 (2016).
- 670 79 Kulda, J., Nohýnková, E. & Čepička, I. in *Handbook of the Protists* (eds John M. Archibald  
671 *et al.*) 1-32 (Springer International Publishing, 2017).
- 672 80 Bogenschutz, N. L., Rodriguez, J. & Tsukiyama, T. Initiation of DNA replication from non-  
673 canonical sites on an origin-depleted chromosome. *PLoS One* **9**, e114545 (2014).
- 674 81 Bell, S. P. Rethinking origin licensing. *Elife* **6** (2017).
- 675 82 Shibata, E. *et al.* Two subunits of human ORC are dispensable for DNA replication and  
676 proliferation. *Elife* **5** (2016).
- 677 83 Park, S. Y. & Asano, M. The origin recognition complex is dispensable for endoreplication in  
678 *Drosophila*. *Proc. Natl. Acad. Sci. U.S.A.* **105**, 12343-12348 (2008).
- 679 84 Okano-Uchida, T. *et al.* Endoreduplication of the mouse genome in the absence of ORC1.  
680 *Genes Dev.* **32**, 978-990 (2018).
- 681 85 Theis, J. F. *et al.* The DNA damage response pathway contributes to the stability of  
682 chromosome III derivatives lacking efficient replicators. *PLoS Genet.* **6**, e1001227 (2010).
- 683 86 Hawkins, M., Malla, S., Blythe, M. J., Nieduszynski, C. A. & Allers, T. Accelerated growth  
684 in the absence of DNA replication origins. *Nature* **503**, 544-547 (2013).
- 685 87 Bader, A. S., Hawley, B. R., Wilczynska, A. & Bushell, M. The roles of RNA in DNA  
686 double-strand break repair. *Br. J. Cancer* **122**, 613-623 (2020).
- 687 88 Gillespie, K. A., Mehta, K. P., Laimins, L. A. & Moody, C. A. Human papillomaviruses  
688 recruit cellular DNA repair and homologous recombination factors to viral replication centers.  
689 *J. Virol.* **86**, 9520-9526 (2012).

- 690 89 Kogoma, T. Stable DNA replication: interplay between DNA replication, homologous  
691 recombination, and transcription. *Microbiol. Mol. Biol. Rev.* **61**, 212-238 (1997).
- 692 90 Gerhold, J. M. *et al.* Replication intermediates of the linear mitochondrial DNA of *Candida*  
693 *parapsilosis* suggest a common recombination based mechanism for yeast mitochondria. *J.*  
694 *Biol. Chem.* **289**, 22659-22670 (2014).
- 695 91 Lydeard, J. R., Jain, S., Yamaguchi, M. & Haber, J. E. Break-induced replication and  
696 telomerase-independent telomere maintenance require Pol32. *Nature* **448**, 820-823 (2007).
- 697 92 Sung, P. & Klein, H. Mechanism of homologous recombination: mediators and helicases take  
698 on regulatory functions. *Nat. Rev. Mol. Cell Biol.* **7**, 739-750 (2006).
- 699 93 Malkova, A., Naylor, M. L., Yamaguchi, M., Ira, G. & Haber, J. E. RAD51-dependent break-  
700 induced replication differs in kinetics and checkpoint responses from RAD51-mediated gene  
701 conversion. *Mol. Cell. Biol.* **25**, 933-944 (2005).
- 702 94 Jain, S. *et al.* A recombination execution checkpoint regulates the choice of homologous  
703 recombination pathway during DNA double-strand break repair. *Genes Dev.* **23**, 291-303  
704 (2009).
- 705 95 Sakofsky, C. J. & Malkova, A. Break induced replication in eukaryotes: mechanisms,  
706 functions, and consequences. *Crit. Rev. Biochem. Mol. Biol.* **52**, 395-413 (2017).
- 707 96 Wilson, M. A. *et al.* Pif1 helicase and Poldelta promote recombination-coupled DNA  
708 synthesis via bubble migration. *Nature* **502**, 393-396 (2013).
- 709 97 Vijayraghavan, S., Tsai, F. L. & Schwacha, A. A checkpoint-related function of the MCM  
710 replicative helicase is required to avert accumulation of RNA:DNA hybrids during S-phase  
711 and ensuing DSBs during G2/M. *PLoS Genet.* **12**, e1006277 (2016).



- 712 98 Drissi, R. *et al.* Destabilization of the minichromosome maintenance (MCM) complex  
713 modulates the cellular response to DNA double strand breaks. *Cell Cycle* **17**, 2593-2609  
714 (2018).
- 715 99 Chen, Y. H. *et al.* Transcription shapes DNA replication initiation and termination in human  
716 cells. *Nat. Struct. Mol. Biol.* **26**, 67-77 (2019).
- 717 100 Keskin, H., Meers, C. & Storici, F. Transcript RNA supports precise repair of its own DNA  
718 gene. *RNA Biol.* **13**, 157-165 (2016).
- 719 101 Bartek, J., Lukas, C. & Lukas, J. Checking on DNA damage in S phase. *Nat. Rev. Mol. Cell*  
720 *Biol.* **5**, 792-804 (2004).
- 721 102 Das, M., Singh, S., Pradhan, S. & Narayan, G. MCM paradox: abundance of eukaryotic  
722 replicative helicases and genomic integrity. *Mol. Biol. Int.* **2014**, 574850 (2014).
- 723 103 Sasaki, T. & Gilbert, D. M. The many faces of the origin recognition complex. *Curr. Opin.*  
724 *Cell Biol.* **19**, 337-343 (2007).
- 725 104 Borlado, L. R. & Mendez, J. CDC6: from DNA replication to cell cycle checkpoints and  
726 oncogenesis. *Carcinogenesis* **29**, 237-243 (2008).
- 727 105 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina  
728 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 729 106 Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting  
730 and repeat separation. *Genome Res.* **27**, 722-736 (2017).
- 731 107 Loman, N. J., Quick, J. & Simpson, J. T. A complete bacterial genome assembled de novo  
732 using only nanopore sequencing data. *Nat. Methods* **12**, 733-735 (2015).
- 733 108 Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial genome  
734 assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**, e1005595 (2017).

- 735 109 Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078-  
736 2079 (2009).
- 737 110 Pedersen, B. S. & Quinlan, A. R. Mosdepth: quick coverage calculation for genomes and  
738 exomes. *Bioinformatics* **34**, 867-868 (2018).
- 739 111 Haas, B. J. *et al.* De novo transcript sequence reconstruction from RNA-seq using the Trinity  
740 platform for reference generation and analysis. *Nat. Protoc.* **8**, 1494-1512 (2013).
- 741 112 Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory  
742 requirements. *Nat. Methods* **12**, 357-360 (2015).
- 743 113 Lomsadze, A., Burns, P. D. & Borodovsky, M. Integration of mapped RNA-seq reads into  
744 automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res.* **42**, e119 (2014).
- 745 114 Stanke, M., Schoffmann, O., Morgenstern, B. & Waack, S. Gene prediction in eukaryotes  
746 with a generalized hidden Markov model that uses hints from external sources. *BMC*  
747 *Bioinformatics* **7**, 62 (2006).
- 748 115 Haas, B. J. *et al.* Improving the *Arabidopsis* genome annotation using maximal transcript  
749 alignment assemblies. *Nucleic Acids Res.* **31**, 5654-5666 (2003).
- 750 116 Waterhouse, R. M. *et al.* BUSCO applications from quality assessments to gene prediction  
751 and phylogenomics. *Mol. Biol. Evol.* **35**, 543-548 (2018).
- 752 117 Slater, G. S. & Birney, E. Automated generation of heuristics for biological sequence  
753 comparison. *BMC Bioinformatics* **6**, 31 (2005).
- 754 118 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,  
755 357-359 (2012).

- 756 119 Augusto Correa Dos Santos, R., Goldman, G. H. & Riano-Pachon, D. M. ploidyNGS:  
757 visually exploring ploidy with next generation sequencing data. *Bioinformatics* **33**, 2575-2576  
758 (2017).
- 759 120 Yamada, K. D., Tomii, K. & Katoh, K. Application of the MAFFT sequence alignment  
760 program to large data-reexamination of the usefulness of chained guide trees. *Bioinformatics*  
761 **32**, 3246-3251 (2016).
- 762 121 Tan-Wong, S. M., Dhir, S. & Proudfoot, N. J. R-Loops promote antisense transcription across  
763 the mammalian genome. *Mol. Cell* **76**, 600-616 e606 (2019).
- 764 122 Mazina, O. M. *et al.* Replication protein A binds RNA and promotes R-loop formation. *J.*  
765 *Biol. Chem.* **295**, 14203-14213 (2020).
- 766 123 Saldivar, J. C., Cortez, D. & Cimprich, K. A. The essential kinase ATR: ensuring faithful  
767 duplication of a challenging genome. *Nat. Rev. Mol. Cell Biol.* **18**, 622-636 (2017).
- 768 124 Longhese, M. P., Plevani, P. & Lucchini, G. Replication factor A is required in vivo for DNA  
769 replication, repair, and recombination. *Mol. Cell. Biol.* **14**, 7884-7890 (1994).
- 770 125 Domingo-Prim, J., Bonath, F. & Visa, N. RNA at DNA double-strand breaks: the challenge of  
771 dealing with DNA:RNA hybrids. *Bioessays* **42**, e1900225 (2020).
- 772 126 Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUASt: quality assessment tool for  
773 genome assemblies. *Bioinformatics* **29**, 1072-1075 (2013).

774

## 775 **Acknowledgments**

776 The majority of this work was supported by a Foundation grant FRN-142349, awarded to A.J.R. by  
777 the Canadian Institutes of Health Research. Archibald Lab contributions to this study were supported  
778 by a Discovery Grant from the Natural Sciences and Engineering Research Council of Canada

779 (RGPIN 05871-2014). E.C.T. is supported by a Herchel Smith Postdoctoral Fellowship at the  
780 University of Cambridge.

781 **Author contributions**

782 D.E.S-L and A.J.R. conceived the study. D.E.S-L and B.A.C. scripted *in house* programs and led the  
783 bioinformatics workflow. J.J-H and M.K. grew cultures, extracted nucleic acids, and carried out *in*  
784 *house* sequencing. D.E.S-L., B.A.C., E.C.T., Z.Y, J.S.S-L., L.G-L., G.J.P.L.K, J.M.A., A.G.B.S. and  
785 A.J.R. analyzed and manually curated the genomic data. E.C.T. and D.E.S-L made the figures.  
786 D.E.S-L and A.J.R. led the writing of the manuscript with input from all authors. All documents were  
787 edited and approved by all authors.

788 **Competing interests**

789 Authors declare no competing interests.

790 **Additional information**

791 Supplementary Information (also containing legends for Supplementary Table 1 and Supplementary  
792 Figures 1 – 10)

793

## 794 **Figure legends**

795 **Figure 1** The distribution of core molecular systems in the replisome and DNA repair across  
796 eukaryotic diversity. A schematic global eukaryote phylogeny is shown on the left with classification  
797 of the major metamonad lineages indicated at right. **A)** The Replisome. Reduction of the replication  
798 machinery complexity and extensive loss of the Orc1-6 subunits are observed in metamonad lineages,  
799 including the unexpected loss of the highly conserved ORC complex and Cdc6 in *Carpediemonas*.  
800 Most metamonad Orc1 and Cdc6 homologs were conservatively named as ‘Orc1/Cdc6-like’ as they  
801 are very divergent, do not have the typical domain architecture and, in phylogenetic reconstructions,  
802 they form clades separate from the main eukaryotic groups, preventing confident orthology  
803 assignments (**Supplementary Figure 1**). Numbers within subunits represent the number of copies and  
804 are only presented for ORC components, additional information in **Supplementary Table 1**. The  
805 polymerase epsilon ( $\epsilon$ ) is composed of 4 subunits, but we included the interacting protein Chrac1  
806 (depicted as ‘4!’ in the figure) as its HMM retrieves the polymerase delta subunit Dbp3 from *S.*  
807 *cerevisiae*. \*Firing and elongation factors, \*\*Protein fusion between the catalytic subunit and subunit 2  
808 of DNA polymerase  $\epsilon$ . + Preaxostyla, ++ Parabasalida, +++*Carpediemonas*-Like Organisms. **B)**  
809 Predicted *Carpediemonas* replisome overlaid on a typical eukaryotic replisome. Origin recognition  
810 (ORC), Cdc6 and replication progression (RPC) complexes are depicted. Grey colour represents the  
811 absence of typical eukaryotic proteins in *C. membranifera* replisome.

## 812 **Figure 2. Pif1 protein family expansion**

813 Pif1 helicase family tree. Three clades are highlighted: at the top, a Pif1-like clade encompassing some  
814 metamonads and at the bottom a *Carpediemonas*-specific Pif1-like clade. The third clade shows the  
815 typical Pif1 orthologs encompassing fornicates. The maximum-likelihood tree was inferred under the  
816 LG+PMSF(C60)+F+  $\Gamma$  model using 100 bootstraps based on an alignment length of 265 sites. The tree

817 was midpoint-rooted and the support values on the branches correspond to SH-aLRT/aBayes/standard  
818 bootstrap (values below 80/0.8/80 are not shown). The scale bar shows the inferred number of amino  
819 acid substitutions per site.

820 **Figure 3 Radical reduction of ancestral kinetochore network complexity in *Carpediemonas***  
821 **species. A)** Schematic of canonical mitotic cell cycle progression in eukaryotes. During mitosis,  
822 duplicated chromosomes each attach to microtubules (MTs) emanating from opposite poles of the  
823 spindle apparatus, in order to be segregated into two daughter cells. Kinetochores (KTs) are built upon  
824 centromeric DNA to attach microtubules to chromosomes. To prevent precocious chromosome  
825 segregation, unattached KT's signal to halt cell cycle progression (STOP), a phenomenon known as the  
826 Spindle Assembly Checkpoint (SAC). The SAC entails the inhibition of the Anaphase Promoting  
827 Complex/Cyclosome (APC/C), a multi-subunit E3 ubiquitin ligase complex that drives the entry of  
828 mitotic cells into anaphase by promoting the proteolysis of its substrates. Once all KT's are correctly  
829 attached to spindle MTs and aligned in the middle of the cell (metaphase), the APC/C is released, its  
830 substrates are degraded, and chromosome segregation is initiated (anaphase). **B)** Cartoon of the  
831 molecular makeup of a single KT unit that was likely present in Last Eukaryotic Common Ancestor  
832 (LECA). Colours indicate the various functional complexes and structures. The primary KT structure  
833 is provided by the Constitutive Centromere Associated Network (CCAN; yellow), which is built upon  
834 centromeric chromatin that contains Centromere protein A (CenpA; orange), a centromere-specific  
835 Histone H3. During mitosis the CCAN recruits the Mis12 complex (linker; light green), which  
836 provides a platform for the recruitment of the SAC signalling (light blue) and microtubule-interacting  
837 complexes. The Chromosomal Passenger Complex (CPC; dark purple) localizes at the inner  
838 centromere and harbours a kinase (aurora) that regulates microtubule attachments. Unattached KT's  
839 catalyse the production of a diffusible cytosolic inhibitor of the APC/C, known as the mitotic

840 checkpoint complex (MCC), which captures the mitotic APC/C co-activator Cdc20. Initial KT-MT  
841 encounters are driven by the kinesin Centromere protein E (CenPE; pink), which binds MTs at the  
842 lateral sides. The Ndc80 complex (dark red) constitutes the main end-on MT binding activity of KTs.  
843 To facilitate the tracking of the plus-end (+) of MT during anaphase, eukaryotes utilize two different  
844 complexes: Dam (light purple; likely not present in LECA) and Ska (red). Once KTs are bound by  
845 MTs, SAC signalling proteins are removed and the SAC is turned off. **C)** Reconstruction of the  
846 evolution of the KT and mitotic signalling in eukaryotes based on protein presence-absence patterns  
847 reveals extensive reduction of ancestral KT complexity and loss of the SAC in most metamonad  
848 lineages, including the striking loss of the highly conserved core MT-binding activity of the KT  
849 (Ndc80) in *Carpodomonas*. On top/bottom of panel C: the number of components per complex and  
850 different structural parts of the KT, SAC signalling and the APC/C. Middle: presence/absence matrix  
851 of KT, SAC and APC/C complexes; one circle per complex, colours correspond to panel A & B; grey  
852 indicates its (partial) loss (for a complete overview see **Supplementary Table 1, Supplementary Fig.**  
853 **4**). The red STOP sign indicates the likely presence of a functional SAC response (see for discussion  
854 **Supplementary Fig. 6**). On the left: cartoon of a phylogenetic tree of metamonad and other selected  
855 eukaryotic species with a projection of the loss and gain events on each branch. Specific loss events of  
856 kinetochore and SAC genes in specific lineages are highlighted in colour.

857 **Figure 4. Hypothesis for Dmc1-dependent DNA replication in *Carpodomonas*.**

858 **A)** R-loop stimulated sense and antisense transcription<sup>121</sup> in a highly transcribed locus results in a  
859 DNA break, triggering DSB checkpoint control systems to assemble HR complexes and the replication  
860 proteins near the lesion<sup>11,37,122-124</sup>. Once the damage is processed into a DSB, end resection by  
861 Mre11/Rad50 creates a 3' overhang and the strands are coated with Replication protein A (RPA),  
862 while resected ends are coated with the recombinase Dmc1. **B)** A recombination checkpoint decides

863 the HR sub-pathway to be used<sup>94</sup>, then strand invasion of a broken end is initiated into a transcript-  
864 RNA or -cDNA template<sup>39,100,125</sup> followed by the initiation and progression of DNA synthesis with the  
865 aid of Pif1 helicase\*. This leads to the establishment of a double Holliday Junction (HJ) which can be  
866 resolved by endonucleases (*e.g.*, Mus81, Flap, Mlh1/Mlh3). The lack of Chk1 may result in mis-  
867 segregation caused by aberrant processing of DNA replication intermediates by Mus81<sup>48</sup>. Given the  
868 shortness of the RNA or cDNA template, most possible HJ resolutions, except for the one depicted in  
869 the figure, would lead to the loss of chromosome fragments. The HJ resolution shown would allow  
870 steps shown in panel C. C) A multimeric *Carpediemonas* Pif1-like helicase is bound to the repaired  
871 DNA as well as to the template. Here, the shortness of the template could resemble a replication  
872 intermediate that could prompt the recruitment of MCM, following the addition of the replisome  
873 proteins and establishing a fully functional replication fork (Dark blue fragments on 3' ends of the  
874 bottom figure represent Okazaki fragments).

875 \*Notes: Polymerases  $\alpha$  and  $\delta$  are able to incorporate the correct nucleotides using RNA template<sup>40</sup>;  
876 RNase H2 would excise ribonucleosides and replace the correct nucleotide.

877

878 **Tables**

879



**Table 1 Summary statistics of nuclear genomes of Metamonada species.**

Description	<i>Trichomonas vaginalis</i>	<i>Monocercomonoides exilis</i>	<i>Carpediemonas membranifera</i>	<i>Carpediemonas frisia</i>	<i>Kipferlia bialata</i>	<i>Spirotrunculus salmonicida</i>	<i>Trepomonas PCI*</i>	<i>Giardia intestinalis A 50803</i>	<i>Giardia intestinalis B 50581</i>	<i>Giardia muris</i>
Genome size (Mb)	176.4	74.7	24.3	12.4	51.0	12.9		11.7	11.0	9.7
Contigs/Scaffolds	64764	2095	68	3232	11563	233		211	2931	59
N50 (bp)	27258	71440	906349	9593	10488	150829		2,762,469	34,141	2,398,647
GC (%)	32.7	37.4	57.19	58.6	47.8	33.5		49.0	46.5	54.71
No. of predicted genes	94255	16780	11883	5695	17389	8354	7980	5901	4470	4936
No. BUSCO genes (percentage)	223 (91)	224 (91)	217 (89)	184 (75)	207(84)	152 (62)	147 (60)	168 (69)	169 (69)	173 (71)
SINEs (%)	0.07	0	0.2	0	0	0.16		0	0.07	0.03
LINEs (%)	0.06	0.79	8.07	0	1.08	0		0.98	0.12	0.59
LTR Elements (%)	0.52	4.44	20.6	0.4	1.34	0.29		0	0	0.79
DNA Elements (%)	50.66	9.96	0.9	0.07	22.7	0.2		0	0	0
Unclassified (%)	15.41	21.76	14.9	4.97	1.22	5.64		8.64	6.76	11.77
Total interspersed repeats (%)	66.72	36.94	43.97	4.45	26.38	6.3		9.62	6.95	13.18
Simple Repeats (%)	0.21	1.03	0.24	0	0.1	0		0	0	0

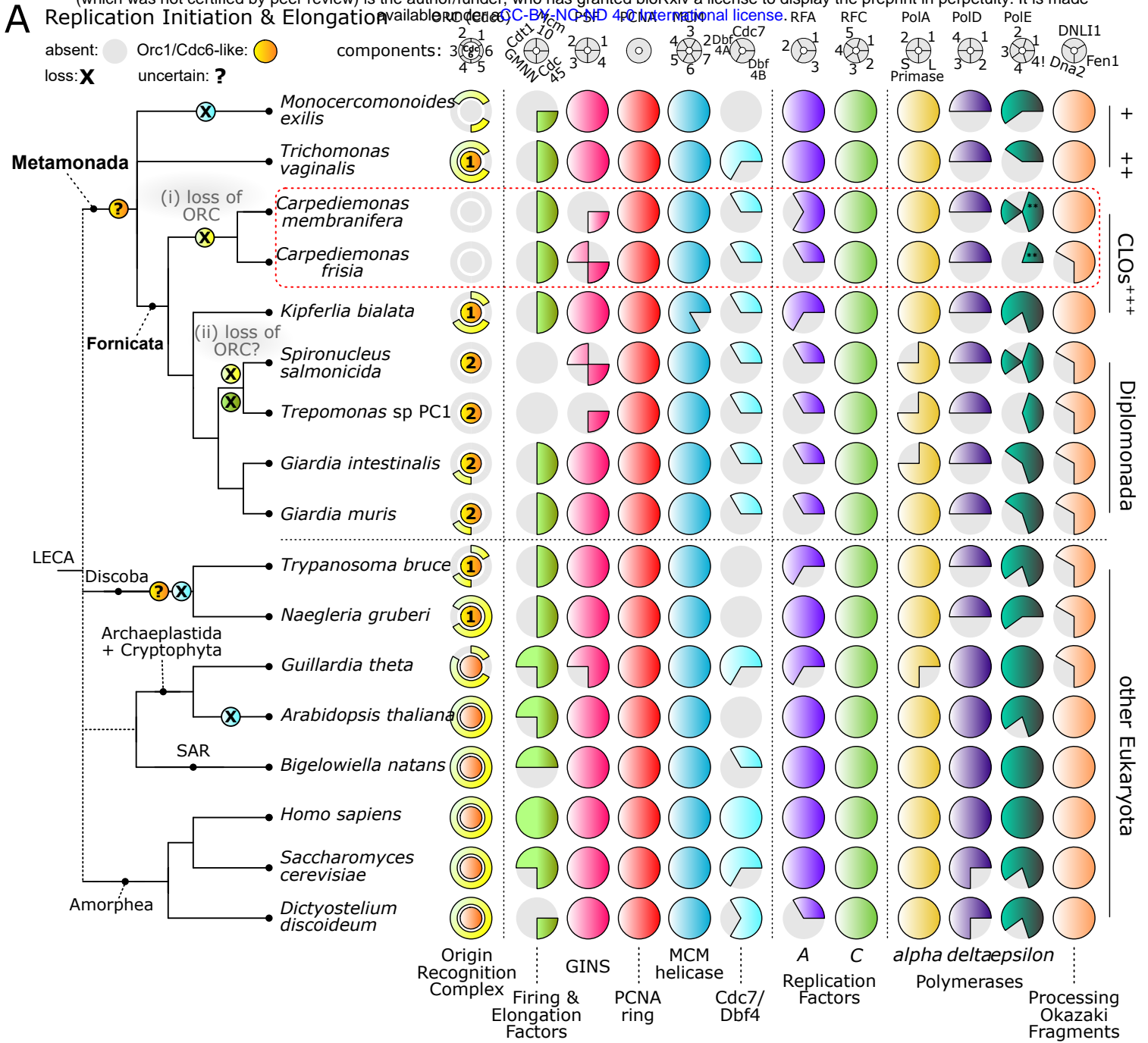
**All the statistics were recalculated with Quast<sup>126</sup> for completion as not all of these were originally reported, and the BUSCO**

**reference protein set corresponds to a maximum of 245 proteins.**

\*transcriptome data only

**Figure 1**

bioRxiv preprint doi: <https://doi.org/10.1101/2021.03.14.435266>; this version posted March 15, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



**B *Carpediemonas membranifera* Replisome**

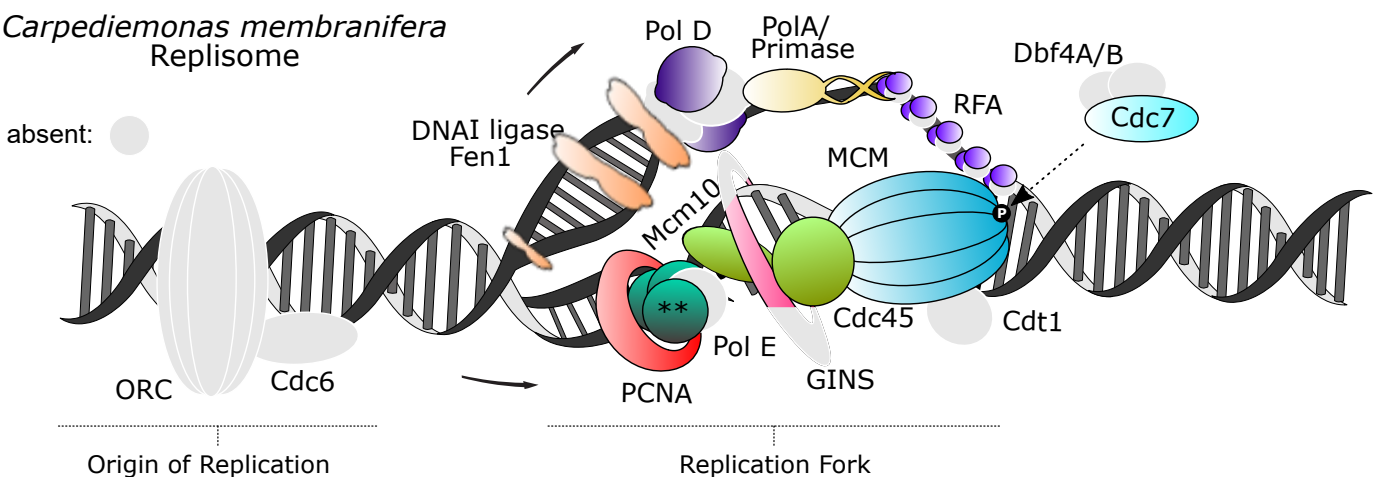
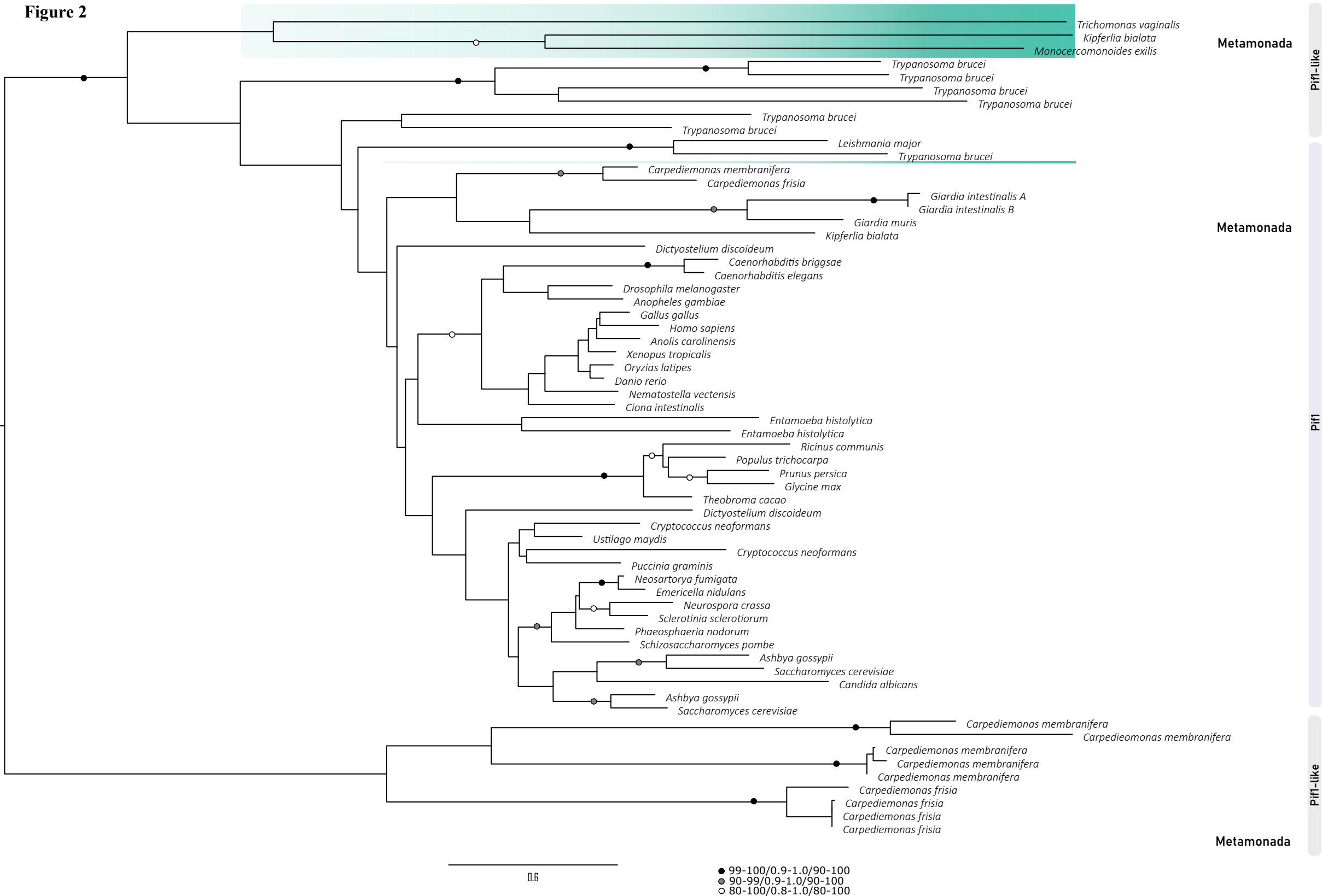
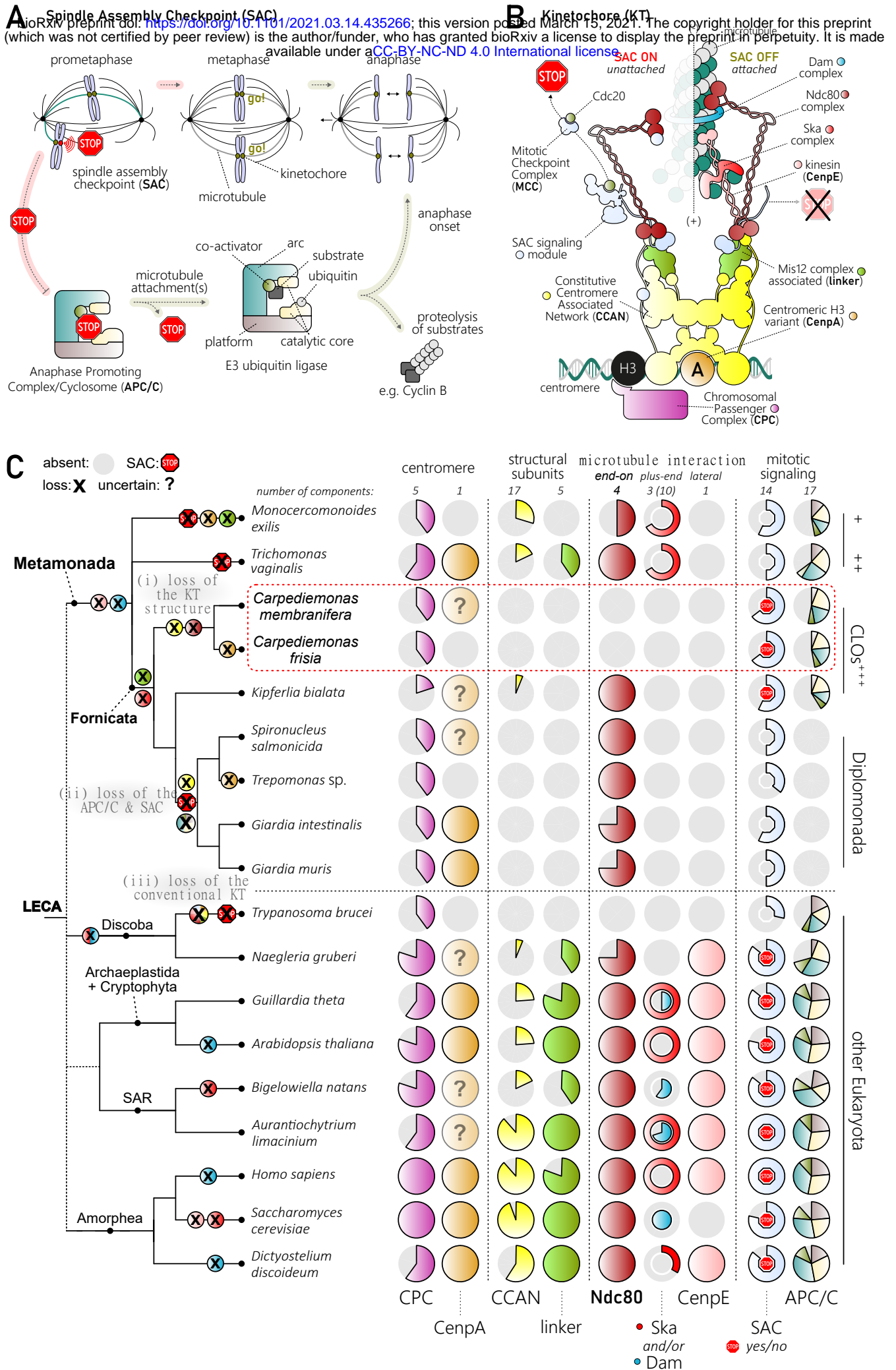


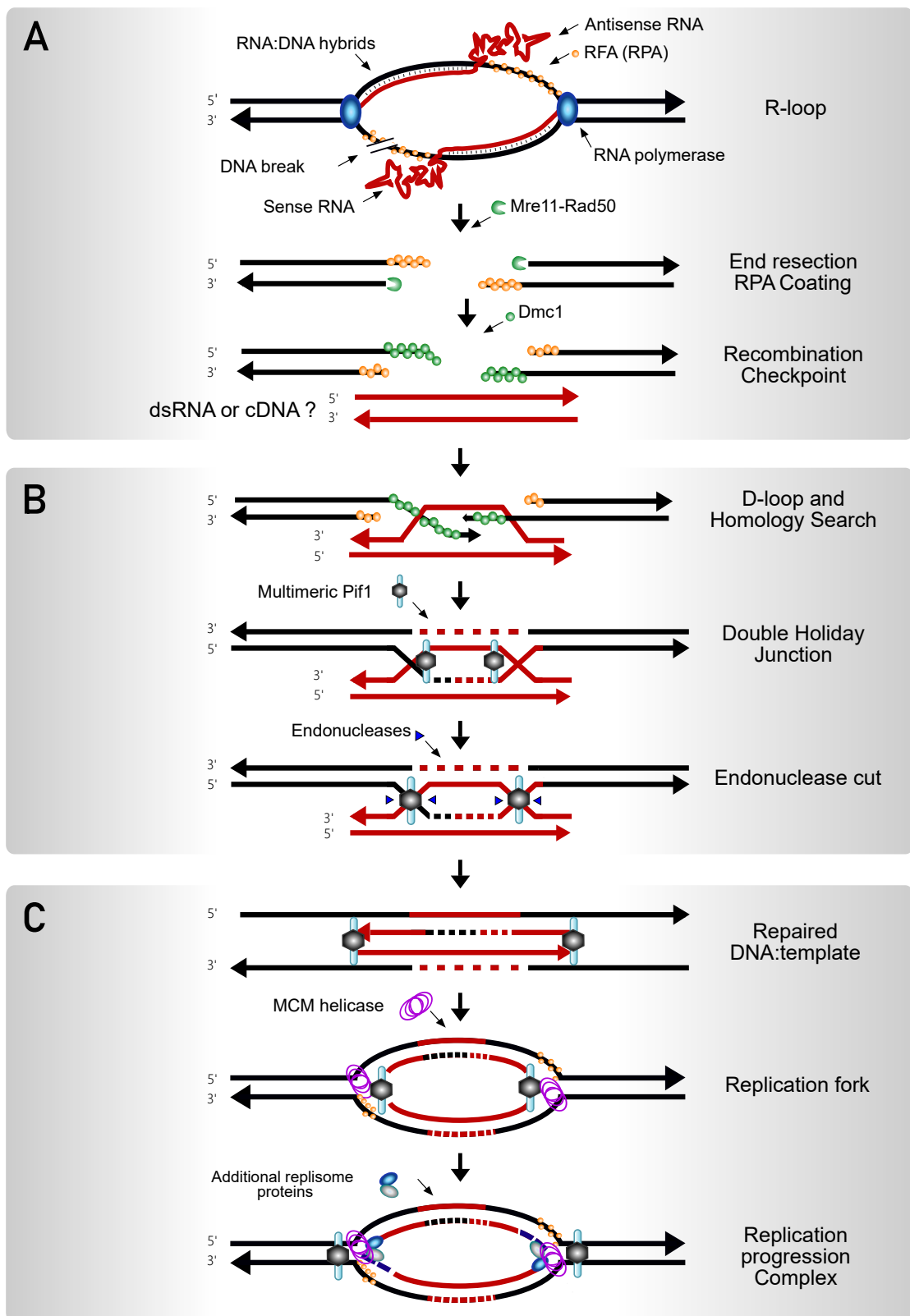
Figure 2



**Figure 3**



**Figure 4**



1	<b>Supplementary information</b>
2	<b>Table of contents</b>
3	<b>A. Supplementary methods</b>
4	<b>A1. Culturing and DNA isolation</b>
5	<b>A2. Genome size and completeness using BUSCO and a phylogenetic guided approach</b>
6	<b>A3. Taxa selected for the comparative genomic analysis.</b>
7	<b>A4. Additional strategies used to search for ORC, Cdc6 ad Ndc80 proteins.</b>
8	
9	<b>B. Supplementary results</b>
10	<b>B1. BUSCO completeness.</b>
11	<b>B2. Additional search strategies to find missing proteins.</b>
12	<b>B3. DNA replication streamlining in nucleomorphs of chlorarachniophytes and cryptophytes</b>
13	<b>B4. Acquisition of Endonuclease IV, RarA and RNase H1 by lateral gene transfer</b>
14	<b>C. Supplementary discussion</b>
15	<b>C1. BUSCO incompleteness</b>
16	<b>D. Supplementary references</b>
17	<b>E. Supplementary Figure legends</b>
18	<b>F. Supplementary Table legends</b>
19	
20	<b>A. Supplementary methods</b>
21	<b>A1. Culturing and DNA isolation</b>

22 Sequencing of *C. membranifera* BICM strain was done with Illumina short paired-end and long  
23 MinION read technologies. The Illumina sequencing employed DNA from a monoxenic culture  
24 grown in 50 ml Falcon tubes in F/2 media enriched with the bacterium *Shewanella frigidimarina* as  
25 food. DNA was isolated from a total of two litres of culture using a salt extraction protocol followed  
26 by CsCl gradient centrifugation. RNA was also extracted from these cultures using TRIzol  
27 (Invitrogen, USA), following the manufacturer's instructions. For MinION sequencing, *C.*  
28 *membranifera* was grown in sterile filtered 50% natural sea water media with 3% LB with either  
29 *Shewanella sp* or *Vibrio sp.* isolate JH43 as food. Cell cultures were harvested at peak density by  
30 centrifugation at 500×g, 8 min, 20 °C. The cells were resuspended in sterile-filtered spent growth  
31 media (SFSGM) and centrifuged again at 500×g, 8 min, 20 °C. The cell pellets were resuspended in  
32 1.5 mL SFSGM, layered on top of 9 mL Histopaque®-1077 (Sigma-Aldrich) and centrifuged at  
33 2000×g, 20 min, 20 °C. The protists were recovered from the media:Histopaque interface by  
34 pipetting, diluted in 10 volumes of SFSGM and centrifuged 500×g, 8 min, 20 °C. High molecular  
35 weight DNA was extracted using MagAttract HMW DNA Kit (Qiagen, Cat No. 67563), purified with  
36 GenomicTip 20/G (Qiagen, Cat No. 10223) and resuspended in 5 mM Tris-HCl (pH 8.5).

## 37 **A2. Genome size and completeness using BUSCO and a phylogeny-guided approach**

38 The BUSCO approach<sup>1</sup> was prone to false negative predictions with our dataset because of the  
39 extreme divergence of metamonad homologs. Therefore, the completeness of the BUSCO set was re-  
40 assessed with a phylogeny-guided search. For this, we eliminated 31 proteins associated with  
41 mitochondria or mitochondrion- related organelles (MROs) as Metamonada have reduced or no  
42 MROs<sup>2</sup>, and employed taxa-enriched Hidden Markov Model (HMM) searches to account for  
43 divergence between the remaining 272 proteins and the studied taxa. In brief: BLASTp was carried  
44 out using the 272 BUSCO proteins as queries for finding their orthologues in a local version of the

45 PANTHER 14.0 database<sup>3</sup> to enable the identification of the most likely Panther subfamily HMM  
46 and its annotation. Then, each corresponding subfamily HMM was searched for in the predicted  
47 proteomes with an e-value cut-off of  $1 \times 10^{-1}$  with HMMER v3.1b2<sup>4</sup>. In cases where these searches did  
48 not produce any result, a broader search was run using the HMM of the Panther family with  $1 \times 10^{-3}$  as  
49 e-value cut-off. Five best hits for each search were retrieved from each proteome, aligned to the  
50 corresponding Panther subfamily or family sequences with MAFFT v7.310<sup>5</sup> and phylogenetic  
51 reconstructions were carried out using IQ-TREE v1.6.5<sup>6</sup> under the LG+C60+F+ $\Gamma$  model with  
52 ultrafast bootstrapping (1000 replicates). Protein domain architectures were visualized by mapping  
53 the respective Pfam accessions onto trees using ETE tools v3.1.1<sup>7</sup>.

54

### 55 **A3. Taxa selected for comparative genomic analysis.**

56 Our analyses included the publicly available genomes and predicted proteomes of *Trichomonas*  
57 *vaginalis* G3 (Parabasalia, [www.trichdb.org](http://www.trichdb.org)), *Monocercomonoides exilis* (Preaxostyla,  
58 [www.protistologie.cz/hampllab](http://www.protistologie.cz/hampllab)), the free-living fornicates *Carpediemonas frisia*<sup>8</sup> (i.e., metagenomic  
59 bin and predicted proteome), *Carpediemonas membranifera* (reported here) and *Kipferlia bialata*<sup>9</sup>,  
60 plus the parasitic diplomonad fornicates: *Giardia intestinalis* Assemblages A and B, *Giardia muris*,  
61 *Spironucleus salmonicida*-ATCC50377 ([www.giardiadb.org](http://www.giardiadb.org)) and *Trepomonas* PC1<sup>10</sup> –the latter was  
62 only available as a transcriptome. We also included a set of genomes that are broadly representative  
63 of eukaryote diversity, such as *Homo sapiens* GRCh38, *Saccharomyces cerevisiae* S288C,  
64 *Arabidopsis thaliana* TAIR10, *Dictyostelium discoideum* AX4, *Trypanosoma brucei* TREU927-rel28  
65 ([www.uniprot.org](http://www.uniprot.org)), *Naegleria gruberi* NEG-M ([www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)), *Guillardia theta* and  
66 *Bigelowiella natans* ([www.genome.jgi.doe.gov/portal/](http://www.genome.jgi.doe.gov/portal/)).

67 Additional analyzed genomes were those of the microsporidia *Encephalitozoon intestinalis*



68 ATCC 50506 (ASM14646v1), *E. cuniculi* GB-M1 (ASM9122v2) and *Trachipleistophora hominis*  
69 (ASM31613v1), the yeasts *Hanseniaspora guilliermondii* (ASM491977v1), *Hanseniaspora opuntiae*  
70 (ASM174979v1), *Hanseniaspora osmophila* (ASM174704v1), *Hanseniaspora uvarum*  
71 (ASM174705v1) and *Hanseniaspora valbyensis* NRRL Y-1626 (GCA\_001664025.1), *Tritrichomonas*  
72 *foetus* (ASM183968v1), the nucleomorphs of *Hemiselmis andersenii* (ASM1864v1), *Cryptomonas*  
73 *paramecium* (ASM19445v1), *Chroomonas mesostigmatica* (ASM28609v1), *Guillardia theta*  
74 (ASM297v1), *Lotharella vacuolata* (AB996599–AB996601), *Amorphochlora amoebiformis*  
75 (AB996602–AB996604) and *Bigelowiella natans* (ASM245v1), the corals *Galaxea fascicularis*,  
76 *Fungia* sp., *Goniastrea aspera*, *Acropora tenuis* and the coral endosymbionts *Symbiodinium kawagutii*  
77 and *Symbiodinium goreau*<sup>11,12</sup>.

#### 78 **A4. Additional strategies used to search for ORC, Cdc6 and Ndc80 proteins.**

79 Strategies included enriched HMMs as mentioned in the main text and HMMs for individual Pfam  
80 domains with e-value thresholds of  $1 \times 10^{-3}$ . 1) Metamonad-specific HMMs were built as described for  
81 kinetochore proteins – containing the newly found hits plus orthologs from additional publicly  
82 available metamonad proteomes or transcriptomes<sup>2,13</sup>, 2) we applied the eggNOG 4.5 profiles  
83 COG1474, COG5575, KOG2538, KOG2228, KOG2543, KOG4557, KOG4762, KOG0995,  
84 KOG4438, KOG4657 and 2S26V which encompass 2774, 495, 452, 466, 464, 225, 383, 504, 515,  
85 403 and 84 taxa, respectively, and 3) the Pfam v33.1 HMMs: PF09079 (Cdc6\_C), PF17872  
86 (AAA\_lid\_10), PF00004 (AAA+), PF13401 (AAA\_22), PF13191 (AAA\_16), PF01426 (BAH),  
87 PF04084 (Orc2), PF07034 (Orc3), PF18137 (ORC\_WH\_C), PF14629 (Orc4\_C), PF14630 (Orc5\_C),  
88 PF05460 (Orc6), PF03801 (Ndc80\_HEC), PF03800 (Nuf2), PF08234 (Spindle\_Spc25) and PF08286  
89 (Spc24). For Ndc80, Nuf2, Spc24 and Spc25 we also applied the HMMs models published in<sup>14</sup>.

91 **B. Supplementary results**

92 **B1. BUSCO completeness.**

93 A subset of 272 BUSCO proteins from the odb9 database was used for a phylogeny-guided search for  
94 divergent orthologs. This revealed that: *i*) 27 out of 272 BUSCO (9.9%) proteins are absent in all  
95 metamonads, *ii*) only 101 (~41%) of the remaining 245 proteins were shared by all metamonad  
96 proteomes, and *iii*) up to 38% are absent in all Fornicata. Metamonad genomes only contained 60% to  
97 91% of the BUSCO proteins (**Table 1, Supplementary Table 1**, Note: the BUSCO presence-absence  
98 patterns of the transcriptomic data from *Trepomonas* sp. PC1 are consistent with those of the  
99 remaining diplomonads). These analyses demonstrate that the Metamonada have secondarily lost a  
100 relatively large number of highly conserved eukaryotic proteins and, therefore, BUSCO analysis  
101 cannot be used on its own to evaluate metamonad genome completeness.

102

103 **B2. Additional search strategies to find missing proteins.**

104 Metamonad-specific HMM retrieved two candidates for Orc1/Cdc6 proteins from *C. frisia* (*i.e.*,  
105 Cfrisia\_2222, Cfrisia\_2845) and one from *C. membranifera* (*i.e.*, c4603.t1), and one Orc4 candidate  
106 from each *Carpediemonas* species (*i.e.*, Cfrisia\_2559, ds58\_16707). Further inspection of these hits  
107 showed that only the AAA+ region shared similarity among all of these proteins, which is expected  
108 as ORC and Cdc6 proteins belong to the ATPase superfamily. However, based on full protein  
109 identity, full profile composition and domain architecture, the proteins retrieved with the Orc1/Cdc6  
110 HMM were confidently annotated as Katanin P60 ATPase-containing subunit A1, Replication factor  
111 C subunits 1 and 5, and proteins retrieved with Orc4 HMM were members of the Dynein heavy chain  
112 and AAA-family ATPase families. The latter is a 744 aa protein that has a C-terminal region with no  
113 sequence similarity or amino acid profile frequencies that resembles a Orc4\_C Pfam domain from

114 other metamonads or model eukaryotes. All the additional search strategies yielded false positives in  
115 *Carpodiemonas* species, as these retrieved AAA-family members lacking sequence similarity to orc  
116 proteins, showed completely different protein domain architecture than the expected one and were  
117 associated with different functional annotation (data not shown). When reconstructing the domain  
118 architecture of ORC and Cdc6 proteins in metamonads, we noted that Fornicata Orc1/Cdc6-like  
119 proteins are remarkably smaller (*i.e.*, 1.5 to 3 times smaller) than Orc1 and Cdc6 from the model  
120 organisms and other protists used later in phylogenetic reconstruction (**Supplementary Figure 1A**  
121 **and B, Supplementary Table 1**). In most cases, the small proteins lack protein domains rendering a  
122 different domain architecture with respect to their homologs in *S. cerevisiae*, *H. sapiens*, *A. thaliana*  
123 and *T. vaginalis* (**Supplementary Figure 1A, Supplementary Table 1**). For example, Orc1 and  
124 Cdc6 paralogs in Fornicata lack BAH, and AAA\_lid10 and Cdc6\_C domains. Protein alignments  
125 show that the conserved areas of these proteins correspond to AAA+ domain that have relatively  
126 conserved Walker domains A and B (except MONOS\_13325 from *M. exilis*), with a few proteins  
127 lacking the arginine finger motif (R-finger) within the Walker B motif (**Supplementary Figure 1B**).  
128 The latter may negatively affect ATPase activity of the R-finger-less proteins. In an attempt to  
129 establish orthology, metamonad Orc1/Cdc6 candidates were used for phylogenetic reconstruction  
130 together with publicly available proteins that have reliable annotations for Orc1 and Cdc6, expected  
131 domain architecture and/or with experimental evidence of their functional activity in the replisome.  
132 Phylogenetic analysis shows that metamonad proteins form separate clades from the *bona fide* Orc1  
133 and Cdc6 sequences (**Supplementary Figure 1C**). One of these separate clades encompasses Orc1-b  
134 from *T. brucei* that has been shown to participate during DNA replication despite lacking the typical  
135 domain architecture<sup>15</sup>.

136

137 **B3. DNA replication streamlining in nucleomorphs**

138 The loss of ORC/Cdc6 accompanied by the partial retention of MCM, PCNA, Cdc45, RCF, GINS  
139 and the homologous recombination (HR) recombinase Rad51 was observed in cryptophyte and  
140 chlorarachniophyte nucleomorphs (**Supplementary table 1**). ORC and Cdc6 were found as single  
141 copies (except Orc2) in the nuclear DNA of these two groups; their predicted proteins lack obvious  
142 signal and targeting peptides which would likely prevent them from participating in a nucleus-  
143 coordinated nucleomorph replication. Hence, nucleomorph DNA replication likely occurs by HR  
144 without the assistance of ORC/Cdc6 origin-binding, but this replication might nonetheless be  
145 regulated at the transcriptional level by the nucleus as shown by<sup>16</sup>. Many of the remaining nuclear-  
146 encoded proteins involved in replication are present in more than one copy in those taxa, with several  
147 of them containing signal and transit peptides (*e.g.*, H2A, POLD, RCF1 and RFA1)<sup>16,17</sup>.

148 **B4. Acquisition of Endonuclease IV, RarA and RNase H1 by lateral gene transfer**

149 The Endonuclease IV (Apl1 in yeast) and exonuclease III (Exo III) function in the removal of  
150 abasic sites in DNA via the BER pathway. Our analyses show that *C. frisia* and *C. membranifera*  
151 have Exo III and have a prokaryotic version of Endo IV (**Supplementary Fig 8**). Interestingly, none  
152 of the parabasalids and *Giardia* spp. have an Endo IV homolog, either eukaryotic or prokaryotic. *S.*  
153 *salmonicida* and *Trepomonas* sp. PC1, by contrast, appear to encode a typical eukaryotic Endo IV.

154 The RarA (Replication-Associated Recombination protein A, also named MgsA) protein is  
155 ubiquitous in bacteria and eukaryotes (*e.g.*, homologs Msg1 in yeast and WRNIP1 in mammals) and  
156 acts in the context of collapsed replication forks<sup>18,19</sup>. *Carpediemonas* possesses a prokaryotic-like  
157 version (**Supplementary Fig 9**) that lacks the ubiquitin-binding Zn finger N-terminal domain typical  
158 of eukaryotic homologs<sup>18</sup>. No canonical eukaryotic RarAs were detected in the remaining

159 metamonads, but it appears that prokaryotic-like RarA proteins in *Giardia*, *S. salmonicida* and  
160 *Trepomonas* sp. PC1 were acquired in an independent event from that of *Carpediemonas*.

161 Both *Carpediemonas* genomes have a eukaryotic RNase H2, lack eukaryotic RNase H1 but  
162 encode up to two copies of a prokaryotic-like RNase H1 (**Supplementary Fig. 10**) which do not  
163 have the typical eukaryotic HBD domain<sup>20</sup>. The HBD domain is thought to be responsible for the  
164 higher affinity of this protein for DNA/RNA duplexes rather than for dsRNA<sup>21,22</sup>. All prokaryotic-  
165 like RNase H1s in metamonads are highly divergent (**Supplementary Fig. 10**) and, in the case of *S.*  
166 *salmonicida* RNaseH1 proteins, these formed very long branches in all of our preliminary trees, that  
167 had to be removed for the final phylogenetic reconstruction. Remarkably, the phylogenetic  
168 reconstruction that includes other metamonad proteins suggests that *Giardia*, *Trepomonas* sp. PC1, *T.*  
169 *foetus* and *T. vaginalis*, also acquired bacterial RNaseH1. *Trepomonas* sp. PC1 and *Giardia*  
170 sequences cluster together but the *T. foetus* and *T. vaginalis* enzymes each emerge amidst different  
171 bacterial branches, suggesting that they have been acquired independently from the *Carpediemonas*  
172 homologs. It should, however, be noted that the support values are overall low, partly due to the fact  
173 that these sequences and their relatives are highly divergent from each other, from *Carpediemonas*  
174 bacterial-like sequences, and from typical eukaryotic RNaseH1.

175

## 176 **C. Supplementary discussion**

### 177 **C1. BUSCO incompleteness**

178 Both eukaryote-wide and protist BUSCO analyses using the BUSCO methods underperformed in our  
179 analyses. Despite using a phylogeny-guided search with the Eukaryota database, a more  
180 comprehensive database than the protist BUSCO database, a remarkably large number of BUSCO  
181 proteins were inconsistently present in Metamonada. This is not surprising, as the clade harbors a very

182 diverse group of taxa with varied lifestyles and many have undergone genome streamlining<sup>9,10,23-25</sup>,  
183 and the BUSCO databases are expected to be more accurate with greater taxonomic proximity to the  
184 studied genome<sup>1,26,27</sup>. While it might be tempting to suggest the 101 BUSCO proteins that are shared  
185 by all metamonads be used to evaluate genome completion in the clade, the overwhelming evidence of  
186 differential genome streamlining strongly indicates that databases should be lineage specific (*e.g.*,  
187 *Carpodomonas*, *Giardia*, etc). Hence, our results highlight the need for constructing such databases  
188 including proteins that showcase the sequence diversity of the groups and genes that are truly single  
189 copy in each of these lineages. Regardless, using only standard BUSCO methods to capture genome  
190 completion will still fall short in such assessments as it will fail to evaluate the most difficult-to-  
191 assemble regions of the genome<sup>27,28</sup>. For that reason, combined approaches such as the ones used here  
192 provide a more comprehensive global overview of genome completeness.

193

#### 194 **D. Supplementary references**

- 195 1 Waterhouse, R. M. *et al.* BUSCO applications from quality assessments to gene prediction  
196 and phylogenomics. *Mol. Biol. Evol.* **35**, 543-548 (2018).
- 197 2 Leger, M. M. *et al.* Organelles that illuminate the origins of *Trichomonas* hydrogenosomes  
198 and *Giardia* mitosomes. *Nat Ecol Evol* **1**, 0092 (2017).
- 199 3 Mi, H. *et al.* PANTHER version 11: expanded annotation data from Gene Ontology and  
200 Reactome pathways, and data analysis tool enhancements. *Nucleic Acids Res.* **45**, D183-D189  
201 (2017).
- 202 4 Eddy, S. R. Accelerated profile HMM searches. *PLoS Comp. Biol.* **7**, e1002195 (2011).

- 203 5 Yamada, K. D., Tomii, K. & Katoh, K. Application of the MAFFT sequence alignment  
204 program to large data-reexamination of the usefulness of chained guide trees. *Bioinformatics*  
205 **32**, 3246-3251 (2016).
- 206 6 Nguyen, L. T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective  
207 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**,  
208 268-274 (2015).
- 209 7 Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and Visualization of  
210 Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635-1638 (2016).
- 211 8 Hamann, E. *et al.* Syntrophic linkage between predatory *Carpodomonas* and specific  
212 prokaryotic populations. *ISME J* **11**, 1205-1217 (2017).
- 213 9 Tanifuji, G. *et al.* The draft genome of *Kipferlia bialata* reveals reductive genome evolution  
214 in fornicate parasites. *PLoS One* **13**, e0194487 (2018).
- 215 10 Xu, F. *et al.* On the reversibility of parasitism: adaptation to a free-living lifestyle via gene  
216 acquisitions in the diplomonad *Trepomonas sp.* PC1. *BMC Biol.* **14**, 62 (2016).
- 217 11 Ying, H. *et al.* Comparative genomics reveals the distinct evolutionary trajectories of the  
218 robust and complex coral lineages. *Genome Biol.* **19**, 175-175 (2018).
- 219 12 Voolstra, C. *et al.* The ReFuGe 2020 Consortium—using “omics” approaches to explore the  
220 adaptability and resilience of coral holobionts to environmental change. *Front. Mar. Sci.* **2**  
221 (2015).
- 222 13 Benchimol, M. *et al.* Draft genome sequence of *Tririchomonas foetus* Strain K. *Genome*  
223 *Announc.* **5**, e00195-00117 (2017).
- 224 14 Tromer, E., van Hooff, J., Kops, G. & Snel, B. Mosaic origin of the eukaryotic kinetochore.  
225 *Proc. Natl. Acad. Sci.*, 201821945 (2019).

- 226 15 Dang, H. Q. & Li, Z. The Cdc45.Mcm2-7.GINS protein complex in trypanosomes regulates  
227 DNA replication and interacts with two Orc1-like proteins in the origin recognition complex.  
228 *J. Biol. Chem.* **286**, 32424-32435 (2011).
- 229 16 Onuma, R., Mishra, N. & Miyagishima, S. Y. Regulation of chloroplast and nucleomorph  
230 replication by the cell cycle in the cryptophyte *Guillardia theta*. *Sci. Rep.* **7**, 2345 (2017).
- 231 17 Suzuki, S., Ishida, K. & Hirakawa, Y. Diurnal transcriptional regulation of endosymbiotically  
232 derived genes in the Chlorarachniophyte *Bigeloviella natans*. *Genome Biol. Evol.* **8**, 2672-  
233 2682 (2016).
- 234 18 Romero, H. *et al.* Single molecule tracking reveals functions for RarA at replication forks but  
235 also independently from replication during DNA repair in *Bacillus subtilis*. *Sci. Rep.* **9**, 1997  
236 (2019).
- 237 19 Yoshimura, A., Seki, M. & Enomoto, T. The role of WRNIP1 in genome maintenance. *Cell*  
238 *Cycle* **16**, 515-521 (2017).
- 239 20 Cerritelli, S. *et al.* Failure to produce mitochondrial DNA results in embryonic lethality in  
240 RNaseH1 null mice. *Mol. Cell* **11**, 807-815 (2003).
- 241 21 Nowotny, M. *et al.* Specific recognition of RNA/DNA hybrid and enhancement of human  
242 RNase H1 activity by HBD. *EMBO J.* **27**, 1172-1181 (2008).
- 243 22 Cerritelli, S. M. & Crouch, R. J. Ribonuclease H: the enzymes in eukaryotes. *FEBS J.* **276**,  
244 1494-1505 (2009).
- 245 23 Morrison, H. G. *et al.* Genomic Minimalism in the Early Diverging Intestinal Parasite *Giardia*  
246 *lamblia*. *Science* **317**, 1921-1926 (2007).
- 247 24 Xu, F. *et al.* The compact genome of *Giardia muris* reveals important steps in the evolution of  
248 intestinal protozoan parasites. *Microb. Genom.* (2020).



- 249 25 Xu, F. *et al.* The genome of *Spironucleus salmonicida* highlights a fish pathogen adapted to  
250 fluctuating environments. *PLoS Genet.* **10**, e1004053 (2014).
- 251 26 Saary, P., Mitchell, A. L. & Finn, R. D. Estimating the quality of eukaryotic genomes  
252 recovered from metagenomic analysis with EukCC. *Genome Biol.* **21**, 244 (2020).
- 253 27 Hanschen, E., Hovde, B. & Starckenburg, S. An evaluation of methodology to determine algal  
254 genome completeness. *Algal Res.* **51**, 102019 (2020).
- 255 28 Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality,  
256 completeness, and phasing assessment for genome assemblies. *Genome Biol.* **21**, 245 (2020).
- 257 29 Waterhouse, A. M., Procter, J. B., Martin, D. M., Clamp, M. & Barton, G. J. Jalview Version  
258 2--a multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189-  
259 1191 (2009).
- 260 30 Musacchio, A. The molecular biology of spindle assembly checkpoint signaling dynamics.  
261 *Curr. Biol.* **25**, R1002-R1018 (2015).
- 262 31 Alfieri, C., Zhang, S. & Barford, D. Visualizing the complex functions and mechanisms of the  
263 anaphase promoting complex/cyclosome (APC/C). *Open Biol.* **7** (2017).
- 264 32 Tromer, E., Bade, D., Snel, B. & Kops, G. J. Phylogenomics-guided discovery of a novel  
265 conserved cassette of short linear motifs in BubR1 essential for the spindle checkpoint. *Open*  
266 *Biol.* **6** (2016).
- 267 33 Vleugel, M. *et al.* Arrayed BUB recruitment modules in the kinetochore scaffold KNL1  
268 promote accurate chromosome segregation. *J. Cell Biol.* **203**, 943-955 (2013).
- 269 34 Shepperd, L. A. *et al.* Phosphodependent recruitment of Bub1 and Bub3 to Spc7/KNL1 by  
270 Mph1 kinase maintains the spindle checkpoint. *Curr. Biol.* **22**, 891-899 (2012).

- 271 35 Tromer, E., Snel, B. & Kops, G. Widespread recurrent patterns of rapid repeat evolution in  
272 the kinetochore scaffold KNL1. *Genome Biol. Evol.* **7**, 2383-2393 (2015).
- 273 36 Moyle, M. W. *et al.* A Bub1-Mad1 interaction targets the Mad1-Mad2 complex to unattached  
274 kinetochores to initiate the spindle checkpoint. *J. Cell Biol.* **204**, 647-657 (2014).
- 275 37 Ji, Z., Gao, H., Jia, L., Li, B. & Yu, H. A sequential multi-target Mps1 phosphorylation  
276 cascade promotes spindle checkpoint signaling. *Elife* **6** (2017).
- 277 38 Zhang, G. *et al.* Bub1 positions Mad1 close to KNL1 MELT repeats to promote checkpoint  
278 signalling. *Nat. Commun.* **8**, 15822 (2017).
- 279 39 Faesen, A. C. *et al.* Basis of catalytic assembly of the mitotic checkpoint complex. *Nature*  
280 **542**, 498-502 (2017).
- 281 40 Izawa, D. & Pines, J. The mitotic checkpoint complex binds a second CDC20 to inhibit active  
282 APC/C. *Nature* **517**, 631 (2014).
- 283 41 Di Fiore, B., Wurzenberger, C., Davey, N. E. & Pines, J. The mitotic checkpoint complex  
284 requires an evolutionary conserved cassette to bind and inhibit active APC/C. *Mol. Cell* **64**,  
285 1144-1153 (2016).
- 286 42 Burton, J. L. & Solomon, M. J. D box and KEN box motifs in budding yeast Hsl1p are  
287 required for APC-mediated degradation and direct binding to Cdc20p and Cdh1p. *Genes Dev.*  
288 **15**, 2381-2395 (2001).

## 289 **E. Supplementary figures**

290 **Supplementary Fig 1** Orc1-6 and Cdc6 proteins. **A)** Left: typical domain architecture observed for  
291 Orc1-6 and Cdc6 in *Saccharomyces cerevisiae*, Right: representative domain architecture of  
292 metamonad proteins drawn to reflect the most common protein size. If no species name is given, then  
293 the depicted domain structure was found in all of the metamonads where present. Numbers on the right

294 of each depiction correspond to the total protein length or its range in the case of metamonads  
295 (additional information in **Supplementary Table 1**). B) Comparison of Orc1, Cdc6 and Orc1/Cdc6-  
296 like protein lengths across 81 eukaryotes encompassing metamonads and non-metamonads protists  
297 (source information in **Supplementary Table 1**). Metamonad proteins are highlighted with green  
298 shaded bubbles in the background. C) Orc1/Cdc6 partial ATPase domain showing Walker A and  
299 Walker B motifs including R-finger. Reference species at the top. Multiple sequence alignment was  
300 visualized with Jalview<sup>29</sup> using the Clustal colouring scheme. D) Phylogenetic reconstruction of Orc1,  
301 Cdc6 and Orc1/Cdc6-like proteins inferred with IQ-TREE<sup>6</sup> under the LG+ C10+F+  $\Gamma$  model using  
302 1000 ultrafast bootstraps (bootstrap value ranges for branches are shown with black and grey dots).  
303 The alignment consists of 81 taxa with 367 sites after trimming. Orc1/Cdc6-like proteins do not form a  
304 clade with *bona fide* Orc1 and Cdc6 proteins making it impossible to definitively establish whether or  
305 not they are orthologs.

306 **Supplementary Fig 2** The distribution of core molecular systems of the replisome, double strand  
307 break repair and endonucleases in nucleomorph genomes of cryptophyte and chlorarachniophytes.

308 **Supplementary Fig 3** The distribution of core molecular systems of DNA repair across eukaryotic  
309 diversity. A schematic global eukaryote phylogeny is shown on the left with classification of the major  
310 metamonad lineages indicated. Double strand break repair and endonuclease sets. \*\*\* *Carpediemonas*-  
311 Like Organisms. '?' is used in cases where correct orthology was difficult to establish, so the protein  
312 name appears with the suffix '-like' in tables.

313 **Supplementary Fig 4. Presence/absence diagram of LECA kinetochore components in**  
314 **eukaryotes, with a greater sampling of metamonads, including *C. membranifera* and *C. frisia*.**

315 Left: matrix of presences (coloured) and absences (light grey) of kinetochore, SAC and APC/C  
316 proteins that were present in LECA. On top: names of the different subunits; single letters (A-X)

317 indicate Centromere protein A-X (*e.g.*, CenpA) and numbers, APC/C subunit 1-15 (*e.g.*, Apc1). E2S  
318 and E2C, refer to E2 ubiquitin conjugases S and C, respectively. Colour schemes correspond to the  
319 kinetochore overview figure on the right and to that used in Figure 1. Right: cartoon of the components  
320 of the kinetochore, SAC signalling, the APC/C and its substrates (Cyclin A/B) in LECA and  
321 *Carpediemonas* species to indicate the loss of components (light grey shading). Blue lines indicate the  
322 presence of proteins that are part of the MCC. Asterisk: Apc10 has three paralogs in *C. membranifera*  
323 and two in *C. frisia*. One is the canonical Apc10, the two others are fused to a BTB-Kelch protein of  
324 which its closest homologs is a likely adapter for the E3 ubiquitin ligase Cullin 3.

325 **Supplementary Fig 5. *Carpediemonas* harbours three different types of Histone H3 proteins, a**  
326 **centromere-specific variant (CenpA).** Multiple sequence alignment of different Histone H3 variants  
327 in eukaryotes and metamonads, including the secondary structure of canonical H3 in humans (pdb:  
328 6ESF\_A). CenpA orthologs are characterized by extended amino and carboxy termini and a large L1  
329 loop. Red names in the CenpA panel indicate for which species centromere/kinetochore localization  
330 has been confirmed. In addition to CenpA and canonical Histone H3-variants, multiple eukaryotes,  
331 including *C. membranifera* and *C. frisia*, harbour other divergent H3 variants. Such divergent variants  
332 make the annotation of Histone H3 homologs ambiguous (see Asterisks; incomplete sequences).  
333 Multiple sequence alignments were visualized with Jalview<sup>29</sup>, using the Clustal colour scheme.  
334 Asterisks indicate two potential CenpA candidates in *T. vaginalis*.

335 **Supplementary Fig 6. Likely presence of SAC signalling in *Carpediemonas*.** A) Short linear motifs  
336 form the basis of SAC signalling. During prometaphase, unattached kinetochores catalyse the  
337 production of inhibitor of the cell cycle machinery, a phenomenon known as the SAC<sup>30</sup>. (I) The main  
338 protein scaffold of SAC signalling is the kinase MadBub (paralogs Mad3/Bub1 exist in eukaryotes),  
339 which consist of many short linear motifs (SLiMs) that mediate the interaction of SAC components

340 and the APC/C (light blue)<sup>31,32</sup>. MadBub itself is recruited to the kinetochore through interaction with  
341 Bub3 (GLEBS), which on its turn binds repeated phosphomotifs in Knl1<sup>33-35</sup>. The CDI or CMI motif  
342 aids to recruit Mad1<sup>36-38</sup>, which has a Mad2-interaction Motif (MIM) that mediated the kinetochore-  
343 dependent conversion of open-Mad2 to Mad2 in a closed conformation<sup>39</sup>. **(II)** Mad2, MadBub, Bub3  
344 and 2x Cdc20 (APC/C co-activator) form the mitotic checkpoint complex (MCC) and block the  
345 APC/C<sup>32,40,41</sup>. MadBub contains 3 different APC/C degrons (D-box, KEN-box and ABBA motif)<sup>31</sup> that  
346 direct its interaction with 2x Cdc20s and effectively make the MCC a pseudo substrate of the APC/C.  
347 **(III)** Increasing amounts of kinetochore-microtubule attachments silence the production of the MCC at  
348 kinetochores and the APC/C is released. Cdc20 now presents its substrates Cyclin A and Cyclin B  
349 (some eukaryotes have other substrates as well, but they are not universally conserved) for  
350 ubiquitination and subsequent degradation through recognition of a Dbox motif<sup>42</sup>. Chromosome  
351 segregation will now be initiated (anaphase). **B)** Presence/absence matrix of motifs involved in SAC  
352 signalling in a selection of Eukaryotes and Metamonads, including *C. membranifera* and *C. frisia*.  
353 Colours correspond to the motifs in panel A, light grey indicates motif loss. *N* signifies the number of  
354 MadBub homologs that are present in each species. ‘Incomplete’ points to sequences that were found  
355 to be incomplete due to gaps in the genome assembly. Question marks indicate the uncertainty in the  
356 presence of that particular motif. Although Metamonads have all four MCC components (Mad2, Bub3,  
357 MadBub and Cdc20), most homologs do not contain the motifs to elicit a canonical SAC signalling  
358 and it is therefore likely that they do not have a SAC response. Exceptions are *C membranifera*, *C.*  
359 *frisia* and *Kipferlia bialata*. They retained the N-terminal KEN-boxes and one ABBA motif, which are  
360 involved in the binding of two Cdc20s and a Mad2-interaction motif (MIM) in Mad1 and Cdc20. **C)**  
361 Multiple sequence alignments of the motifs from panel A and B. Coloured motif boxes correspond to

362 panel A and B. Multiple sequence alignments were visualized with Jalview<sup>29</sup>, using the Clustal  
363 colouring scheme. Asterisks indicate ambiguous motifs in *Carpediemonas membranifera*.

364 **Supplementary Fig 7 Histogram showing the frequency distribution of single nucleotide variants**  
365 **in the genome of *C. membranifera*.** Diagram showing the typical distribution of a haploid genome.

366 **Supplementary Fig 8 Maximum likelihood reconstruction of Endo IV.** The unrooted tree contains  
367 eukaryotic and prokaryotic Endo IV sequences, showing *Carpediemonas* sequences emerging within  
368 bacterial proteins. The tree was inferred with IQ-TREE under the LG+I+C20 model with 1000  
369 ultrafast bootstraps; alignment length was 276. Scale bar shows the inferred number of amino acid  
370 substitutions per site.

371 **Supplementary Fig 9 Maximum likelihood reconstruction of RarA.** The unrooted tree contains  
372 eukaryotic and prokaryotic sequences, showing *Carpediemonas* sequences emerging within bacterial  
373 proteins. The tree was inferred with IQ-TREE under the LG+I+C20 model with 1000 ultrafast  
374 bootstraps; alignment length was 414. Scale bar shows the inferred number of amino acid substitutions  
375 per site.

376 **Supplementary Fig 10 Maximum likelihood reconstruction of RNase H1.** *Carpediemonas* RarA-  
377 like proteins emerge within bacterial proteins. Parabasalia and diplomonada proteins highlighting the  
378 proteins have been acquired in different events. The tree was inferred with IQ-TREE under the  
379 LG+I+G+C20 model with 1000 ultrafast bootstraps; alignment length was 149. Scale bar shows the  
380 inferred number of amino acid substitutions per site.

## 381 **F. Supplementary tables**

382 Secure download link: <http://perun.biochem.dal.ca/downloads/dsalas/SuppInfo.tar.gz>

383 **Supplementary Table 1:**

384 **Supplementary Table 1A** BUSCO proteins found in Metamonada based on searches for 245 proteins

385 present in at least one taxon

386 **Supplementary Table 1B** DNA replication and repair orthologs in 18 diverse eukaryotic genomes

387 **Supplementary Table 1C** Spindle assembly, kinetochore and APC/C orthologs in 18 diverse

388 eukaryotic genomes

389 **Supplementary Table 1D** Additional genomes queried during the searches for ORC, Cdc6 and Ndc80

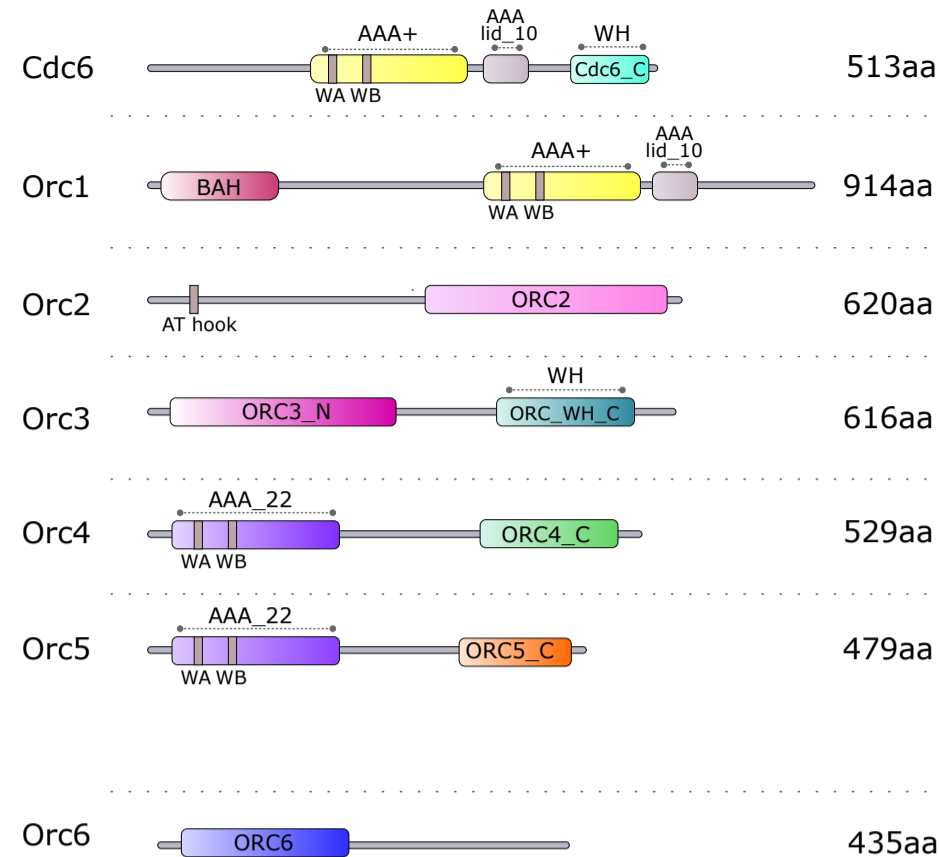
390 proteins

391 **Supplementary Table 1E** Lengths of Orc1-6, Cdc6 and Orc1/Cdc6-like proteins and domain

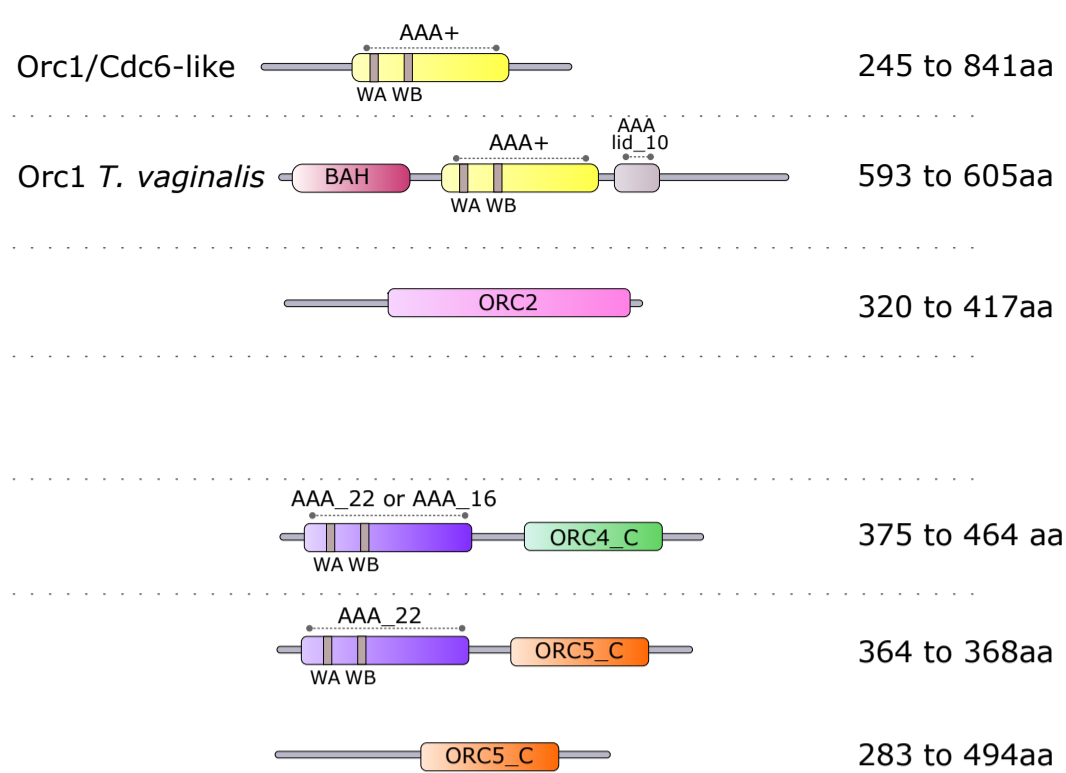
392 architecture comparisons between metamonads and other eukaryotes.

393 **Supplementary Table 1F** Orc1, Cdc6 and Orc1/Cdc6-like proteins. Information used in

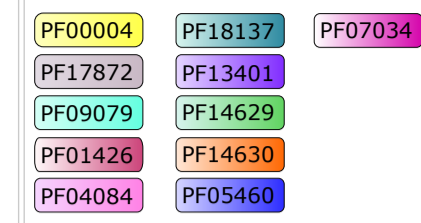
394 Supplementary Figure 1 panels B and D

**A****Supplementary Figure 1**Fungi (*Saccharomyces cerevisiae*)

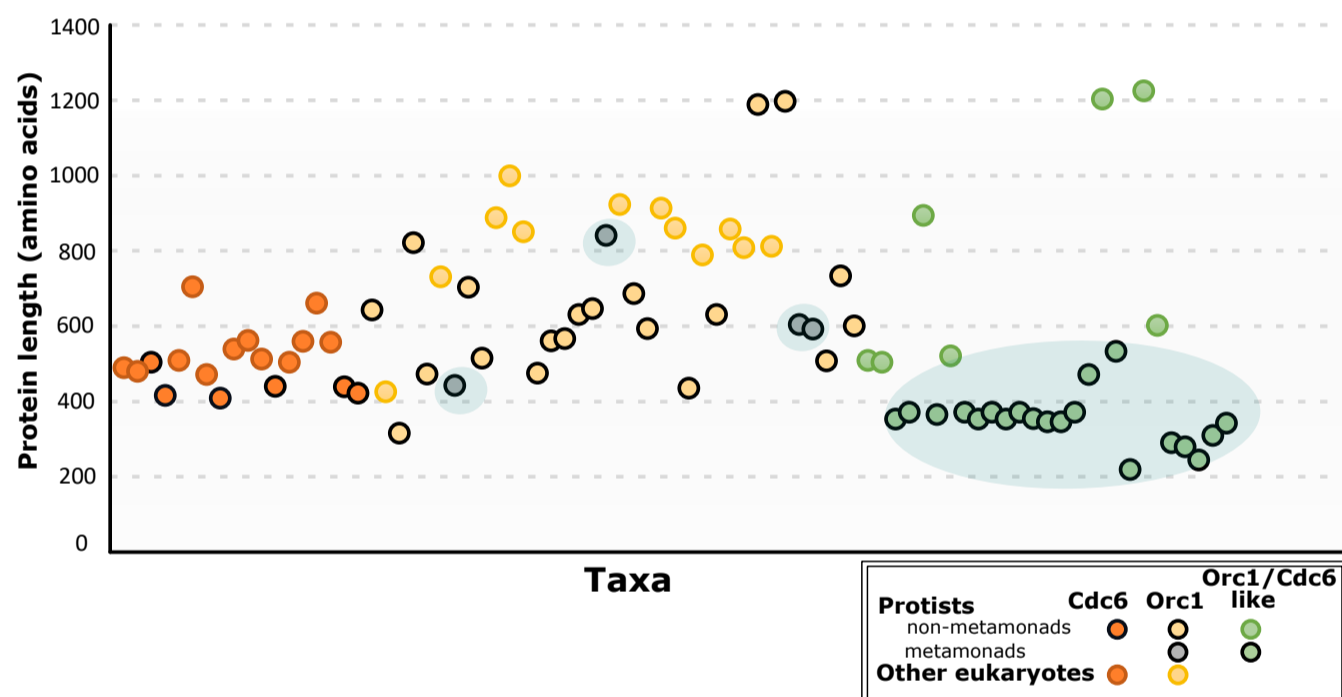
Metamonada



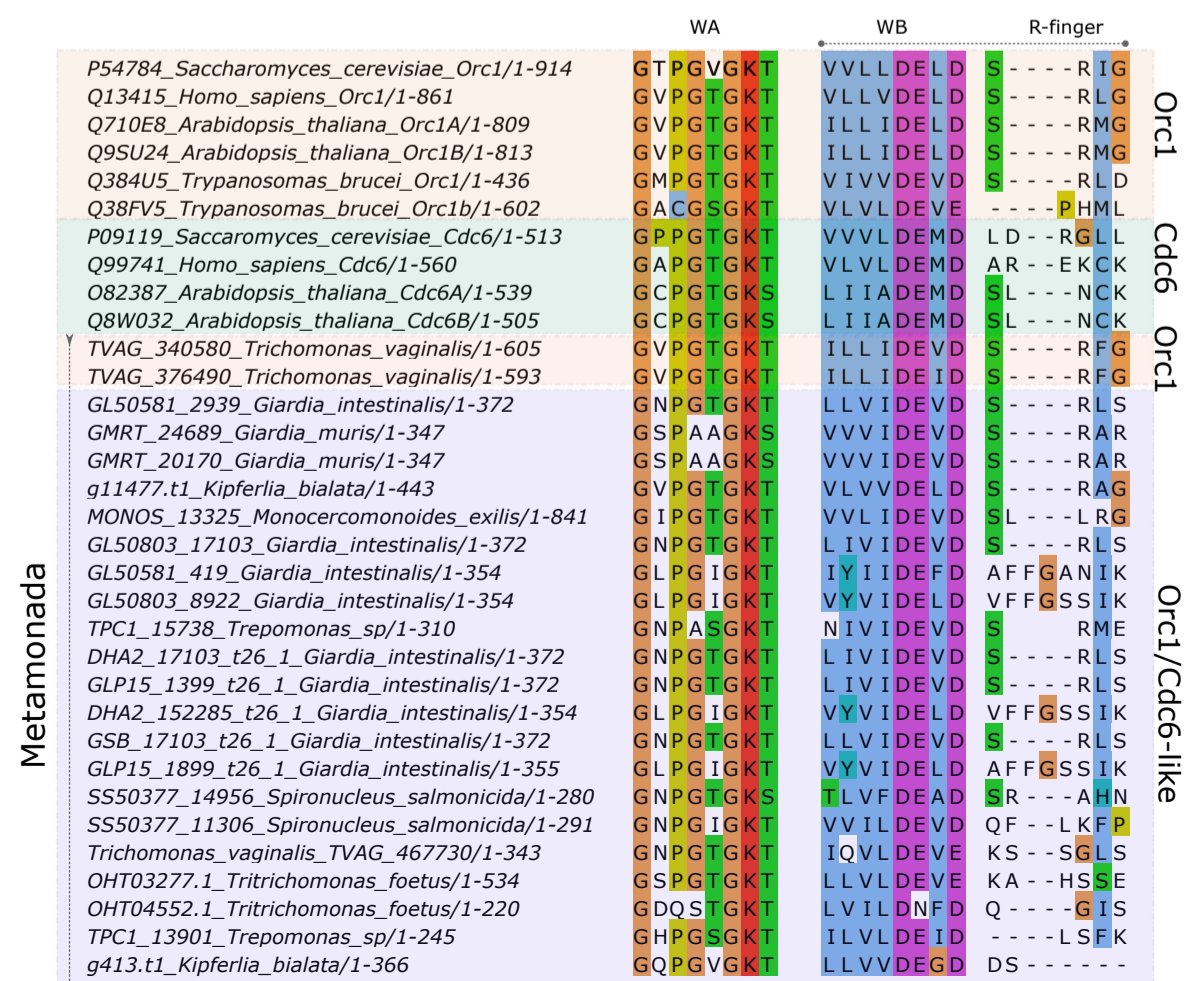
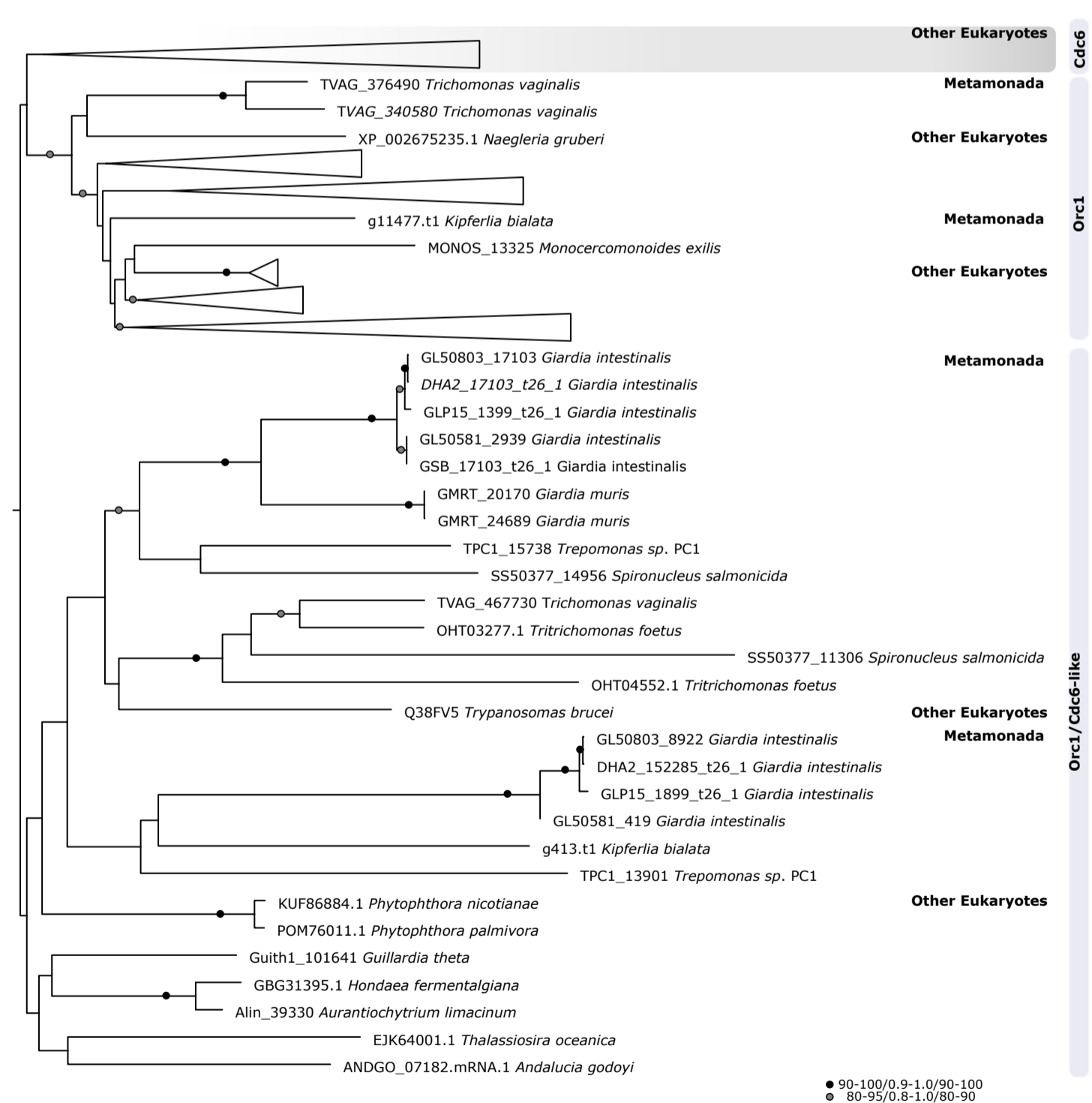
PFAM Accessions



bioRxiv preprint doi: <https://doi.org/10.1101/2021.03.14.435266>; this version posted March 15, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

**B****C**

ATPase domain (partial) Orc1, Cdc6 and Orc1/Cdc6-like

**D**

● 90-100/0.9-1.0/90-100  
● 80-95/0.8-1.0/80-90

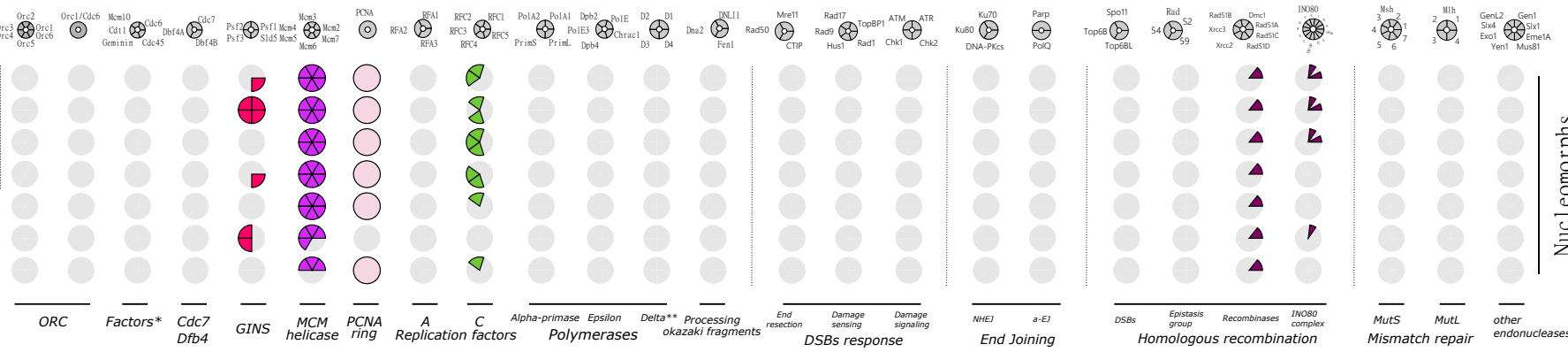


# Supplementary Figure 2

## DNA Replication & DNA Repair

absent: components:

<b>Cryptophyte</b>	<i>Cryptomonas paramecium</i>
	<i>Hemiselmis andersenii</i>
	<i>Chroomonas mesostigmatica</i>
	<i>Guillardia theta</i>
<b>Chlorarachniophyte</b>	<i>Bigelowiella natans</i>
	<i>Lotharella vacuolate</i>
	<i>Amorphochlora amoebiformis</i>



Replisome

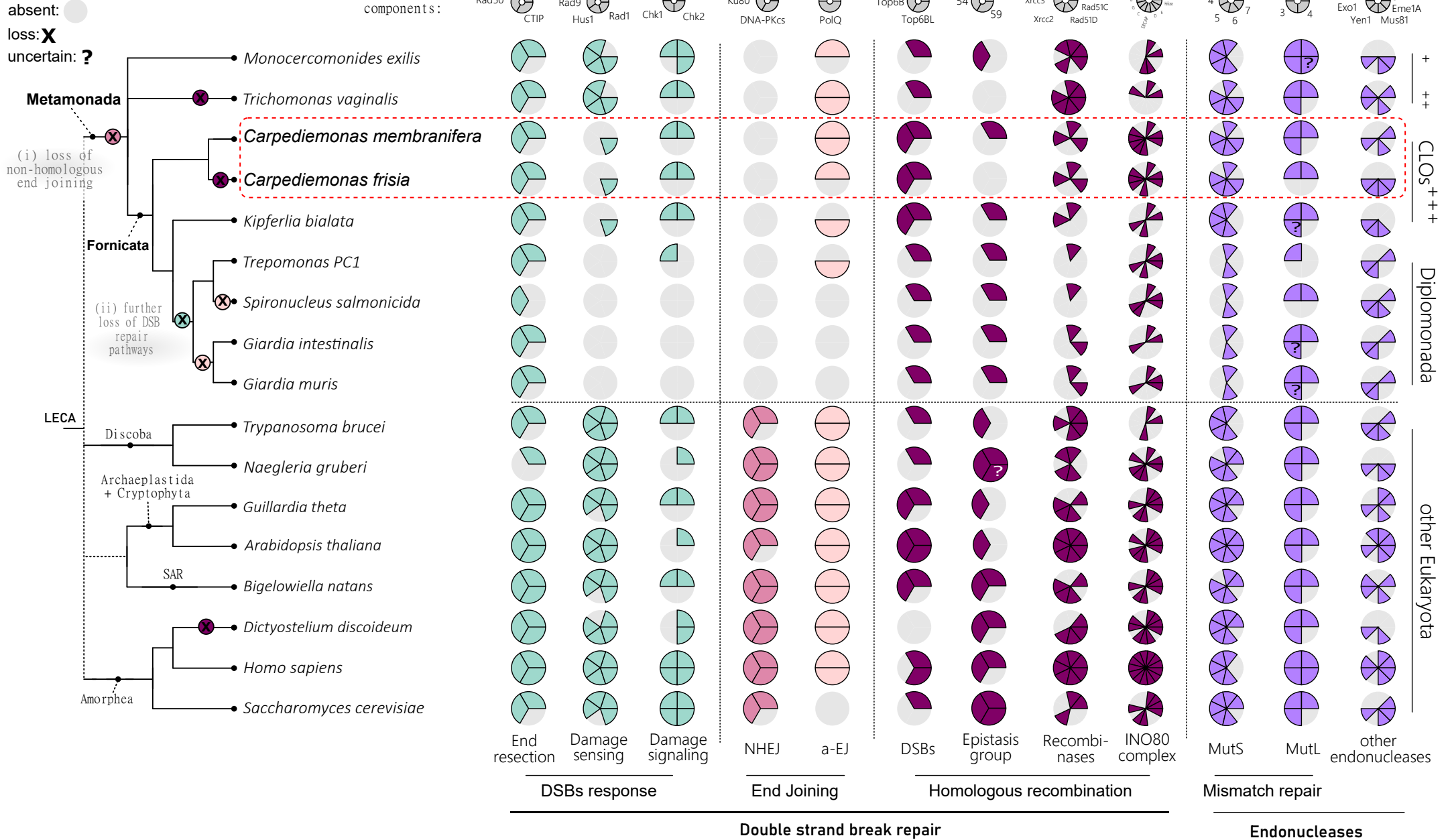
Double strand break repair

Endonucleases

Nucleomorphs

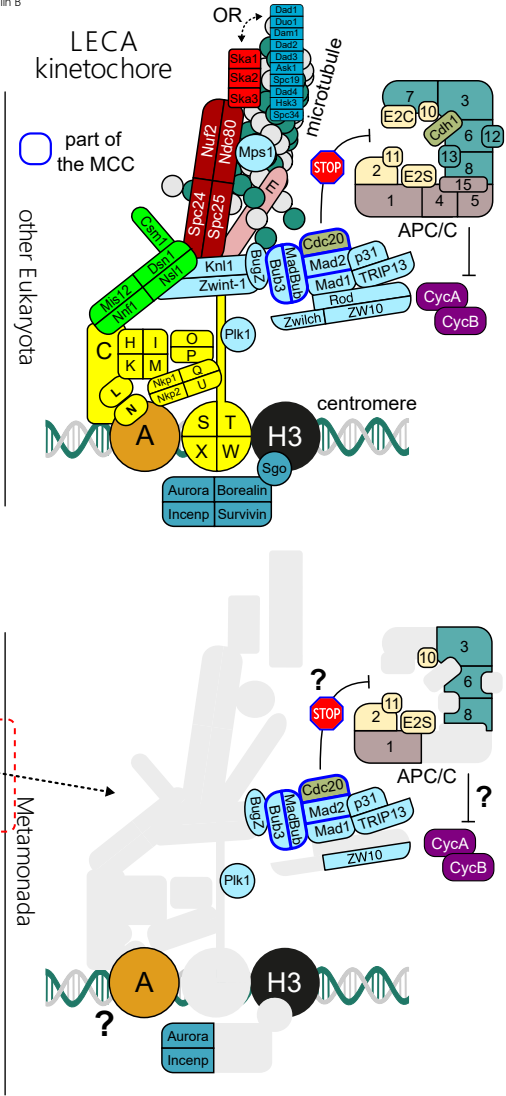
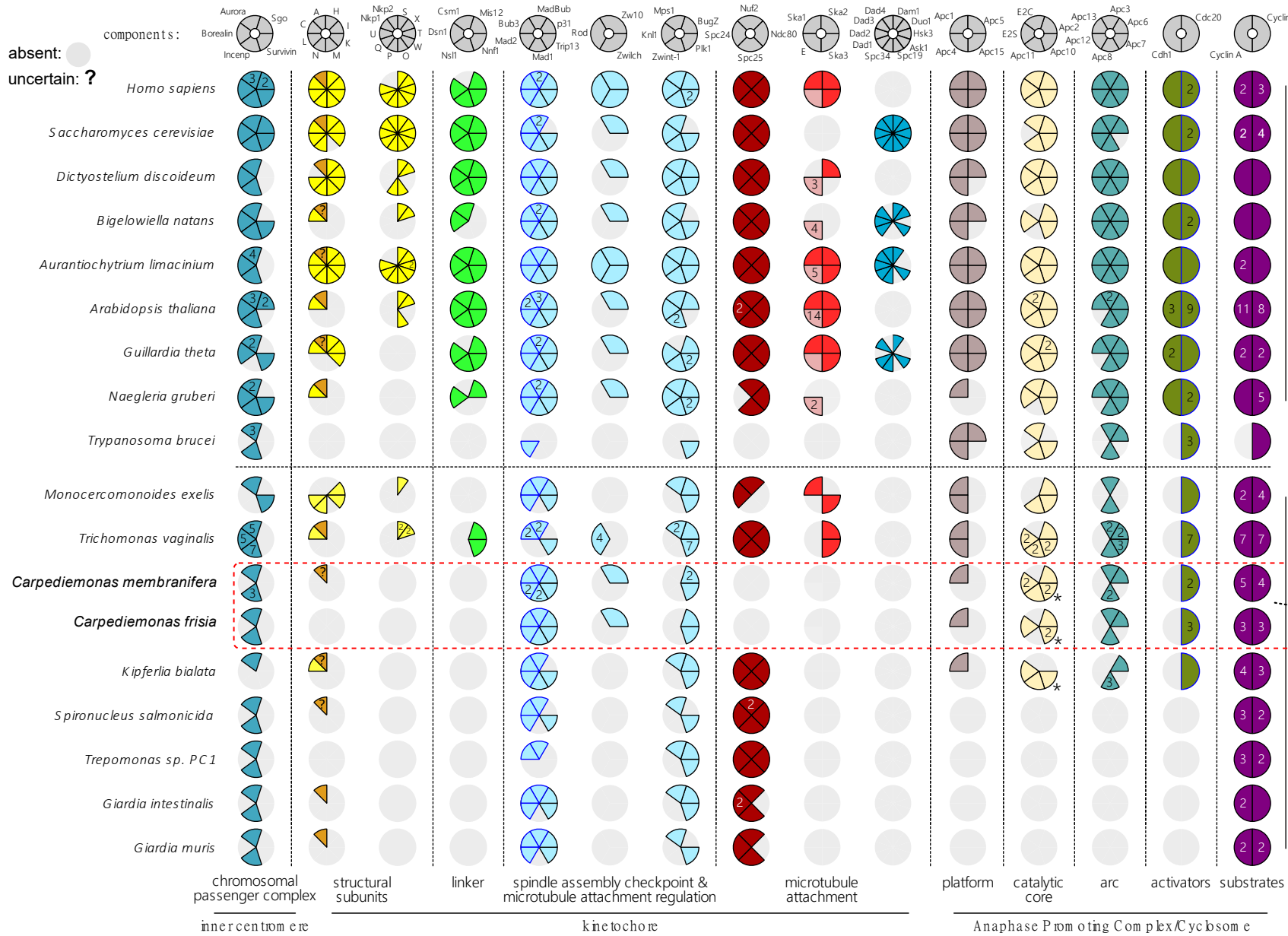
# Supplementary Figure 3

## DNA repair & Endonucleases



# Supplementary Figure 4

## Kinetochores Network and Spindle Assembly Checkpoint Signaling

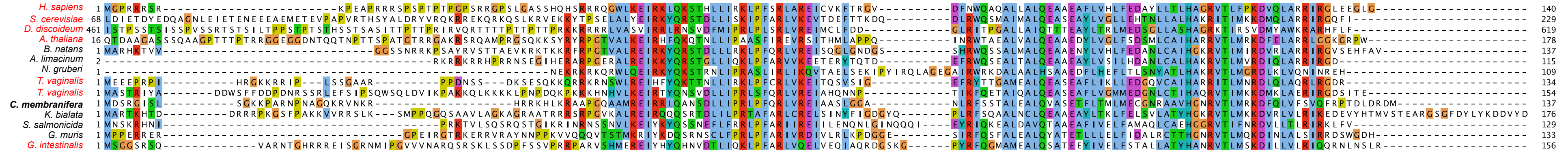


# Supplementary Figure 5

## canonical H3



## centromeric H3 (CenpA)

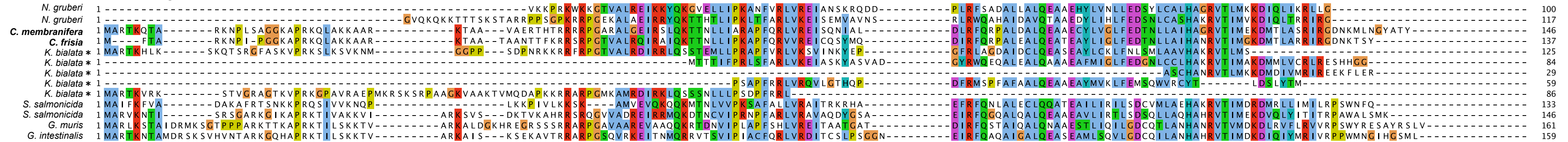


←----- divergent N-terminus ----->

←----- Loop 1 ----->

←-- extended C-terminus -->

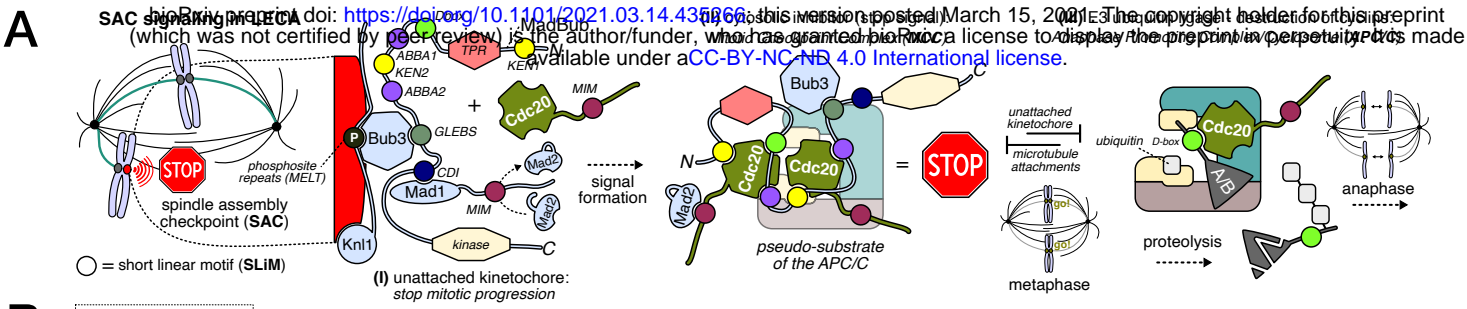
## other divergent H3-variants



Consensus: MARTKQTA I+R+KSSR+RKST+GGKAPRKQLASKAARS I+SG+R+RARKTASTPSTGGVKKPRRYRPGTVALREIRKYQKSTELLIRKLPFQRLVREIAQDFKT+DGYQ+Q+AGEGDLRFQSAALQEAEEAYLVGLFEDTNLCAIHAQRVTIMKDKMQLARRIRGERWP+Y+K+++L+GSGFDYLYKDDVYD



# Supplementary Figure 6



**B**

X = protein lost  
 ● = motif/domain lost

	MadBub						Knl1			Cdc20	Mad1	CycA	CycB	
Conserved features:	KEN1	TPR	Dbox	ABBA1	KEN2	ABBA2	GLEBS	CDI	kinase	MELT	MIM	D-box/KEN	presence of the SAC	
Interacting with:	Cdc20						Bub3	Mad1	-	Bub3	Mad2	Cdc20		
<i>Homo sapiens</i>	2	●	●	●	●	●	●	●	●	●	●	●	STOP	
<i>Saccharomyces cerevisiae</i>	2	●	●	●	●	●	●	●	●	●	●	●	STOP	
<i>Dictyostelium discoideum</i>	1	●	●	●	●	●	●	●	●	●	●	●	STOP	
<i>Bigeloviella natans</i>	2	●	●	incomplete sequence	●	●	●	●	●	●	●	●	STOP	
<i>Aurantiochytrium limacinium</i>	1	●	●	●	●	●	●	●	●	●	●	●	STOP	
<i>Arabidopsis thaliana</i>	3	●	●	●	●	●	●	●	●	●	●	●	STOP	
<i>Naegleria gruberi</i>	2	●	●	●	●	●	●	●	●	●	●	●	STOP	
<i>Trypanosoma brucei</i>	0	x	x	x	x	x	x	x	x	x	x	x		
<i>Monocercomonoides sp. PA203</i>	1	●	●	●	●	●	●	●	x	x	x	●	●	
<i>Trichomonas vaginalis</i>	2	●	●	●	●	●	●	●	x	x	x	●	●	
<i>Carpediemonas membranifera</i>	1	●	●	●	●	●	●	●	x	?	?	●	?	STOP
<i>Carpediemonas frisia</i>	1	●	●	●	●	●	●	●	x	●	●	●	●	STOP
<i>Kipferlia bialata</i>	1	●	●	●	●	●	incomplete sequence	●	x	incomplete sequence	●	●	●	STOP
<i>Spironucleus salmonicida</i>	1	incomplete sequence	●	●	●	●	incomplete sequence	●	x	x	x	●	●	
<i>Trepomonas sp. PC1</i>	1	incomplete sequence	●	●	●	●	incomplete sequence	●	x	x	x	●	●	
<i>Giardia intestinalis</i>	1	●	●	●	●	●	●	●	x	x	x	●	●	
<i>Giardia muris</i>	1	●	●	●	●	●	●	●	x	x	x	●	●	

other Eukaryota  
 Metamonada

**C**

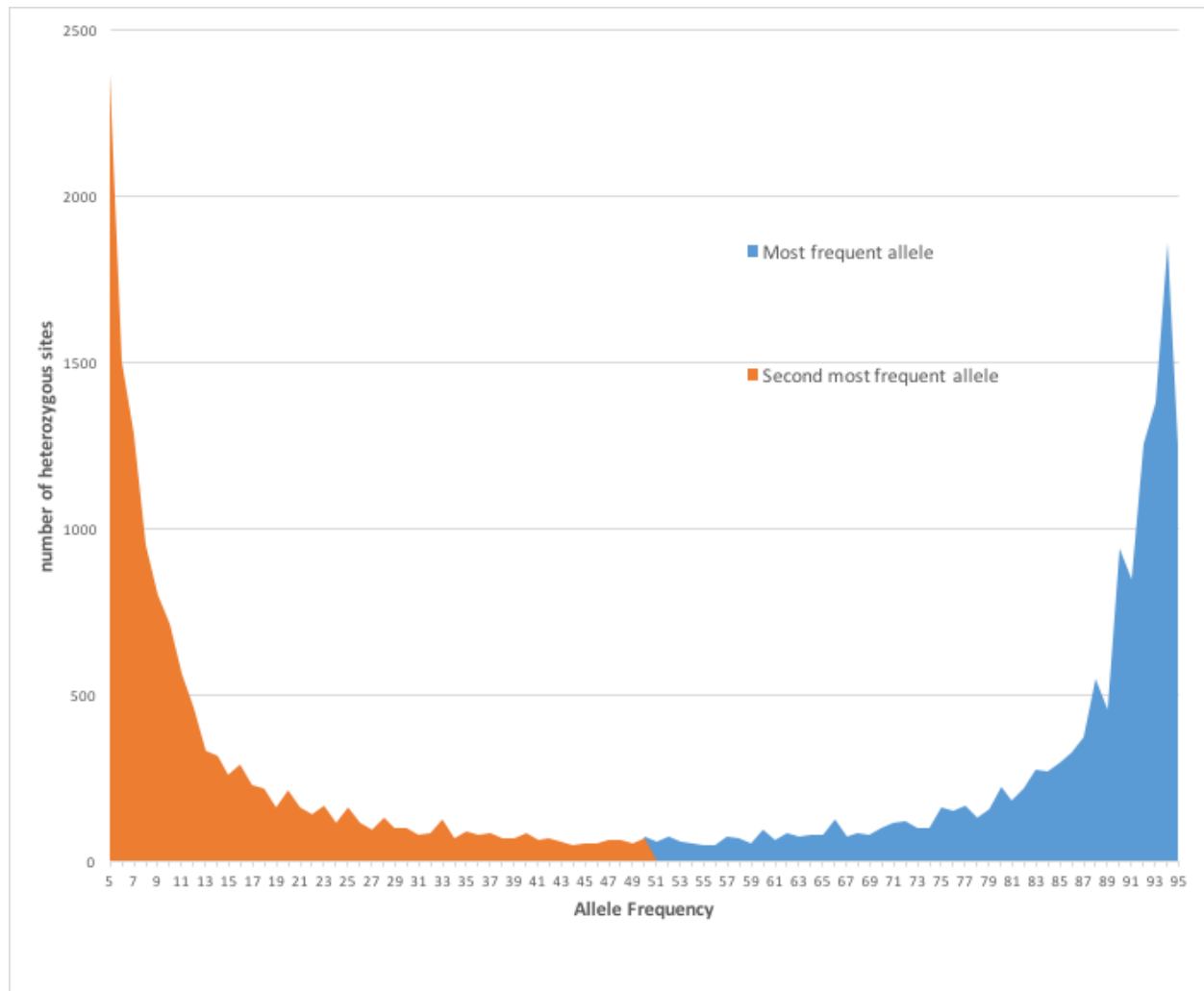
	MadBub	KEN1	Dbox-like	ABBA1	KEN2	ABBA2	GLEBS	CDI
<i>H. sapiens</i>	EWELSKENVQPLRQGRIR	QRSTLAELEKSK	RITVFDE	PRAKENELQ	SFTPYVEE	GFESFEEIRAEVF	PSPTVHTKEALGF	IMNMF
<i>S. cerevisiae</i>	EIETQKENILPLKKEGRS	KNNVFVD	ERNKENLR	KLPFRDS	EEFNTEEILAMIK	PTVTAFSKDA	INEVFSMF	
<i>D. discoideum</i>	DWENTKENIVPLKTRGRD	TRSALEGTISS	GQIFDD	NKHKENTQL	GFIYCDQ	HEISFEEYRASLF	PTMTIHTKQAF	HDMAMF
<i>A. thaliana</i>	EWELFKENVRPLKGRGRN	PSRSFGLLSR	PFAIYAD	ERNKENNSL	TFEYFVDE		VDPPTINMKEAMNT	INNMF
<i>B. natans</i>	EIEGSKENIQPLKRRGRN					RETSEFEYRAMIY	PSPTIHTKEAM	SDVLAMF
<i>A. limacinium</i>	DWEQCKENAMPKRRGRD	SRRSLQPIRSR	RFEVFEQ	ELQKENSID	AFDVYVEE		EDMTINTRLAL	DDMFEMF
<i>N. gruberi</i>	EFEAAKENIIPLKGGRRQ	QRLPFNSLAKK	NEVIYDE	EQEKENEK	HFSIFVDE	AEFTFEELRAMHY		
<i>Monocercomonoides sp.</i>						ARRSVEERERMTN		no motifs
<i>T. vaginalis</i>						TSRSFVEARLADM		<i>S. salmonicida</i>
<i>C. membranifera</i>	AWEEKKENVLPRAAQRD	VRAALMDLDAQ		DTVKENFGR	DFEVFI DD			<i>Trepomonas sp.</i>
<i>C. frisia</i>	incomplete sequence			DVTKENIGR	DFDILVDD			
<i>K. bialata</i>	EAAARDKENLQPLDGGRRQ			SFDVHDD	KAKRENSKA			
<i>G. intestinalis</i>						SDVSFEEQRLLDI		
<i>G. muris</i>						SFISPEEQHLLDL		

	MIM (Mad1)	MIM (Cdc20)	Dbox/KEN (Cyclin A)	Dbox (Cyclin B)	phosphosite repeats (MELT)
<i>H. sapiens</i>	RTKVLHMSLNPT	EAKILRLSGKPKQ	PQRTVLGLLTAN	RPRTALGDI GNK	
<i>S. cerevisiae</i>	KIRILQLRDGP	NKRILQYMPPEP	VQRLALNNVTNT	VPRTILGNVTNN	
<i>D. discoideum</i>	KTKEVLHFSNPS	ESKILSFKSKAP	HKRVALYDVT HQ	SHRGALSDLTNN	
<i>A. thaliana</i>	NTRVLRMVT	HTRILAFRNKPKQ	KKRVLVGLGELPNL	QNRKVLGDI GNL	
<i>B. natans</i>	TTKVLHMRFN PQ	ESRVLAFKFKAP		TSRRALGDISNA	
<i>A. limacinium</i>	KTKEVLHLLTFNPE	SSKILAYRQKAP			
<i>N. gruberi</i>	KFKILHMKLNPE	DAKILALTEKAP			
<i>Monocercomonoides sp.</i>			QSRREILRDVINS	TQRTCLGDISNQ	
<i>T. vaginalis</i>			RPRAASLTVAGT	RRPKALEQITNT	
<i>C. membranifera</i>	SYIVMHMSLNPT*	TSSIVHAKSAAP*	HQRVALSRVSNIT		
<i>C. frisia</i>	EKVLHATHNP	EQKVIHSRSEAP	DKKENIRKLOEA	EHKAPLGDISV*	
<i>K. bialata</i>	incomplete sequence	APKVLALSGKAP		LTRAPLRDVGRT	

*H. sapiens* (21) *D. discoideum* (5)  
*S. cerevisiae* (6) *B. natans* (8) *N. gruberi* (1)  
*A. thaliana* (5)  
*A. limacinium* (11)

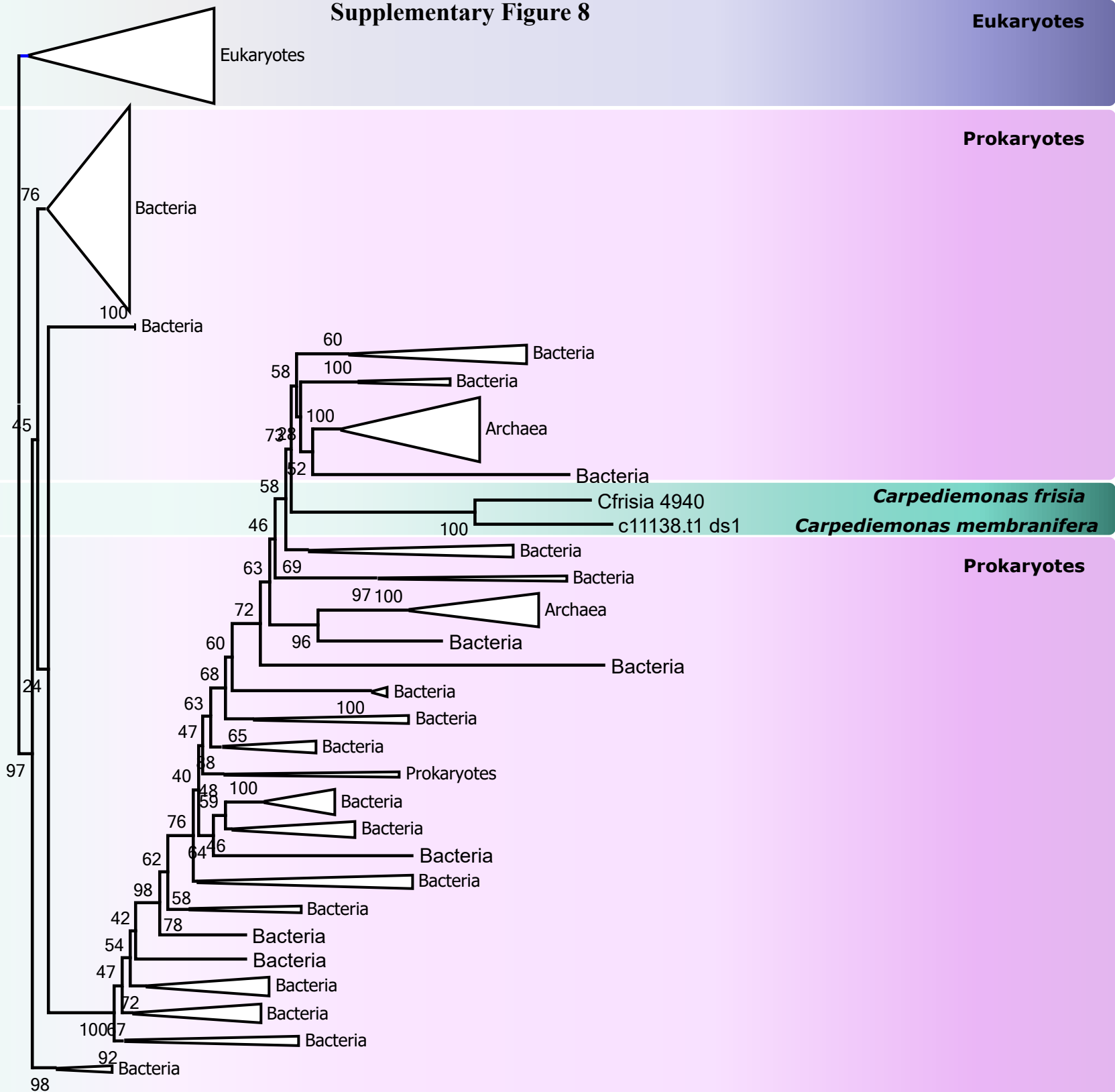
## Supplementary Figure 7



# Supplementary Figure 8

Eukaryotes

Prokaryotes



*Carpediemonas frisia*

*Carpediemonas membranifera*

Prokaryotes

0.5

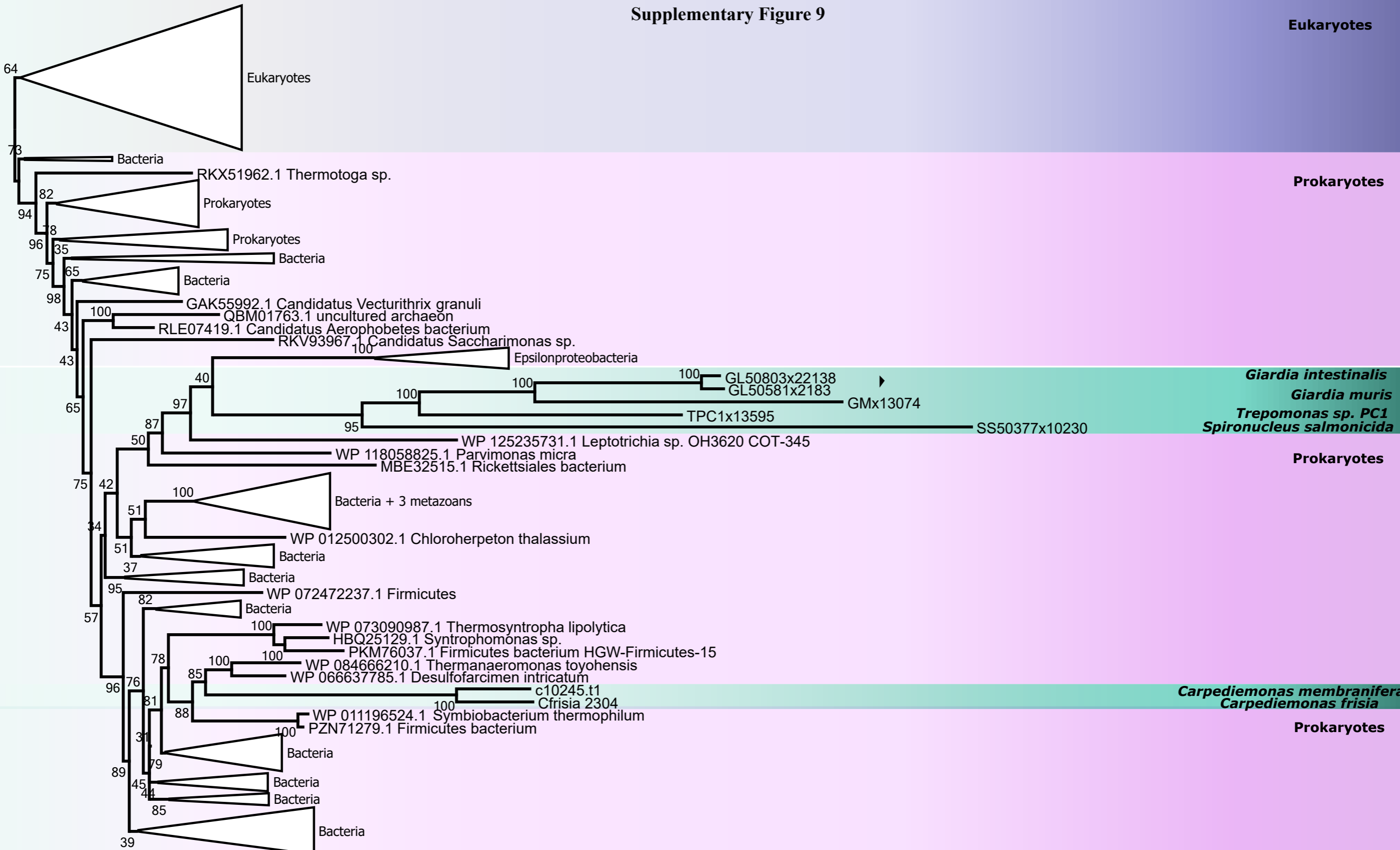
Supplementary Figure 9

Eukaryotes

Prokaryotes

Prokaryotes

Prokaryotes



0.5



