

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

2 Dayana E. Salas-Leiva^{1*}, Eelco C. Tromer^{2,3}, Bruce A. Curtis¹, Jon Jerlström-Hultqvist¹, Martin

3 Kolisko⁴, Zhenzhen Yi⁵, Joan S. Salas-Leiva⁶, Lucie Gallot-Lavallée¹, Geert J. P. L. Kops³, John M.

4 Archibald¹, Alastair G. B. Simpson⁷ and Andrew J. Roger^{1*}

5 ¹Centre for Comparative Genomics and Evolutionary Bioinformatics (CGEB), Department of

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

7 ²Department of Biochemistry, University of Cambridge, Cambridge, United Kingdom

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

9 and University Medical Centre Utrecht, Utrecht, The Netherlands

10 ⁴Institute of Parasitology Biology Centre, Czech Acad. Sci, České Budějovice, Czech Republic

11 ⁵Guangzhou Key Laboratory of Subtropical Biodiversity and Biomonitoring, School of Life Science,

12 South China Normal University, Guangzhou 510631, China

13 ⁶CONACyT-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 ⁷Centre for Comparative Genomics and Evolutionary Bioinformatics (CGEB), Department of

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

17 *corresponding authors: Andrew.Roger@dal.ca and Dayana.Salas@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: 0000-0001-6763-3388
- 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, 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 novel or alternative set of mechanisms to carry out
47 these 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 organisms¹. 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 mutations². 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) complexes³⁻⁶. 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 it⁷, while recombination-dependent synthesis uses non-homologous
60 or homologous templates for repair (reviewed in refs.^{8,9}). Recombination-dependent synthesis occurs
61 in response to single- or double-strand DNA breakage^{8,10,11}. 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 chromatin¹². 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)¹². Failure to pass any of these checkpoints (*e.g.*,
70 G1/S, S, G2/M and SAC checkpoints reviewed in refs.¹²⁻¹⁴) 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-level¹⁵⁻¹⁷. Indeed, recent studies show
79 that metamonad parasites have secondarily lost parts of the ancestral DNA replication and
80 segregation apparatuses^{18,19}. Furthermore, metamonad proteins are often highly divergent compared
81 to other eukaryotic orthologs, indicating a high substitution rate in these organisms that is suggestive
82 of error-prone replication and/or DNA repair^{20,21}. Yet, it is unclear whether the divergent nature of
83 proteins studied in metamonads is the result from the host-associated lifestyle or is a more ancient
84 feature of Metamonada. To increase the representation of free-living metamonads in our analyses, we
85 have generated a high-quality draft genome assembly of *Carpediemonas membranifera*, a flagellate
86 isolated from hypoxic marine sediments²². Our analyses of genomes from across the tree of
87 eukaryotes show that many systems for DNA replication, repair, segregation, and cell cycle control
88 are ancestral to eukaryotes and highly conserved. However, metamonads have secondarily lost an
89 extraordinarily large number of components. Most remarkably, the free-living *Carpediemonas*
90 species have been drastically reduced further, having lost a large set of key proteins from the
91 replisome and cell-cycle checkpoints (*i.e.*, including several from the kinetochore and repair
92 pathways). We propose a hypothesis of how DNA replication may be achieved in this organism.

93 **Results**

94 **The *C. membranifera* genome assembly is complete.**

95 Our assembly for *C. membranifera* (a member of the Fornicata clade within metamonads, **Fig. 1**)
96 is highly contiguous (**Table 1**) and has deep read coverage (*i.e.*, median coverage of 150× with short
97 reads and 83× with long-reads), with an estimated genome completeness of 99.27% based on the
98 Merqury²³ method. 97.6% of transcripts mapped to the genome along their full length with an identity
99 of ≥ 95% while a further 2.04% mapped with an identity between 90 - 95%. The high contiguity of the

100 assembly is underscored by the large number of transcripts mapped to single contigs (90.2%), and
101 since the proteins encoded by transcripts were consistently found in the predicted proteome, the latter
102 is also considered to be of high quality. We also conducted BUSCO analyses, with the foreknowledge
103 that genomic streamlining typical in Metamonada has led to the loss of many conserved proteins^{16,17,24}.
104 Our analyses show that previously completed metamonad genomes only encoded between 60% to 91%
105 BUSCO proteins, while *C. membranifera* exhibits a relatively high 89% (**Table 1, Supplementary**
106 **Information**). In any case, our coverage estimates for the *C. membranifera* genome for short and long
107 read sequencing technologies are substantially greater than those found to be sufficient to capture
108 genic regions that otherwise would have been missed (*i.e.*, coverage >52× for long reads and >60× for
109 short paired-end reads, see ref.²⁵). All these various data indicate that the draft genome of *C.*
110 *membranifera* is nearly complete; if any genomic regions are missing, they are likely confined to
111 difficult-to-sequence highly repetitive regions such as telomeres and centromeres.

112 Note that a previous study conducted a metagenomic assembly of a related
113 species, *Carpediemonas frisia*, together with its associated prokaryotic microbiota²⁶. For completeness,
114 we have included these data in our comparative genomic analyses (**Table 1, Supplementary**
115 **Information**), although we note that the *C. frisia* metagenomic bin is based on only short-read data
116 and might be partial.

117

118 **Extreme streamlining of the DNA replication apparatus in metamonads**

119 The first step in the replication of DNA is the assembly of ORC which serves to nucleate the pre-RC
120 formation. The initiator protein Orc1 first binds an origin of replication, followed by the recruitment
121 of Orc 2-6 proteins, which associate with chromatin²⁷. As the cell transitions to G1 phase, the
122 initiator Cdc6 binds to the ORC, forming a checkpoint control²⁸. Cdt1 then joins Cdc6, promoting the

123 loading of the replicative helicase MCM forming the pre-RC, a complex that remains inactive until
124 the onset of S-phase when the ‘firing’ factors are recruited to convert the pre-RC into the pre-IC³⁻⁵.
125 Additional factors join to form the RPC to stimulate replication elongation²⁹. While the precise
126 replisome protein complement varies somewhat between different eukaryotes, metamonads show
127 dramatic variation in ORC, pre-RC and replicative polymerases (**Fig. 1**). The presence-absence of
128 ORC and Cdc6 proteins is notably patchy across Metamonada. Strikingly, whereas all most
129 metamonads retain up to two paralogs of the core protein family Orc1/Cdc6 (here called Orc1 and
130 Orc1/Cdc6-like, **Supplementary Figure 1**), plus some orthologs of Orc 2-6, all these proteins are
131 absent in *C. membranifera* and *C. frisia* (**Fig. 1, Supplementary Table 1**). The lack of these proteins
132 in a eukaryote is unexpected and unprecedented, since their absence would be expected to make the
133 genome prone to DSBs and impair DNA replication, as well as interfere with other non-replicative
134 processes³⁰. To rule out false negatives, we conducted further analyses using metamonad-specific
135 HMMs (Hidden Markov Models), various other profile-based search strategies (**Supplementary**
136 **Information**), tBLASTn³¹ searches (*i.e.*, on the genome assembly and unassembled long-reads), and
137 applied HMMER³² on 6-frame assembly translations. These additional methods were sufficiently
138 sensitive to identify these proteins in all nuclear genomes we examined, with the exception of the
139 *Carpediemonas* species and the highly reduced, endosymbiotically-derived nucleomorphs of
140 cryptophytes and chlorarachniophytes (**Supplementary Information, Supplementary Table 1,**
141 **Supplementary Fig. 1 and 2**). *Carpediemonas* species are, therefore, the only known eukaryotes to
142 completely lack ORC and Cdc6.

143 **DNA damage repair systems have undergone several modifications**

144 DNA repair occurs continuously during the cell cycle depending on the type or specificity of the
145 lesion. Among the currently known mechanisms are base-excision repair (BER), nucleotide excision

146 repair (NER), mismatch repair (MMR), and double strand break repair, with the latter conducted by
147 either homologous recombination (HR), canonical non-homologous end joining (NHEJ) or alternative
148 end joining (a-EJ)^{8,14}. MMR can be coupled directly to replication or play a role in HR. MMR, BER
149 and NER are present in all studied taxa (**Supplementary Table 1**), although our analyses indicate that
150 damage sensing and downstream functions in NER seem to be modified in the metamonad taxa
151 Parabasalia and Fornicata due to the absence of the XPG and XPC sensor proteins.

152 Double strand breaks (DSBs) are extremely dangerous for cells and can occur as a result of
153 damaging agents or from self-inflicted cuts during DNA repair and meiosis. NHEJ requires the
154 heterodimer Ku70-Ku80 to recruit the catalytic kinase DNA-PKcs and accessory proteins.
155 Metamonads lack all of these proteins, as do a number of other eukaryotes investigated here and in
156 ref.³³. The a-EJ system seems to be fully present in metamonads like *C. membranifera*, partial in
157 others, and completely absent in parasitic diplomonads. NHEJ is thought to be the predominant
158 mechanism for repairing DSBs in eukaryotes³⁴, but since our analyses indicate this pathway is absent
159 in metamonads and a-EJ is highly mutagenic⁸, the HR pathway is likely to be essential for DSB repair
160 in most metamonads. Repair by the HR system occurs through multiple sub-pathways that are
161 influenced by the extent of the similarity of the DNA template or its flanking sequences to the
162 sequences near the break. HR complexes are recruited during DNA replication and transcription, and
163 utilize DNA, transcript-RNA or newly synthesized transcript-cDNA as a homologous template^{11,35-40}.
164 These complexes are formed by recombinases from the RecA/Rad51 family that interact with
165 members of the Rad52 family and chromatin remodeling factors of the SNF2/SWI2 sub-family^{41,42}.
166 Although the recombinases Rad51A-D are all present in most eukaryotes, we found a patchy
167 distribution in metamonads (**Supplementary Table 1, Supplementary Fig. 3**). All examined
168 Fornicata have lost the major recombinase Rad51A and have two paralogs of the meiosis-specific

169 recombinase Dmc1, as first noted in *Giardia intestinalis*⁴³. Dmc1 has been reported to provide high
170 stability to recombination due to strong D-loop resistance to strand dissociation⁴⁴. The recombination
171 mediator Rad52 is present in metamonads but Rad59 or Rad54 are not. Metamonads have no
172 components of an ISWI remodeling complex yet retain a reduced INO80 complex. Therefore,
173 replication fork progression and HR are likely to occur under the assistance of INO80 alone. HR
174 requires endonucleases and exonucleases, and our searches for proteins additional to those from the
175 MMR pathway revealed a gene expansion of the Flap proteins from the Rad2/XPG family in some
176 metamonads. We also found proteins of the PIF1 helicase family that encompasses homologs that
177 resolve R-loop structures, unwind DNA–RNA hybrids and assists in fork progression in regular
178 replication and HR^{45,46}. Phylogenetic analysis reveals that although *Carpodiemonas* species have
179 orthologs that branch within a metamonad group in the main PIF1 clade (**Fig. 2**), they also possess a
180 highly divergent clade of PIF1-like proteins. Each *Carpodiemonas* species has multiple copies of PIF1-
181 like proteins that have independently duplicated within each species; these may point to the *de novo*
182 emergence of specialized functions in HR and DNA replication for these proteins. Metamonads appear
183 capable of using all of the HR sub-pathways (*e.g.*, classical DSB repair, single strand annealing, break
184 induced replication), but these are modified (**Supplementary Table 1, Supplementary Figure 3**).
185 Overall, the presence-absence patterns of the orthologs involved in DSB repair in Fornicata point to
186 the existence of a highly specialized HR pathway which is presumably not only essential for the cell
187 cycle of metamonads but is also likely the major pathway for replication-related DNA repair and
188 recombination.

189 **Modified DSB damage response checkpoints in metamonads**

190 Checkpoints constitute a cascade of signaling events that delay replication until DNA lesions are
191 resolved¹³. The ATR-Chk1, ATM-Chk2 and DNA-PKcs pathways are activated by the interaction of

192 TopBP1 and the 9-1-1 complex (Rad9-Hus1-Rad1) for DNA repair regulation during replication stress
193 and response to DSBs⁴⁷. The ATR-Chk1 signaling pathway is the initial response to ssDNA damage
194 and is responsible for the coupling of DNA replication with mitosis, but when it is defective, the
195 ssDNA is converted into DSBs to activate the ATM-Chk2 pathway. The DNA-PKcs act as sensors of
196 DSBs to promote NHEJ, but we found no homologs of DNA-PKcs in metamonads (**Supplementary**
197 **Fig. 3**), which is consistent with the lack of a NHEJ repair pathway in the group. All the checkpoint
198 pathways described are present in humans and yeasts, while the distribution of core checkpoint
199 proteins in the remaining taxa is patchy. Notably, Fornicata lack several of the proteins thought to be
200 needed to activate the signaling kinase cascades and, while orthologs of ATM or ATR kinases are
201 present in some fornicates, there are no clear orthologs of Chk1 or Chk2 in metamonads except in
202 *Monocercomonoides exilis* (**Supplementary Table 1, Supplementary Fig. 3**). *Carpediemonas* species
203 and *Kipferlia bialata* contain ATM and ATR but lack Chk1, Chk2 and Rad9. Diplomonads possess
204 none of these proteins. The depletion of Chk1 has been shown to increase the incidence of
205 chromosomal breaks and mis-segregation⁴⁸. All these absences reinforce the idea that the checkpoint
206 controls in Fornicata are non-canonical.

207 **Reduction of mitosis and meiosis machinery in metamonads**

208 Eukaryotes synchronize cell cycle progression with chromosome segregation by a kinetochore based
209 signaling system called the spindle assembly checkpoint (SAC)^{49,50} that is ancestral to all eukaryotes
210 (**Fig. 3A, B**). Kinetochores primarily form microtubule attachments through the Ndc80 complex,
211 which is connected through a large network of structural subunits to a histone H3-variant CenpA that
212 is specifically deposited at centromeres¹². To prevent premature chromosome segregation, unattached
213 kinetochores catalyse the production of the Mitotic Checkpoint Complex (MCC)⁴⁹, a cytosolic
214 inhibitor of the Anaphase Promoting Complex/Cyclosome (APC/C), a large multi-subunit E3 ubiquitin

215 ligase that drives progression into anaphase by promoting the proteolysis of its substrates such as
216 various Cyclins⁵¹ (**Fig. 3A**). Our analysis indicates the reduction of ancestral complexity of these
217 proteins in metamonads (**Fig. 3C, Supplementary Table 1, Supplementary Fig. 4**). Surprisingly,
218 such reduction is most extensive in *Carpodiemonas* species. We found that most structural kinetochore
219 subunits, a microtubule plus-end tracking complex and all four subunits of the Ndc80 complex are
220 absent (**Fig. 3C, Supplementary Fig. 4**). None of our additional search strategies led to the
221 identification of Ndc80 complex members, making *Carpodiemonas* the only known eukaryotic lineage
222 without it, except for kinetoplastids, which appear to have lost the canonical kinetochore and replaced
223 it by an analogous molecular system, although there is still some controversy about this loss^{52,53}. With
224 such widespread absence of kinetochore components it might be possible that *Carpodiemonas*
225 underwent a similar replacement process to that of kinetoplastids⁵². We did however find a potential
226 candidate for the centromeric Histone H3-variant (CenpA) in *C. membranifera*. CenpA forms the basis
227 of the canonical kinetochore in most eukaryotes⁵⁴ (**Supplementary Fig. 5**). On the other hand, the
228 presence or absence of CenpA is often correlated with the presence/absence of its direct interactor
229 CenpC¹⁹. Similar to diplomonads, *C. membranifera* lacks CenpC and therefore the molecular network
230 associated with kinetochore assembly on CenpA chromatin may be very different.

231 Most metamonads encode all MCC components, but diplomonads lost the SAC response and
232 the full APC/C complex⁵⁵. In contrast, only *Carpodiemonas* species and *K. bialata* have MCC subunits
233 that contain the conserved short linear motifs to potentially elicit a canonical SAC signal^{51,56}
234 (**Supplementary Fig. 6**). Interestingly, not all of these motifs are present, and most are seemingly
235 degenerate compared to their counterparts in other eukaryotic lineages (**Supplementary Fig. 6C**).
236 Also, many other SAC-related genes are conserved, even in diplomonads (*e.g.*, Mad2, MadBub)⁵⁵.
237 Furthermore, the cyclins in *C. membranifera*, the main target of SAC signalling, have a diverged

238 destruction motif (D-box) in their N-termini (**Supplementary Fig. 6C**). Collectively, our observations
239 indicate that *Carpediemonas* species could elicit a functional SAC response, but whether this would be
240 kinetochore-based is unclear. Alternatively, SAC-related genes could have been repurposed for another
241 cellular function(s) as in diplomonads⁵⁵. Given that ORC has been observed to interact with the
242 kinetochore (throughout chromosome condensation and segregation), centrioles and promotes
243 cytokinesis³⁰, the lack of Ncd80 and ORC complexes suggest that *Carpediemonas* species possess
244 radically unconventional cell division systems.

245 Neither sexual nor parasexual processes have been directly observed in Metamonada⁴³.
246 Nonetheless, our surveys confirm the conservation of the key meiotic proteins in metamonads⁴³,
247 including Hap2 (for plasmogamy) and Gex1 (karyogamy). Unexpectedly, *Carpediemonas* species have
248 homologs from the tmcB family that acts in the cAMP signaling pathway specific for sexual
249 development in *Dictyostelium*⁵⁷, and sperm-specific channel subunits (*i.e.*, CatSper α , β , δ and γ)
250 reported previously only in Opisthokonta and three other protists⁵⁸. In opisthokonts, the CatSper
251 subunits enable the assembly of specialized Ca²⁺ influx channels and are involved in the signaling for
252 sperm maturation and motility⁵⁸. In *Carpediemonas*, the tmcB family and CatSper subunits could
253 similarly have a role in signaling and locomotion pathways required for a sexual cycle. As proteins in
254 the cAMP pathway and Ca²⁺ signaling cooperate to generate a variety of complex responses, the
255 presence of these systems in *Carpediemonas* species but absence in all other sampled metamonads is
256 intriguing and deserves further investigation. Even if these systems are not directly involved in a
257 sexual cycle, the presence of Hap2 and Gex1 proteins is strong evidence that *C. membranifera* can
258 reproduce sexually. Interestingly, based on the frequencies of single nucleotide polymorphisms, *C.*
259 *membranifera* is predicted to be haploid (**Supplementary Fig. 7**). If this is correct, its sexual

260 reproduction should include the formation of a zygote followed by a meiotic division to regain its
261 haploid state⁵⁹.

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

263 The unprecedented absence of many components of canonical DNA replication, repair, and
264 segregation systems in *Carpediemonas* species led us to investigate whether they had been replaced
265 by analogous systems acquired by lateral gene transfer (LGT) from viruses or prokaryotes. We
266 detected four Geminivirus-like replication initiation protein sequences in the *C. membranifera*
267 genome but not in *C. frisia*, and helitron-related helicase endonucleases in both *Carpediemonas*
268 genomes. All these genes were embedded in high-coverage eukaryotic scaffolds, yet all of them lack
269 introns and show no evidence of gene expression in the RNA-Seq data. As RNA was harvested from
270 log-phase actively replicating cell cultures, their lack of expression suggests it is unlikely that these
271 acquired proteins were coopted to function in the replication of the *Carpediemonas* genomes.
272 Nevertheless, the presence of Geminivirus protein-coding genes is intriguing as these viruses are
273 known, in other systems (*e.g.*, plants, insects), to alter host transcriptional controls and reprogram the
274 cell-cycle to induce the host DNA replication machinery^{60,61}. We also detected putative LGTs of
275 Endonuclease IV, RarA and RNase H1 from prokaryotes into a *Carpediemonas* ancestor
276 (**Supplementary Information, Supplementary Fig. 8, 9 and 10**). Of these, RarA is ubiquitous in
277 bacteria and eukaryotes and acts during replication and recombination in the context of collapsed
278 replication forks^{62,63}. Interestingly, *Carpediemonas* appears to have lost the eukaryotic ortholog, and
279 only retains the acquired prokaryotic-like RarA, a gene that is expressed (*i.e.*, transcripts are present
280 in the RNA-Seq data). RNase Hs are involved in the cleavage of RNA from RNA:DNA hybrid
281 structures that form during replication, transcription, and repair, and, while eukaryotes have a
282 monomeric RNase H1 and a heterotrimeric RNase H2, prokaryotes have either one or both types.

283 Eukaryotic RNase H1 removes RNA primers during replication and R-loops during transcription,
284 and also participates in HR-mediated DSB repair^{64,65}. The prokaryotic homologs have similar roles
285 during replication and transcription⁶⁶. *C. membranifera* lacks a typical eukaryotic RNase H1 but has
286 two copies of prokaryotic homologs. Both are located in scaffolds comprising intron-containing
287 genes and have RNA-Seq coverage, clearly demonstrating that they are not from prokaryotic
288 contaminants in the assembly.

289

290 **Discussion**

291 **Genome streamlining in metamonads**

292 The reductive evolution of the DNA replication and repair, and segregation systems and the low
293 retention of proteins in the BUSCO dataset in metamonads demonstrate that substantial gene loss has
294 occurred (**Supplementary information**), providing additional evidence for streamlining of gene
295 content prior to the last common ancestor of Metamonada¹⁵⁻¹⁷. However, the patchy distribution of
296 genes within the group suggests ongoing differential reduction in different metamonad groups. Such
297 reduction – especially the unprecedented complete absence of systems such as the ORC, Cdc6 and
298 kinetochore Ndc80 complexes in *Carpediemonas* species – demands an explanation. Whereas the loss
299 of genes from varied metabolic pathways is well known in lineages with different lifestyles⁶⁷⁻⁷², loss of
300 cell cycle, DNA damage sensing and repair genes in eukaryotes is extremely rare. New evidence from
301 yeasts of the genus *Hanseniaspora* suggests that the loss of proteins in these systems can lead to
302 genome instability and long-term hypermutation leading to high rates of sequence substitution⁶⁷. This
303 could also apply to metamonads, especially fornicates, which are well known to have undergone rapid
304 sequence evolution; these taxa form a highly divergent clade with very long branches in phylogenetic
305 trees^{20,73}. Most of the genes that were retained by Metamonada in the various pathways we examined

306 were divergent in sequence relative to homologs in other eukaryotes and many of the gene losses
307 correspond to proteins that are essential in model system eukaryotes. Gene essentiality appears to be
308 relative and context-dependent, and some studies have shown that the loss of ‘indispensable’ genes
309 could be permitted by evolving divergent pathways that provide similar activities via chromosome
310 stoichiometry changes and compensatory gene loss^{67-69,74}.

311 The patchy distribution of genes from different ancestral eukaryotic pathways suggests that the
312 last common ancestor of Metamonada had a broad gene repertoire for maintaining varied metabolic
313 functions under fluctuating environmental conditions offered by diverse oxygen-depleted habitats.
314 Although the loss of proteins and genomic streamlining are well known in parasitic diplomonads^{15,16},
315 the Fornicata, as a whole, tend to have a reduced subset of the genes that are commonly found in core
316 eukaryotic pathways. In general, such gene content reduction can partially be explained as the result of
317 historical and niche-specific adaptations⁷⁵. Yet, given that 1) genome maintenance mostly depends on
318 the cell cycle checkpoints, DNA repair pathways, and their interactions^{14,76}, 2) the lack of several
319 proteins related to these pathways that were present in the last common ancestor of metamonads, 3)
320 aneuploidy and high overall rates of sequence evolution have been observed in metamonads^{77,78}, and,
321 4) the loss of DNA repair genes can be associated with substantial gene loss and sequence instability
322 that apparently boosts the rates of sequence evolution⁶⁷, it is likely that genome evolution in the
323 Fornicata clade has been heavily influenced by their error-prone DNA maintenance mechanisms.

324
325 **Non-canonical replication initiation and replication licensing in *Carpediemonas*.**

326 Origin-independent replication has been observed in the context of DNA repair (reviewed in ref.¹⁰) and
327 in origin-deficient or -depleted chromosomes in yeast⁷⁹. These studies have highlighted the lack of (or
328 reduction in) the recruitment of ORC and Cdc6 onto the DNA, but no study to date has documented

329 regular eukaryotic DNA replication in the absence of genes encoding these proteins. While it is
330 possible that extremely divergent versions of ORC and Cdc6 are governing the recognition of origins
331 of replication and replication licensing in *Carpediemonas* species, we have no evidence for this.
332 Instead, our findings suggest the existence of an as-yet undiscovered underlying eukaryotic system that
333 can accomplish eukaryotic DNA replication initiation and licensing. The existence of such a system
334 has in fact already been suspected given that: 1) Orc1- or Orc2-depleted human cells and mouse-Orc1
335 and fruit-fly ORC mutants are viable and capable of undergoing replication and endoreplication⁸⁰⁻⁸³,
336 and 2) origin-independent replication at the chromosome level has been reported^{79,84,85}. We propose
337 that *Carpediemonas* species utilize an alternative DNA replication system based on a Dmc1-dependent
338 HR mechanism that is origin-independent and mediated by RNA:DNA hybrids. Here we summarize
339 evidence that such a mechanism is possible based on what is known in model systems and present a
340 hypothetical model as to how it might occur in *Carpediemonas*.

341 During replication and transcription, the HR complexes, RNase H1 and RNA-interacting
342 proteins are recruited onto the DNA to assist in its repair^{36,37,86}. Remarkably, experiments show that
343 HR is able to carry out full genome replication in archaea, bacteria, viruses, and linear mtDNA^{85,87-89},
344 with replication fork progression rates that are comparable to those of regular replication⁹⁰. A variety
345 of *cis* and *trans* homologous sequences (*e.g.*, chromatids, transcript-RNA or -cDNA) can be used as
346 templates^{27,36,40}, and their length as well as the presence of one or two homologous ends likely
347 influence a recombination execution checkpoint that decides which HR sub-pathway is utilized⁹¹. For
348 example, in the absence of a second homologous end, HR by Rad51-dependent break-induced
349 replication (BIR) can either use a newly synthesized DNA strand or independently invade donor
350 sequences, such that the initial strand invasion intermediate creates a migrating D-loop and DNA is
351 synthesized conservatively^{27,91,92}. Studies have found that BIR does not require the assembly of an

352 ORC complex and Cdc6 but the recruitment of the Cdc7, loading of MCM helicase, firing factors and
353 replicative polymerases are needed for assembling the pre-RC complex^{27,91}. The requirement of MCM
354 for BIR was questioned, as PIF1 helicase was found to be essential for long-range BIR⁹³. However,
355 recent evidence shows that MCM is typically recruited for unwinding DNA strands during HR^{94,95} and
356 is likely needed together with PIF1 to enhance processivity. All these proteins are also suspected to
357 operate during origin-independent transcription-initiated replication (TIR), a still-enigmatic
358 mechanism that is triggered by R-loops resulting from RNA:DNA hybrids during transcription^{10,11,96}.

359 Considering the complement of proteins in *Carpediemonas* species discussed above, and that
360 RNA:DNA hybrids are capable of promoting origin-independent replication in model systems^{11,39,97},
361 we suggest that a Dmc1-dependent HR replication mechanism is enabled by excess of RNA:DNA
362 hybrids in these organisms. In such a system, DSBs generated in stressed transcription-dependent R-
363 loops could be repaired by HR with either transcript-RNA- or transcript-cDNA-templates and the *de*
364 *novo* assembly of the replisome as in BIR (**Fig. 4**). The establishment of a replication fork could be
365 favored by the presence of *Carpediemonas*-specific PIF1-like homologs, as these raise the possibility
366 of the assembly of a multimeric PIF1 helicase with increased capability to bind multiple sites on the
367 DNA, thereby facilitating DNA replication processivity and regulation⁴⁵. Note that the foregoing
368 mechanisms will work even if *Carpediemonas* species are haploid as seems likely based on the SNP
369 data. The loss of Rad51A and the duplication of Dmc1 recombinases suggests that a Dmc1-
370 dependent HR mechanism was likely enabled in the last common ancestor of Fornicata and this
371 mechanism may have become the predominant replication pathway in the *Carpediemonas* lineage
372 after its divergence from the other fornicates, ultimately leading to the loss of ORC and Cdc6
373 proteins.

374

375 **The impact of cell cycle dysregulation on genome evolution.**

376 DNA replication licensing and firing are temporally separated (*i.e.*, they occur at G1 and S phases
377 respectively) and are the principal ways to counteract damaging over-replication⁶. As S-phase is
378 particularly vulnerable to DNA errors and lesions, its checkpoints are likely more important for
379 preventing genome instability than those of G1, G2 or SAC⁹⁸. Dysregulation is anticipated if no
380 ORC/Cdc6 are present as licensing would not take place and replication would be blocked²⁸. Yet this
381 clearly does not happen in *Carpediemonas*. This implies that during late G1 phase, activation by
382 loading the MCM helicase has to occur by an alternative mechanism that is still unknown but might
383 already be in place in eukaryotes. Such a mechanism has long been suspected as it could explain the
384 over-abundance and distribution patterns of MCM on the DNA (*i.e.*, the MCM paradox; reviewed in
385 ⁹⁹).

386 In terms of the regulation of M-phase progression, the extremely divergent nature of the
387 kinetochore in *C. membranifera* could suggest that it uses different mechanisms to execute mitosis
388 and meiosis. It is known that in *Carpediemonas*-related fornicates such as retortamonads and in
389 diplomonads, chromosome segregation proceeds inside a persisting nuclear envelope, with the aid of
390 intranuclear microtubules, but with the mitotic spindle nucleated outside the nucleus (*i.e.*, semi-open
391 mitosis)⁷⁸. Although mitosis in *Carpediemonas* has not been directly observed, these organisms may
392 also possess a semi-open mitotic system such as the ones found in other fornicates. Yet how the
393 *Carpediemonas* kinetochore functions in the complete absence of the microtubule-binding Ndc80
394 complex remains a mystery; it is possible that, like in kinetoplastids⁴⁸, other molecular complexes have
395 evolved in this lineage that fulfill the roles of Ndc80 and other kinetochore complexes.

396 Interestingly, a potential repurposing of SAC proteins seems to have occurred in the
397 diplomonad *G. intestinalis*, as it does not arrest under treatment with microtubule-destabilizing drugs

398 and Mad2 localizes to a region of the intracytoplasmic axonemes of the caudal flagella⁵⁵. Other
399 diplomonads have a similar SAC protein complement that may have a similar non-canonical function.
400 In contrast to diplomonads, our investigations (**Fig. 3**) suggest that *Carpediemonas* species could elicit
401 a functional SAC response, although microtubule-disrupting experiments during mitosis will be
402 needed to prove its existence.

403 In addition to the aforementioned apparent dysregulation of checkpoint controls in
404 *Carpediemonas* species, alternative mechanisms for chromosome condensation, spindle attachment,
405 sister chromatid cohesion, cytokinesis, heterochromatin formation, and silencing and transcriptional
406 regulation can also be expected in this organism due to the absence of ORC and Cdc6 (reviewed in
407 refs^{30,100,101}). All of the absences of canonical eukaryotic systems we have described for
408 *Carpediemonas* suggest that a radically different cell cycle has evolved in this free-living protistan
409 lineage. This underscores the fact that our concepts of universality and essentiality rely on studies of
410 a very small subset of organisms. The development of *Carpediemonas* as a model system thus has
411 great potential to enhance our understanding of fundamental DNA replication, repair and cell cycle
412 processes. It could even reveal widely conserved alternative, but as-yet unknown, mechanisms
413 underpinning the evolutionary plasticity of these systems across the eukaryote tree of life.

414

415 **Methods**

416 **Sequencing, assembly, and protein prediction for *C. membranifera***

417 DNA and RNA were isolated from log-phase cultures of *C. membranifera* BICM strain (see details in
418 **Supplementary Information**). Sequencing employed Illumina short paired-end and long read
419 (Oxford Nanopore MinION) technologies. For Illumina, extracted, purified DNA and RNA (*i.e.*,
420 cDNA) were sequenced on the Hiseq 2000 (150 x 2 paired-end) at the Genome Québec facility.

421 Illumina reads were quality trimmed (Q=30) and filtered for length (>40 bp) with Trimmomatic¹⁰².
422 For MinION, the library was prepared using the 1D native barcoding genomic DNA (SQK-LSK108
423 with EXP-NBD103) protocol (NBE_9006_v103_revP_21Dec2016). The final library (1070 ng) was
424 loaded on a R9.4 flow cell and sequenced for 48 h on the MinION Mk1B nanopore sequencer. The
425 long reads were base-called and trimmed with Albacore v2.3.3 (www.nanoporetech.com) and
426 Porechop v0.2.3 (www.github.com/rrwick/Porechop), respectively. ABruijn v1.0
427 (www.github.com/fenderglass/Flye/releases/tag/1.0) with default parameters and max genome size of
428 30Mb produced an assembly that was polished with Nanopolish v0.10.1¹⁰³. The latter was iteratively
429 error-corrected with the genomic paired-end Illumina reads using Unicycler¹⁰⁴. The identification and
430 removal of prokaryotic contigs was assisted by BLASTn searches against the nt database. Read-depth
431 coverage at each position of the genomic scaffolds were obtained with samtools¹⁰⁵ and mosdepth
432 v0.2.5¹⁰⁶.

433 RNA-Seq reads were used for genome-independent assessments of the presence of the proteins
434 of interest and to generate intron junction hints for gene prediction. For the independent assessments
435 we obtained both a *de novo* and a genome-guided transcriptome assembly with Trinity v2.5.0¹⁰⁷. Open
436 reading frames were translated with TransDecoder v5.5.0 (www.github.com/TransDecoder) and were
437 included in all of our analyses. Gene predictions were carried out as follows: repeat libraries were
438 obtained and masked with RepeatModeler 1.0 and RepeatMasker (<http://www.repeatmasker.org>).
439 Then, RNA-Seq reads were mapped onto the assembly using Hisat2¹⁰⁸, generating a bam file for
440 GenMarkET¹⁰⁹. This resulted in a list of intron hints used to train Augustus v3.2.3¹¹⁰. The genome-
441 guided assembled transcriptome, genomic scaffolds and the newly predicted proteome were fed into
442 the PASA pipeline¹¹¹ to yield a more accurate set of predicted proteins. Finally, the predicted proteome
443 was manually curated for the proteins of interest.

444 **Genome size, completeness, and ploidy assessments**

445 We estimated the completeness of the draft genome by 1) using the k-mer based and reference free
446 method Merqury²³, 2) calculating the percentage of transcripts that aligned to the genome, and 3)
447 employing the BUSCO¹¹² framework. For method 1, all paired-end reads were used to estimate the
448 best k-mer and create ‘meryl’ databases necessary to apply Merqury²³. For method 2, transcripts were
449 mapped onto the genome using BLASTn and exonerate¹¹³. For method 3, the completeness of the
450 draft genome was evaluated in a comparative setting by including the metamonads and using the
451 universal single copy orthologs (BUSCO) from the Eukaryota (odb9) and protist databases
452 (<https://busco.ezlab.org/>), which contain 303 and 215 proteins, respectively. Each search was run
453 separately on the assembly and the predicted proteome for all these taxa. Unfortunately, both
454 BUSCO database searches yielded false negatives in that several conserved proteins publicly
455 reported for *T. vaginalis*, *G. intestinalis* and *Spironucleus salmonicida* were not detected due to the
456 extreme divergence of metamonad homologs. Therefore, genome completeness was re-assessed with
457 a phylogeny-guided search (**Supplementary Information**).

458 The ploidy of *C. membranifera* was inferred by *i*) counting k-mers with Merqury²³, and *ii*)
459 mapping 613,266,290 Illumina short reads to the assembly with Bowtie 2.3.1¹¹⁴ and then using
460 ploidyNGS¹¹⁵ to calculate the distribution of allele frequencies across the genome. A site was deemed
461 to be heterozygous if at least two different bases were present and there were at least two reads with
462 the different bases. Positions with less than 10× coverage were ignored.

463

464 **Functional annotation of the predicted proteins**

465 Our analyses included the genomes and predicted proteomes of *C. membranifera* (reported here) as
466 well as publicly available data for nine additional metamonads and eight other eukaryotes

467 representing diverse groups across the eukaryotic tree of life (**Fig. 1, Supplementary Information**).

468 Orthologs from each of these 18 predicted proteomes were retrieved for the assessment of core

469 cellular pathways, such as DNA replication and repair, mitosis and meiosis and cell cycle

470 checkpoints. For *C. membranifera*, we included the predicted proteomes derived from the assembly

471 plus the 6-frame translated transcriptomes. Positive hits were manually curated in the *C.*

472 *membranifera* draft genome. A total of 367 protein queries were selected based on an extensive

473 literature review and prioritizing queries from taxa in which they had been experimentally

474 characterized. The identification of orthologs was as described for the BUSCO proteins but using

475 these 367 queries for the initial BLASTp (**Supplementary Information**), except for kinetochore

476 (KT), Spindle assembly check point (SAC) and anaphase-promoting complex-related genes (APC/C).

477 For these, previously published refined HMMs with cut-offs specific to each orthologous group

478 (see⁵⁸) were used to query the proteomes with HMMER v3.1b2³². A multiple sequence alignment

479 that included the newly-found hits was subsequently constructed with MAFFT v7.310¹¹⁶ and was

480 used in HMM searches for more divergent homologs. This process was iterated until no new

481 significant hits could be found. As we were unable to retrieve orthologs of a number of essential

482 proteins in the *C. membranifera* and *C. frisia* genomes, we embarked on additional more sensitive

483 strategies to detect them using multiple different HMMs based on aligned homologs from archaea,

484 metamonads, and broad samplings of taxa. Individual PFAM domains were searched for in the

485 genomes, proteome and transcriptomes with e-value thresholds of 10^{-3} (**Supplementary**

486 **Information**). To rule out that failure to detect these proteins was due to insufficient sensitivity of

487 our methods when applied them to highly divergent taxa, we queried 22 extra eukaryotic genomes

488 with demonstrated high rates of sequence evolution, genome streamlining or unusual genomic

489 features (**Supplementary Table 1, Supplementary Information**). Possible non-predicted or mis-

490 predicted genes were investigated using tBLASTn searches of the genomic scaffolds and
491 unassembled reads and 6-frame translation searches with HMMER. Also, as DNA replication and
492 repair genes could have been acquired by lateral gene transfer into *Carpodomonas* species from
493 prokaryotes or viruses, proteins from the DNA replication and repair categories whose best matches
494 were to prokaryotic and viral homologs were subjected to phylogenetic analysis using the methods
495 described for the phylogeny-guided BUSCO analysis and using substitution models specified in the
496 legend of each tree (**Supplementary Information**).

497 **Data availability**

498 Genome assembly will be available at NCBI under BioProject PRJNA719540, biosample number
499 SAMN18612951, accession numbers <XXXX>.

500

501

502

503 **References**

- 504 1 Yeeles, J. T., Deegan, T. D., Janska, A., Early, A. & Diffley, J. F. Regulated eukaryotic DNA
505 replication origin firing with purified proteins. *Nature* **519**, 431-435 (2015).
- 506 2 Parker, M. W., Botchan, M. R. & Berger, J. M. Mechanisms and regulation of DNA
507 replication initiation in eukaryotes. *Crit. Rev. Biochem. Mol. Biol.* **52**, 107-144 (2017).
- 508 3 Shen, Z. & Prasanth, S. G. Emerging players in the initiation of eukaryotic DNA replication.
509 *Cell Div* **7**, 22 (2012).
- 510 4 Burgers, P. M. J. & Kunkel, T. A. Eukaryotic DNA replication fork. *Annu. Rev. Biochem.* **86**,
511 417-438 (2017).

- 512 5 Riera, A. *et al.* From structure to mechanism-understanding initiation of DNA replication.
513 *Genes Dev.* **31**, 1073-1088 (2017).
- 514 6 Reusswig, K. U. & Pfander, B. Control of eukaryotic DNA replication initiation-mechanisms
515 to ensure smooth transitions. *Genes (Basel)* **10** (2019).
- 516 7 Waters, L. S. *et al.* Eukaryotic translesion polymerases and their roles and regulation in DNA
517 damage tolerance. *Microbiol. Mol. Biol. Rev.* **73**, 134-154 (2009).
- 518 8 Chang, H. H. Y., Pannunzio, N. R., Adachi, N. & Lieber, M. R. Non-homologous DNA end
519 joining and alternative pathways to double-strand break repair. *Nat. Rev. Mol. Cell Biol.* **18**,
520 495-506 (2017).
- 521 9 Wright, W. D., Shah, S. S. & Heyer, W. D. Homologous recombination and the repair of
522 DNA double-strand breaks. *J. Biol. Chem.* **293**, 10524-10535 (2018).
- 523 10 Ravoityte, B. & Wellinger, R. E. Non-canonical replication initiation: you're fired! *Genes*
524 *(Basel)* **8** (2017).
- 525 11 Stuckey, R., Garcia-Rodriguez, N., Aguilera, A. & Wellinger, R. E. Role for RNA:DNA
526 hybrids in origin-independent replication priming in a eukaryotic system. *Proc. Natl. Acad.*
527 *Sci. U.S.A* **112**, 5779-5784 (2015).
- 528 12 Musacchio, A. & Desai, A. A molecular view of kinetochore assembly and function. *Biology*
529 *(Basel)* **6** (2017).
- 530 13 Hustedt, N., Gasser, S. M. & Shimada, K. Replication checkpoint: tuning and coordination of
531 replication forks in S phase. *Genes (Basel)* **4**, 388-434 (2013).
- 532 14 Hakem, R. DNA-damage repair; the good, the bad, and the ugly. *EMBO J.* **27**, 589-605
533 (2008).

- 534 15 Adam, R. D. *et al.* Genome sequencing of *Giardia lamblia* genotypes A2 and B isolates (DH
535 and GS) and comparative analysis with the genomes of genotypes A1 and E (WB and Pig).
536 *Genome Biol. Evol.* **5**, 2498-2511 (2013).
- 537 16 Xu, F. *et al.* The genome of *Spironucleus salmonicida* highlights a fish pathogen adapted to
538 fluctuating environments. *PLoS Genet.* **10**, e1004053 (2014).
- 539 17 Tanifuji, G. *et al.* The draft genome of *Kipferlia bialata* reveals reductive genome evolution
540 in fornicate parasites. *PLoS One* **13**, e0194487 (2018).
- 541 18 Ocaña-Pallares, E. *et al.* Origin recognition complex (ORC) evolution is influenced by global
542 gene duplication/loss patterns in eukaryotic genomes. *Genome Biol. Evol.* **12**, 3878-3889
543 (2020).
- 544 19 van Hooff, J. J., Tromer, E., van Wijk, L. M., Snel, B. & Kops, G. J. Evolutionary dynamics
545 of the kinetochore network in eukaryotes as revealed by comparative genomics. *EMBO Rep.*
546 **18**, 1559-1571 (2017).
- 547 20 Hampl, V. *et al.* Phylogenomic analyses support the monophyly of Excavata and resolve
548 relationships among eukaryotic "supergroups". *Proc. Natl. Acad. Sci. U.S.A.* **106**, 3859-3864
549 (2009).
- 550 21 Karnkowska, A. *et al.* A eukaryote without a mitochondrial organelle. *Curr. Biol.* **26**, 1274-
551 1284 (2016).
- 552 22 Simpson, A. G. B. & Patterson, D. J. The ultrastructure of *Carpodionomonas membranifera*
553 (Eukaryota) with reference to the "excavate hypothesis". *Eur. J. Protistol.* **35**, 353-370
554 (1999).
- 555 23 Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality,
556 completeness, and phasing assessment for genome assemblies. *Genome Biol.* **21**, 245 (2020).

- 557 24 Morrison, H. G. *et al.* Genomic minimalism in the early diverging intestinal parasite *Giardia*
558 *lamblia*. *Science* **317** (2007).
- 559 25 Ebbert, M. T. W. *et al.* Systematic analysis of dark and camouflaged genes reveals disease-
560 relevant genes hiding in plain sight. *Genome Biol.* **20**, 97 (2019).
- 561 26 Hamann, E. *et al.* Syntrophic linkage between predatory *Carpediemonas* and specific
562 prokaryotic populations. *ISME J* **11**, 1205-1217 (2017).
- 563 27 Lydeard, J. R. *et al.* Break-induced replication requires all essential DNA replication factors
564 except those specific for pre-RC assembly. *Genes Dev.* **24**, 1133-1144 (2010).
- 565 28 Liu, J. *et al.* Structure and function of Cdc6/Cdc18: implications for origin recognition and
566 checkpoint control. *Mol. Cell* **6**, 637-648 (2000).
- 567 29 Georgescu, R. E. *et al.* Reconstitution of a eukaryotic replisome reveals suppression
568 mechanisms that define leading/lagging strand operation. *Elife* **4**, e04988 (2015).
- 569 30 Popova, V. V., Brechalov, A. V., Georgieva, S. G. & Kopytova, D. V. Nonreplicative
570 functions of the origin recognition complex. *Nucleus* **9**, 460-473 (2018).
- 571 31 Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421
572 (2009).
- 573 32 Eddy, S. R. Accelerated profile HMM searches. *PLoS Comp. Biol.* **7**, e1002195 (2011).
- 574 33 Nenarokova, A. *et al.* Causes and effects of loss of classical non-homologous end joining
575 pathway in parasitic eukaryotes. *MBio* (2019).
- 576 34 van den Berg, J. *et al.* DNA end-resection in highly accessible chromatin produces a toxic
577 break. *BioRxiv* (2019).
- 578 35 Wei, L., Levine, A. S. & Lan, L. Transcription-coupled homologous recombination after
579 oxidative damage. *DNA Repair* **44**, 76-80 (2016).

- 580 36 Meers, C., Keskin, H. & Storici, F. DNA repair by RNA: templated, or not templated, that is
581 the question. *DNA Repair (Amst)* **44**, 17-21 (2016).
- 582 37 Aymard, F. *et al.* Transcriptionally active chromatin recruits homologous recombination at
583 DNA double-strand breaks. *Nat. Struct. Mol. Biol.* **21**, 366-374 (2014).
- 584 38 Wei, L. *et al.* DNA damage during the G0/G1 phase triggers RNA-templated, Cockayne
585 syndrome B-dependent homologous recombination. *Proc. Natl. Acad. Sci. U.S.A* **112**, E3495-
586 3504 (2015).
- 587 39 Keskin, H. *et al.* Transcript-RNA-templated DNA recombination and repair. *Nature* **515**, 436-
588 439 (2014).
- 589 40 Storici, F., Bebenek, K., Kunkel, T. A., Gordenin, D. A. & Resnick, M. A. RNA-templated
590 DNA repair. *Nature* **447**, 338-341 (2007).
- 591 41 Ceballos, S. J. & Heyer, W. D. Functions of the Snf2/Swi2 family Rad54 motor protein in
592 homologous recombination. *Biochim. Biophys. Acta.* **1809**, 509-523 (2011).
- 593 42 Mazin, A. V. & Mazina, O. M. in *Molecular Life Sciences: An Encyclopedic Reference* (eds
594 Robert D. Wells, Judith S. Bond, Judith Klinman, & Bettie Sue Siler Masters) 1009-1016
595 (Springer New York, 2018).
- 596 43 Ramesh, M. A., Malik, S. B. & Logsdon, J. M., Jr. A phylogenomic inventory of meiotic
597 genes; evidence for sex in *Giardia* and an early eukaryotic origin of meiosis. *Curr. Biol.* **15**,
598 185-191 (2005).
- 599 44 Bugreev, D. V. *et al.* The resistance of DMC1 D-loops to dissociation may account for the
600 DMC1 requirement in meiosis. *Nat. Struct. Mol. Biol.* **18**, 56-60 (2011).
- 601 45 Byrd, A. K. & Raney, K. D. Structure and function of Pif1 helicase. *Biochem. Soc. Trans.* **45**,
602 1159-1171 (2017).

- 603 46 Hill, J., Eickhoff, P., Drury, L. S., Costa, A. & Diffley, J. F. X. The eukaryotic replisome
604 requires an additional helicase to disarm dormant replication origins. *BioRxiv*,
605 2020.2009.2017.301366 (2020).
- 606 47 Blackford, A. N. & Jackson, S. P. ATM, ATR, and DNA-PK: the trinity at the heart of the
607 DNA damage response. *Mol. Cell.* **66**, 801-817 (2017).
- 608 48 Calzetta, N. L., Gonzalez Besteiro, M. A. & Gottifredi, V. Mus81-Eme1-dependent aberrant
609 processing of DNA replication intermediates in mitosis impairs genome integrity. *Sci. Adv.* **6**
610 (2020).
- 611 49 Sacristan, C. & Kops, G. J. Joined at the hip: kinetochores, microtubules, and spindle
612 assembly checkpoint signaling. *Trends Cell Biol.* **25**, 21-28 (2015).
- 613 50 Kops, G. J. P. L., Snel, B. & Tromer, E. C. Evolutionary dynamics of the spindle assembly
614 checkpoint in eukaryotes. *Curr. Biol.* **30**, R589-R602 (2020).
- 615 51 Alfieri, C., Zhang, S. & Barford, D. Visualizing the complex functions and mechanisms of the
616 anaphase promoting complex/cyclosome (APC/C). *Open Biol.* **7** (2017).
- 617 52 Akiyoshi, B. & Gull, K. Discovery of unconventional kinetochores in kinetoplastids. *Cell*
618 **156**, 1247-1258 (2014).
- 619 53 D'Archivio, S. & Wickstead, B. *Trypanosome* outer kinetochore proteins suggest conservation
620 of chromosome segregation machinery across eukaryotes. *J. Cell Biol.* **216**, 379-391 (2017).
- 621 54 Drinnenberg, I. A., Henikoff, S. & Malik, H. S. Evolutionary turnover of kinetochore
622 proteins: a ship of theseus? *Trends Cell Biol.* **26**, 498-510 (2016).
- 623 55 Markova, K. *et al.* Absence of a conventional spindle mitotic checkpoint in the binucleated
624 single-celled parasite *Giardia intestinalis*. *Eur. J. Cell Biol.* **95**, 355-367 (2016).

- 625 56 Tromer, E., Bade, D., Snel, B. & Kops, G. J. Phylogenomics-guided discovery of a novel
626 conserved cassette of short linear motifs in BubR1 essential for the spindle checkpoint. *Open*
627 *Biol.* **6** (2016).
- 628 57 Muramoto, T., Takeda, S., Furuya, Y. & Urushihara, H. Reverse genetic analyses of gamete-
629 enriched genes revealed a novel regulator of the cAMP signaling pathway in *Dictyostelium*
630 *discoideum*. *Mech. Dev.* **122**, 733-743 (2005).
- 631 58 Cai, X., Wang, X. & Clapham, D. E. Early evolution of the eukaryotic Ca²⁺ signaling
632 machinery: conservation of the CatSper channel complex. *Mol. Biol. Evol.* **31**, 2735-2740
633 (2014).
- 634 59 von Dassow, P. & Montresor, M. Unveiling the mysteries of phytoplankton life cycles:
635 patterns and opportunities behind complexity. *J. Plankton Res.* **33**, 3-12 (2010).
- 636 60 Hanley-Bowdoin, L., Bejarano, E. R., Robertson, D. & Mansoor, S. Geminiviruses: masters at
637 redirecting and reprogramming plant processes. *Nat. Rev. Microbiol.* **11**, 777-788 (2013).
- 638 61 He, Y.-Z. *et al.* A plant DNA virus replicates in the salivary glands of its insect vector via
639 recruitment of host DNA synthesis machinery. *Proceedings of the National Academy of*
640 *Sciences* **117**, 16928-16937 (2020).
- 641 62 Romero, H. *et al.* Single molecule tracking reveals functions for RarA at replication forks but
642 also independently from replication during DNA repair in *Bacillus subtilis*. *Sci. Rep.* **9**, 1997
643 (2019).
- 644 63 Yoshimura, A., Seki, M. & Enomoto, T. The role of WRNIP1 in genome maintenance. *Cell*
645 *Cycle* **16**, 515-521 (2017).
- 646 64 Parajuli, S. *et al.* Human ribonuclease H1 resolves R-loops and thereby enables progression
647 of the DNA replication fork. *J. Biol. Chem.* **292**, 15216-15224 (2017).

- 648 65 Posse, V. *et al.* RNase H1 directs origin-specific initiation of DNA replication in human
649 mitochondria. *PLoS Genet.* **15**, e1007781 (2019).
- 650 66 Tadokoro, T. & Kanaya, S. Ribonuclease H: molecular diversities, substrate binding domains,
651 and catalytic mechanism of the prokaryotic enzymes. *FEBS J.* **276**, 1482-1493 (2009).
- 652 67 Steenwyk, J. L. *et al.* Extensive loss of cell-cycle and DNA repair genes in an ancient lineage
653 of bipolar budding yeasts. *PLoS Biol.* **17**, e3000255 (2019).
- 654 68 Grohme, M. A. *et al.* The genome of *Schmidtea mediterranea* and the evolution of core
655 cellular mechanisms. *Nature* **554**, 56-61 (2018).
- 656 69 Sekelsky, J. DNA repair in *Drosophila*: mutagens, models, and missing genes. *Genetics* **205**,
657 471-490 (2017).
- 658 70 Corradi, N. Microsporidia: eukaryotic intracellular parasites shaped by gene loss and
659 horizontal gene transfers. *Annu. Rev. Microbiol.* **69**, 167-183 (2015).
- 660 71 Galindo, L. J. *et al.* Evolutionary genomics of *Metchnikovella incurvata* (Metchnikovellidae):
661 an early branching microsporidium. *Genome Biol. Evol.* **10**, 2736-2748 (2018).
- 662 72 Albalat, R. & Canestro, C. Evolution by gene loss. *Nat. Rev. Genet.* **17**, 379-391 (2016).
- 663 73 Roger, A. J., Kolisko, M. & Simpson, A. G. B. in *Evolution of Virulence in Eukaryotic*
664 *Microbes* 44-69 (2013).
- 665 74 Rancati, G. *et al.* Aneuploidy underlies rapid adaptive evolution of yeast cells deprived of a
666 conserved cytokinesis motor. *Cell* **135**, 879-893 (2008).
- 667 75 Mendonca, A. G., Alves, R. J. & Pereira-Leal, J. B. Loss of genetic redundancy in reductive
668 genome evolution. *PLoS Comput. Biol.* **7**, e1001082 (2011).
- 669 76 Lindahl, T. & Wood, R. D. Quality control by DNA repair. *Science* **286**, 1897-1905 (1999).

- 670 77 Tumova, P., Uzlikova, M., Jurczyk, T. & Nohynkova, E. Constitutive aneuploidy and
671 genomic instability in the single-celled eukaryote *Giardia intestinalis*. *MicrobiologyOpen* **5**,
672 560-574 (2016).
- 673 78 Kulda, J., Nohýnková, E. & Čepička, I. in *Handbook of the Protists* (eds John M. Archibald
674 *et al.*) 1-32 (Springer International Publishing, 2017).
- 675 79 Bogenschutz, N. L., Rodriguez, J. & Tsukiyama, T. Initiation of DNA replication from non-
676 canonical sites on an origin-depleted chromosome. *PLoS One* **9**, e114545 (2014).
- 677 80 Shibata, E. *et al.* Two subunits of human ORC are dispensable for DNA replication and
678 proliferation. *Elife* **5** (2016).
- 679 81 Bell, S. P. Rethinking origin licensing. *Elife* **6** (2017).
- 680 82 Park, S. Y. & Asano, M. The origin recognition complex is dispensable for endoreplication in
681 *Drosophila*. *Proc. Natl. Acad. Sci. U.S.A.* **105**, 12343-12348 (2008).
- 682 83 Okano-Uchida, T. *et al.* Endoreduplication of the mouse genome in the absence of ORC1.
683 *Genes Dev.* **32**, 978-990 (2018).
- 684 84 Theis, J. F. *et al.* The DNA damage response pathway contributes to the stability of
685 chromosome III derivatives lacking efficient replicators. *PLoS Genet.* **6**, e1001227 (2010).
- 686 85 Hawkins, M., Malla, S., Blythe, M. J., Nieduszynski, C. A. & Allers, T. Accelerated growth
687 in the absence of DNA replication origins. *Nature* **503**, 544-547 (2013).
- 688 86 Bader, A. S., Hawley, B. R., Wilczynska, A. & Bushell, M. The roles of RNA in DNA
689 double-strand break repair. *Br. J. Cancer* **122**, 613-623 (2020).
- 690 87 Gillespie, K. A., Mehta, K. P., Laimins, L. A. & Moody, C. A. Human papillomaviruses
691 recruit cellular DNA repair and homologous recombination factors to viral replication centers.
692 *J. Virol.* **86**, 9520-9526 (2012).

- 693 88 Kogoma, T. Stable DNA replication: interplay between DNA replication, homologous
694 recombination, and transcription. *Microbiol. Mol. Biol. Rev.* **61**, 212-238 (1997).
- 695 89 Gerhold, J. M. *et al.* Replication intermediates of the linear mitochondrial DNA of *Candida*
696 *parapsilosis* suggest a common recombination based mechanism for yeast mitochondria. *J.*
697 *Biol. Chem.* **289**, 22659-22670 (2014).
- 698 90 Malkova, A., Naylor, M. L., Yamaguchi, M., Ira, G. & Haber, J. E. RAD51-dependent break-
699 induced replication differs in kinetics and checkpoint responses from RAD51-mediated gene
700 conversion. *Mol. Cell. Biol.* **25**, 933-944 (2005).
- 701 91 Jain, S. *et al.* A recombination execution checkpoint regulates the choice of homologous
702 recombination pathway during DNA double-strand break repair. *Genes Dev.* **23**, 291-303
703 (2009).
- 704 92 Sakofsky, C. J. & Malkova, A. Break induced replication in eukaryotes: mechanisms,
705 functions, and consequences. *Crit. Rev. Biochem. Mol. Biol.* **52**, 395-413 (2017).
- 706 93 Wilson, M. A. *et al.* Pif1 helicase and Poldelta promote recombination-coupled DNA
707 synthesis via bubble migration. *Nature* **502**, 393-396 (2013).
- 708 94 Vijayraghavan, S., Tsai, F. L. & Schwacha, A. A checkpoint-related function of the MCM
709 replicative helicase is required to avert accumulation of RNA:DNA hybrids during S-phase
710 and ensuing DSBs during G2/M. *PLoS Genet.* **12**, e1006277 (2016).
- 711 95 Drissi, R. *et al.* Destabilization of the minichromosome maintenance (MCM) complex
712 modulates the cellular response to DNA double strand breaks. *Cell Cycle* **17**, 2593-2609
713 (2018).
- 714 96 Chen, Y. H. *et al.* Transcription shapes DNA replication initiation and termination in human
715 cells. *Nat. Struct. Mol. Biol.* **26**, 67-77 (2019).

- 716 97 Keskin, H., Meers, C. & Storici, F. Transcript RNA supports precise repair of its own DNA
717 gene. *RNA Biol.* **13**, 157-165 (2016).
- 718 98 Bartek, J., Lukas, C. & Lukas, J. Checking on DNA damage in S phase. *Nat. Rev. Mol. Cell*
719 *Biol.* **5**, 792-804 (2004).
- 720 99 Das, M., Singh, S., Pradhan, S. & Narayan, G. MCM paradox: abundance of eukaryotic
721 replicative helicases and genomic integrity. *Mol. Biol. Int.* **2014**, 574850 (2014).
- 722 100 Sasaki, T. & Gilbert, D. M. The many faces of the origin recognition complex. *Curr. Opin.*
723 *Cell Biol.* **19**, 337-343 (2007).
- 724 101 Borlado, L. R. & Mendez, J. CDC6: from DNA replication to cell cycle checkpoints and
725 oncogenesis. *Carcinogenesis* **29**, 237-243 (2008).
- 726 102 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
727 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 728 103 Loman, N. J., Quick, J. & Simpson, J. T. A complete bacterial genome assembled de novo
729 using only nanopore sequencing data. *Nat. Methods* **12**, 733-735 (2015).
- 730 104 Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial genome
731 assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**, e1005595 (2017).
- 732 105 Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078-
733 2079 (2009).
- 734 106 Pedersen, B. S. & Quinlan, A. R. Mosdepth: quick coverage calculation for genomes and
735 exomes. *Bioinformatics* **34**, 867-868 (2018).
- 736 107 Haas, B. J. *et al.* De novo transcript sequence reconstruction from RNA-seq using the Trinity
737 platform for reference generation and analysis. *Nat. Protoc.* **8**, 1494-1512 (2013).

- 738 108 Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory
739 requirements. *Nat. Methods* **12**, 357-360 (2015).
- 740 109 Lomsadze, A., Burns, P. D. & Borodovsky, M. Integration of mapped RNA-seq reads into
741 automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res.* **42**, e119 (2014).
- 742 110 Stanke, M., Schoffmann, O., Morgenstern, B. & Waack, S. Gene prediction in eukaryotes
743 with a generalized hidden Markov model that uses hints from external sources. *BMC*
744 *Bioinformatics* **7**, 62 (2006).
- 745 111 Haas, B. J. *et al.* Improving the *Arabidopsis* genome annotation using maximal transcript
746 alignment assemblies. *Nucleic Acids Res.* **31**, 5654-5666 (2003).
- 747 112 Waterhouse, R. M. *et al.* BUSCO applications from quality assessments to gene prediction
748 and phylogenomics. *Mol. Biol. Evol.* **35**, 543-548 (2018).
- 749 113 Slater, G. S. & Birney, E. Automated generation of heuristics for biological sequence
750 comparison. *BMC Bioinformatics* **6**, 31 (2005).
- 751 114 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,
752 357-359 (2012).
- 753 115 Augusto Correa Dos Santos, R., Goldman, G. H. & Riano-Pachon, D. M. ploidyNGS:
754 visually exploring ploidy with next generation sequencing data. *Bioinformatics* **33**, 2575-2576
755 (2017).
- 756 116 Yamada, K. D., Tomii, K. & Katoh, K. Application of the MAFFT sequence alignment
757 program to large data-reexamination of the usefulness of chained guide trees. *Bioinformatics*
758 **32**, 3246-3251 (2016).
- 759 117 Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUASt: quality assessment tool for
760 genome assemblies. *Bioinformatics* **29**, 1072-1075 (2013).

- 761 118 Tan-Wong, S. M., Dhir, S. & Proudfoot, N. J. R-Loops promote antisense transcription across
762 the mammalian genome. *Mol. Cell* **76**, 600-616 e606 (2019).
- 763 119 Mazina, O. M. *et al.* Replication protein A binds RNA and promotes R-loop formation. *J.*
764 *Biol. Chem.* **295**, 14203-14213 (2020).
- 765 120 Saldivar, J. C., Cortez, D. & Cimprich, K. A. The essential kinase ATR: ensuring faithful
766 duplication of a challenging genome. *Nat. Rev. Mol. Cell Biol.* **18**, 622-636 (2017).
- 767 121 Longhese, M. P., Plevani, P. & Lucchini, G. Replication factor A is required in vivo for DNA
768 replication, repair, and recombination. *Mol. Cell. Biol.* **14**, 7884-7890 (1994).
- 769 122 Domingo-Prim, J., Bonath, F. & Visa, N. RNA at DNA double-strand breaks: the challenge of
770 dealing with DNA:RNA hybrids. *Bioessays* **42**, e1900225 (2020).

771 **Acknowledgments**

772 The majority of this work was supported by a Foundation grant FRN-142349, awarded to A.J.R. by
773 the Canadian Institutes of Health Research. Archibald Lab contributions to this study were supported
774 by a Discovery Grant from the Natural Sciences and Engineering Research Council of Canada
775 (RGPIN 05871-2014). E.C.T. is supported by a Herchel Smith Postdoctoral Fellowship at the
776 University of Cambridge.

777 **Author contributions**

778 D.E.S-L and A.J.R. conceived the study. J.J-H and M.K. grew cultures, extracted nucleic acids, and
779 carried out *in 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.,
780 A.G.B.S. and A.J.R. analyzed and manually curated the genomic data. E.C.T. and D.E.S-L made the
781 figures. D.E.S-L and A.J.R. led the writing of the manuscript with input from all authors. All
782 documents were edited and approved by all authors.

783 **Competing interests**

784 Authors declare no competing interests.

785 **Additional information**

786 Supplementary Information (also containing legends for Supplementary Table 1 and Supplementary

787 Figures 1 – 10)

788 **Tables**

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

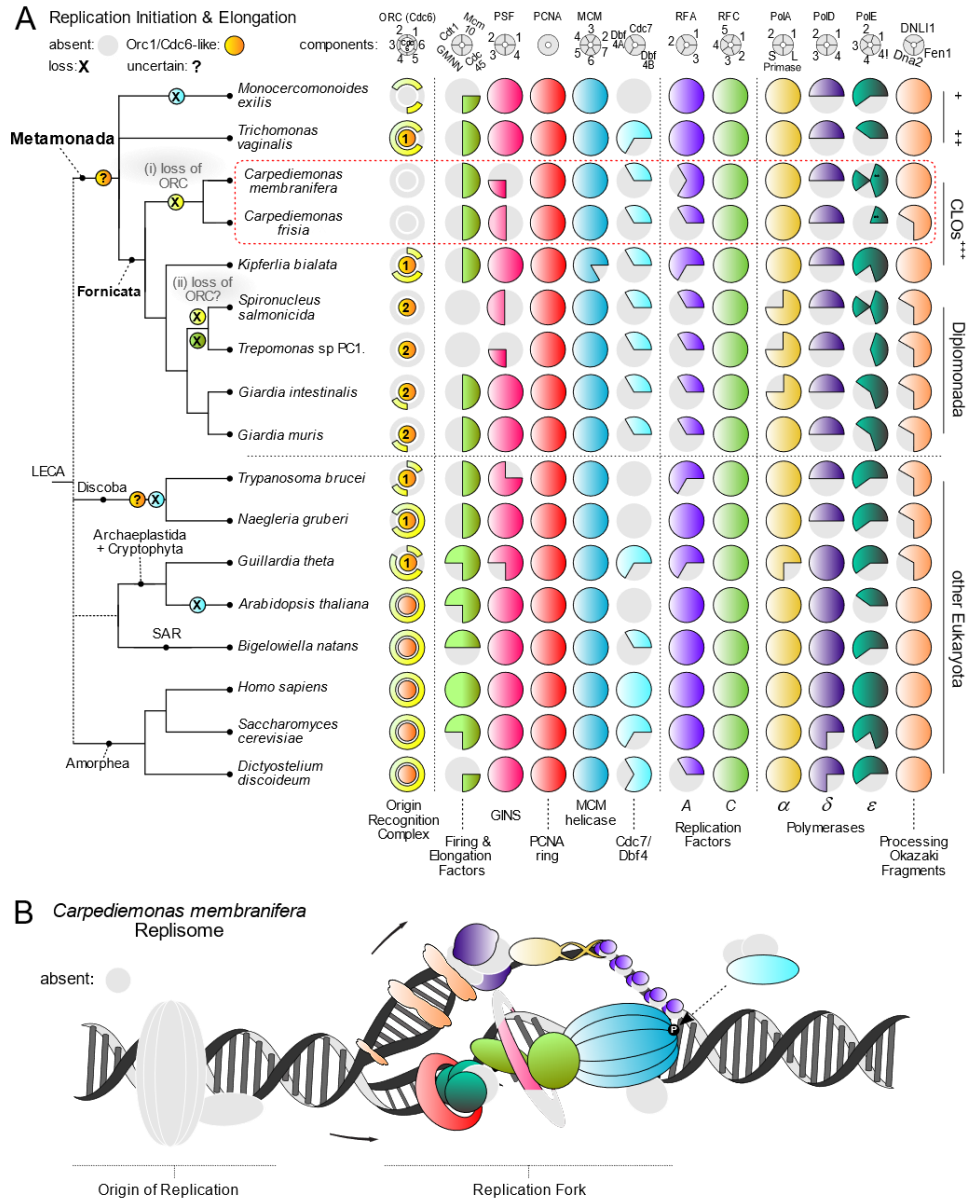
Description	<i>Trichomonas vaginalis</i>	<i>Monocercomonoides exilis</i>	<i>Carpediemonas membranifera</i>	<i>Carpediemonas frisia</i>	<i>Kipferlia bialata</i>	<i>Spirotrunculus salmonicida</i>	<i>Trepomonas PCI*</i>	<i>Giardia intestinalis A 50803</i>	<i>Giardia intestinalis B 50581</i>	<i>Giardia muris</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

790 **All the statistics were recalculated with Quast¹¹⁷ for completion as not all of these were originally reported, and the BUSCO**791 **reference protein set corresponds to a maximum of 245 proteins.**

792 *transcriptome data only

793

794 **Main Figures** (Note: Any reference in main Figure legends can be found in the reference section main text)

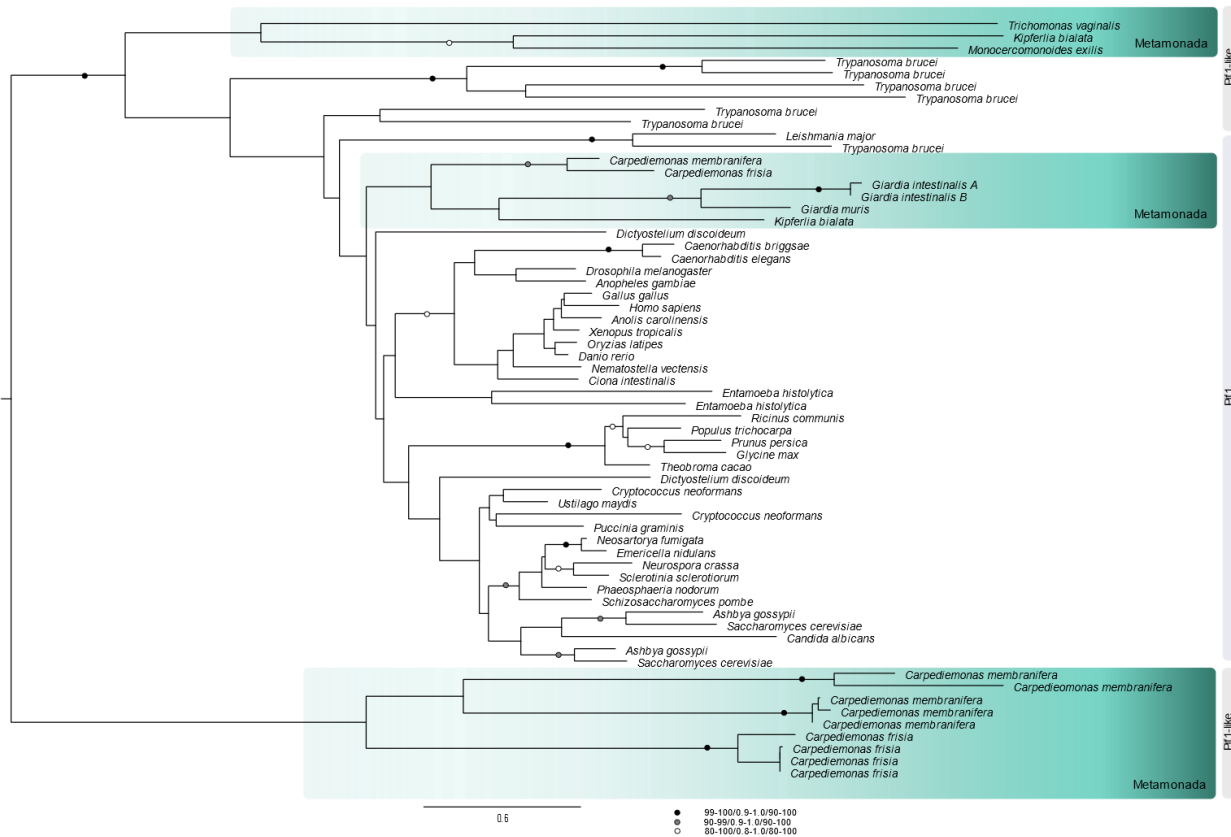


795

796 **Figure 1** The distribution of core molecular systems in the replisome and DNA repair across
 797 eukaryotic diversity. A schematic global eukaryote phylogeny is shown on the left with classification
 798 of the major metamonad lineages indicated at right. **A)** The Replisome. Reduction of the replication

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

813

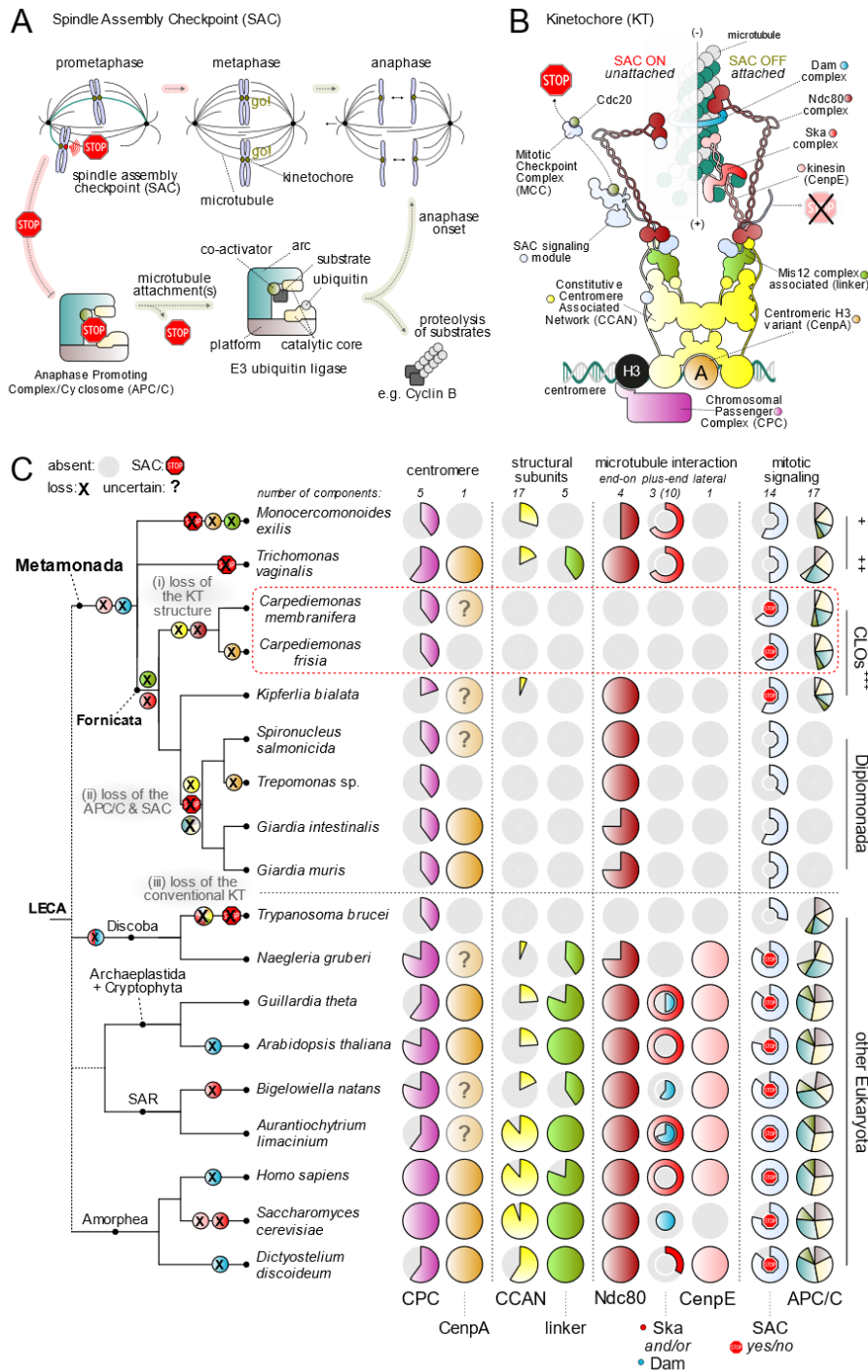


814

815 **Figure 2. Pif1 protein family expansion**

816 Pif1 helicase family tree. Three clades are highlighted: at the top, a Pif1-like clade encompassing some
817 metamonads and at the bottom a *Carpediemonas*-specific Pif1-like clade. The third clade shows the
818 typical Pif1 orthologs encompassing fornicates. The maximum-likelihood tree was inferred under the
819 LG+PMSF(C60)+F+ Γ model using 100 bootstraps based on an alignment length of 265 sites. The tree
820 was midpoint-rooted and the support values on the branches correspond to SH-aLRT/aBayes/standard
821 bootstrap (values below 80/0.8/80 are not shown). The scale bar shows the inferred number of amino
822 acid substitutions per site.

823



824

825 **Figure 3 Radical reduction of ancestral kinetochore network complexity in *Carpediemonas***

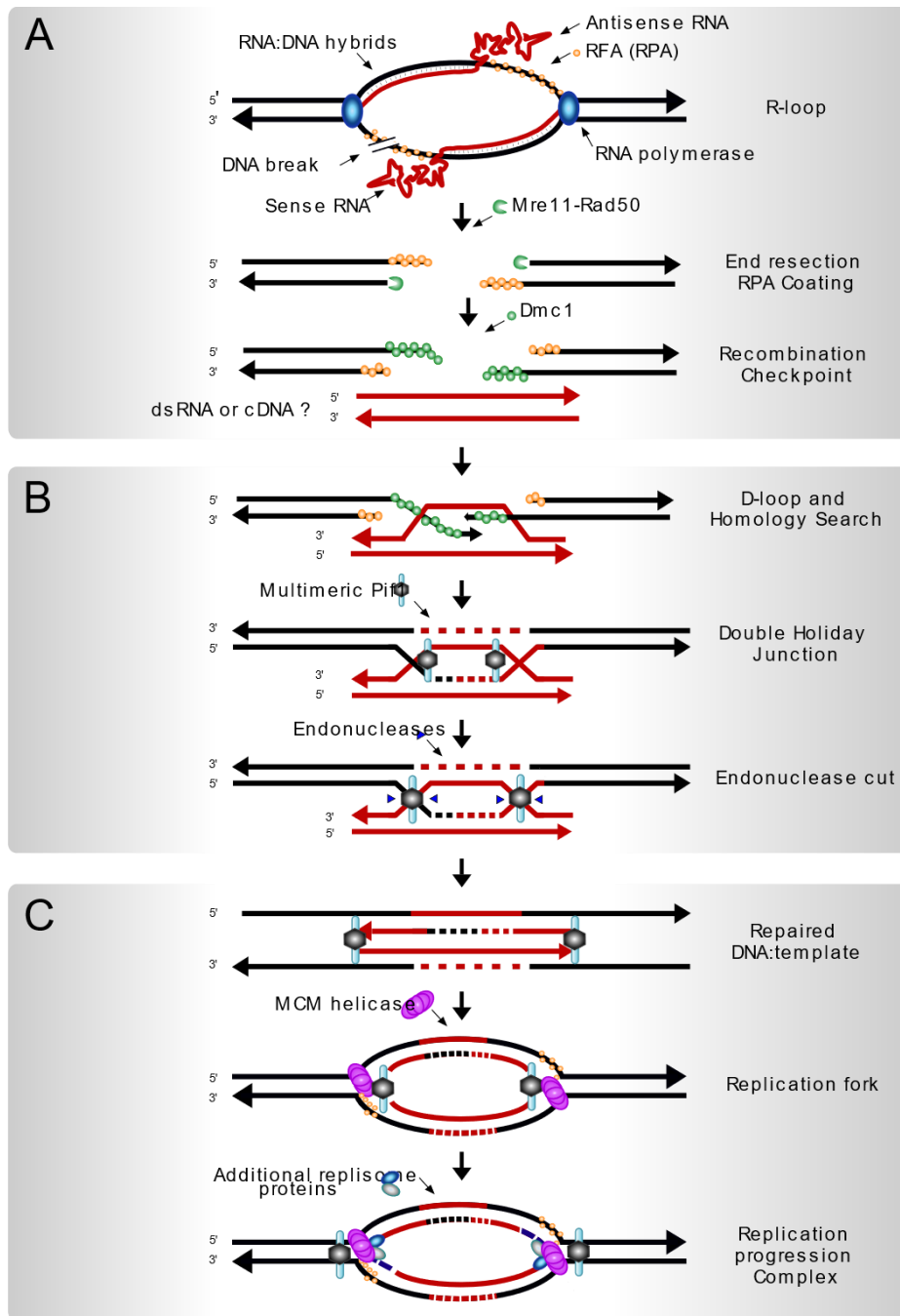
826 **species. A)** Schematic of canonical mitotic cell cycle progression in eukaryotes. During mitosis,

827 duplicated chromosomes each attach to microtubules (MTs) emanating from opposite poles of the

828 spindle apparatus, in order to be segregated into two daughter cells. Kinetochores (KTs) are built upon
829 centromeric DNA to attach microtubules to chromosomes. To prevent precocious chromosome
830 segregation, unattached KT's signal to halt cell cycle progression (STOP), a phenomenon known as the
831 Spindle Assembly Checkpoint (SAC). The SAC entails the inhibition of the Anaphase Promoting
832 Complex/Cyclosome (APC/C), a multi-subunit E3 ubiquitin ligase complex that drives the entry of
833 mitotic cells into anaphase by promoting the proteolysis of its substrates. Once all KT's are correctly
834 attached to spindle MTs and aligned in the middle of the cell (metaphase), the APC/C is released, its
835 substrates are degraded, and chromosome segregation is initiated (anaphase). **B)** Cartoon of the
836 molecular makeup of a single KT unit that was likely present in Last Eukaryotic Common Ancestor
837 (LECA). Colours indicate the various functional complexes and structures. The primary KT structure
838 is provided by the Constitutive Centromere Associated Network (CCAN; yellow), which is built upon
839 centromeric chromatin that contains Centromere protein A (CenP; orange), a centromere-specific
840 Histone H3. During mitosis the CCAN recruits the Mis12 complex (linker; light green), which
841 provides a platform for the recruitment of the SAC signalling (light blue) and microtubule-interacting
842 complexes. The Chromosomal Passenger Complex (CPC; dark purple) localizes at the inner
843 centromere and harbours a kinase (aurora) that regulates microtubule attachments. Unattached KT's
844 catalyse the production of a diffusible cytosolic inhibitor of the APC/C, known as the mitotic
845 checkpoint complex (MCC), which captures the mitotic APC/C co-activator Cdc20. Initial KT-MT
846 encounters are driven by the kinesin Centromere protein E (CenPE; pink), which binds MTs at the
847 lateral sides. The Ndc80 complex (dark red) constitutes the main end-on MT binding activity of KT's.
848 To facilitate the tracking of the plus-end (+) of MT during anaphase, eukaryotes utilize two different
849 complexes: Dam (light purple; likely not present in LECA) and Ska (red). Once KT's are bound by
850 MTs, SAC signalling proteins are removed and the SAC is turned off. **C)** Reconstruction of the
851 evolution of the KT and mitotic signalling in eukaryotes based on protein presence-absence patterns

852 reveals extensive reduction of ancestral KT complexity and loss of the SAC in most metamonad
853 lineages, including the striking loss of the highly conserved core MT-binding activity of the KT
854 (Ndc80) in *Carpodimonas*. On top/bottom of panel C: the number of components per complex and
855 different structural parts of the KT, SAC signalling and the APC/C. Middle: presence/absence matrix
856 of KT, SAC and APC/C complexes; one circle per complex, colours correspond to panel A & B; grey
857 indicates its (partial) loss (for a complete overview see **Supplementary Table 1, Supplementary Fig.**
858 **4**). The red STOP sign indicates the likely presence of a functional SAC response (see for discussion
859 **Supplementary Fig. 6**). On the left: cartoon of a phylogenetic tree of metamonad and other selected
860 eukaryotic species with a projection of the loss and gain events on each branch. Specific loss events of
861 kinetochore and SAC genes in specific lineages are highlighted in colour.

862



863

864 **Figure 4. Hypothesis for Dmc1-dependent DNA replication in *Carpediemonas*.**

865 **A)** R-loop stimulated sense and antisense transcription¹¹⁸ in a highly transcribed locus results in a
866 DNA break, triggering DSB checkpoint control systems to assemble HR complexes and the replication
867 proteins near the lesion^{11,37,119-121}. Once the damage is processed into a DSB, end resection by
868 Mre11/Rad50 creates a 3' overhang and the strands are coated with Replication protein A (RPA),

869 while resected ends are coated with the recombinase Dmc1. **B)** A recombination checkpoint decides
870 the HR sub-pathway to be used⁹¹, then strand invasion of a broken end is initiated into a transcript-
871 RNA or -cDNA template^{39,97,122} followed by the initiation and progression of DNA synthesis with the
872 aid of Pif1 helicase*. This leads to the establishment of a double Holliday Junction (HJ) which can be
873 resolved by endonucleases (*e.g.*, Mus81, Flap, Mlh1/Mlh3). The lack of Chk1 may result in mis-
874 segregation caused by aberrant processing of DNA replication intermediates by Mus81⁴⁸. Given the
875 shortness of the RNA or cDNA template, most possible HJ resolutions, except for the one depicted in
876 the figure, would lead to the loss of chromosome fragments. The HJ resolution shown would allow
877 steps shown in panel C. **C)** A multimeric *Carpodionas* Pif1-like helicase is bound to the repaired
878 DNA as well as to the template. Here, the shortness of the template could resemble a replication
879 intermediate that could prompt the recruitment of MCM, following the addition of the replisome
880 proteins and establishing a fully functional replication fork (Dark blue fragments on 3' ends of the
881 bottom figure represent Okazaki fragments).

882 *Notes: Polymerases α and δ are able to incorporate the correct nucleotides using RNA template⁴⁰;
883 RNase H2 would excise ribonucleosides and replace the correct nucleotide.

884

SUPPLEMENTARY INFORMATION

885 **Table of contents**

886 **A. Supplementary methods**

887 **A1. Culturing and DNA isolation**

888 **A2. Genome size and completeness using BUSCO and a phylogenetic guided approach**

889 **A3. Taxa selected for the comparative genomic analysis.**

890 **A4. Additional strategies used to search for ORC, Cdc6 ad Ndc80 proteins.**

891 **B. Supplementary results**

892 **B1. BUSCO completeness.**

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

894 **B3. DNA replication streamlining in nucleomorphs of chlorarachniophytes and cryptophytes**

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

896 **C. Supplementary discussion**

897 **C1. BUSCO incompleteness**

898 **D. Supplementary references**

899 **E. Supplementary Table legends**

900 **F. Supplementary Figures**

901 **A. Supplementary methods**

902 **A1. Culturing and DNA isolation**

903 Sequencing of *C. membranifera* BICM strain was done with Illumina short paired-end and long

904 MinION read technologies. The Illumina sequencing employed DNA from a monoxenic culture

905 grown in 50 ml Falcon tubes in F/2 media enriched with the bacterium *Shewanella frigidimarina* as

906 food. DNA was isolated from a total of two litres of culture using a salt extraction protocol followed
907 by CsCl gradient centrifugation. RNA was also extracted from these cultures using TRIzol
908 (Invitrogen, USA), following the manufacturer's instructions. For MinION sequencing, *C.*
909 *membranifera* was grown in sterile filtered 50% natural sea water media with 3% LB with either
910 *Shewanella sp* or *Vibrio sp.* isolate JH43 as food. Cell cultures were harvested at peak density by
911 centrifugation at 500×g, 8 min, 20 °C. The cells were resuspended in sterile-filtered spent growth
912 media (SFSGM) and centrifuged again at 500×g, 8 min, 20 °C. The cell pellets were resuspended in
913 1.5 mL SFSGM, layered on top of 9 mL Histopaque®-1077 (Sigma-Aldrich) and centrifuged at
914 2000×g, 20 min, 20 °C. The protists were recovered from the media:Histopaque interface by
915 pipetting, diluted in 10 volumes of SFSGM and centrifuged 500×g, 8 min, 20 °C. High molecular
916 weight DNA was extracted using MagAttract HMW DNA Kit (Qiagen, Cat No. 67563), purified with
917 GenomicTip 20/G (Qiagen, Cat No. 10223) and resuspended in 5 mM Tris-HCl (pH 8.5).

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

919 The BUSCO approach¹ was prone to false negative predictions with our dataset because of the
920 extreme divergence of metamonad homologs. Therefore, the completeness of the BUSCO set was re-
921 assessed with a phylogeny-guided search. For this, we eliminated 31 proteins associated with
922 mitochondria or mitochondrion- related organelles (MROs) as Metamonada have reduced or no
923 MROs², and employed taxa-enriched Hidden Markov Model (HMM) searches to account for
924 divergence between the remaining 272 proteins and the studied taxa. In brief: BLASTp was carried
925 out using the 272 BUSCO proteins as queries for finding their orthologues in a local version of the
926 PANTHER 14.0 database³ to enable the identification of the most likely Panther subfamily HMM
927 and its annotation. Then, each corresponding subfamily HMM was searched for in the predicted
928 proteomes with an e-value cut-off of 1×10^{-1} with HMMER v3.1b2⁴. In cases where these searches did
929 not produce any result, a broader search was run using the HMM of the Panther family with 1×10^{-3} as

930 e-value cut-off. Five best hits for each search were retrieved from each proteome, aligned to the
931 corresponding Panther subfamily or family sequences with MAFFT v7.310⁵ and phylogenetic
932 reconstructions were carried out using IQ-TREE v1.6.5⁶ under the LG+C60+F+ Γ model with
933 ultrafast bootstrapping (1000 replicates). Protein domain architectures were visualized by mapping
934 the respective Pfam accessions onto trees using ETE tools v3.1.1⁷.

935 **A3. Taxa selected for comparative genomic analysis.**

936 Our analyses included the publicly available genomes and predicted proteomes of *Trichomonas*
937 *vaginalis* G3 (Parabasalia, www.trichdb.org), *Monocercomonoides exilis* (Preaxostyla,
938 www.protistologie.cz/hampllab), the free-living fornicates *Carpediemonas frisia*⁸ (i.e., metagenomic
939 bin and predicted proteome), *Carpediemonas membranifera* (reported here) and *Kipferlia bialata*⁹,
940 plus the parasitic diplomonad fornicates: *Giardia intestinalis* Assemblages A and B, *Giardia muris*,
941 *Spironucleus salmonicida*-ATCC50377 (www.giardiadb.org) and *Trepomonas* PC1¹⁰ –the latter was
942 only available as a transcriptome. We also included a set of genomes that are broadly representative
943 of eukaryote diversity, such as *Homo sapiens* GRCh38, *Saccharomyces cerevisiae* S288C,
944 *Arabidopsis thaliana* TAIR10, *Dictyostelium discoideum* AX4, *Trypanosoma brucei* TREU927-rel28
945 (www.uniprot.org), *Naegleria gruberi* NEG-M (www.ncbi.nlm.nih.gov), *Guillardia theta* and
946 *Bigelowiella natans* (www.genome.jgi.doe.gov/portal/).

947 Additional analyzed genomes were those of the microsporidia *Encephalitozoon intestinalis* ATCC
948 50506 (ASM14646v1), *E. cuniculi* GB-M1 (ASM9122v2) and *Trachipleistophora hominis*
949 (ASM31613v1), the yeasts *Hanseniaspora guilliermondii* (ASM491977v1), *Hanseniaspora opuntiae*
950 (ASM174979v1), *Hanseniaspora osmophila* (ASM174704v1), *Hanseniaspora uvarum*
951 (ASM174705v1) and *Hanseniaspora valbyensis* NRRL Y-1626 (GCA_001664025.1), *Tritrichomonas*
952 *foetus* (ASM183968v1), the nucleomorphs of *Hemiselmis andersenii* (ASM1864v1), *Cryptomonas*
953 *paramecium* (ASM19445v1), *Chroomonas mesostigmatica* (ASM28609v1), *Guillardia theta*

954 (ASM297v1), *Lotharella vacuolata* (AB996599–AB996601), *Amorphochlora amoebiformis*
955 (AB996602–AB996604) and *Bigellowiella natans* (ASM245v1), the corals *Galaxea fascicularis*,
956 *Fungia sp.*, *Goniastrea aspera*, *Acropora tenuis* and the coral endosymbionts *Symbiodinium kawagutii*
957 and *Symbiodinium goreau*^{11,12}.

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

959 Strategies included enriched HMMs as mentioned in the main text and HMMs for individual Pfam
960 domains with e-value thresholds of 1×10^{-3} . 1) Metamonad-specific HMMs were built as described for
961 kinetochore proteins – containing the newly found hits plus orthologs from additional publicly
962 available metamonad proteomes or transcriptomes^{2,13}, 2) we applied the eggNOG 4.5 profiles
963 COG1474, COG5575, KOG2538, KOG2228, KOG2543, KOG4557, KOG4762, KOG0995,
964 KOG4438, KOG4657 and 2S26V which encompass 2774, 495, 452, 466, 464, 225, 383, 504, 515,
965 403 and 84 taxa, respectively, and 3) the Pfam v33.1 HMMs: PF09079 (Cdc6_C), PF17872
966 (AAA_lid_10), PF00004 (AAA+), PF13401 (AAA_22), PF13191 (AAA_16), PF01426 (BAH),
967 PF04084 (Orc2), PF07034 (Orc3), PF18137 (ORC_WH_C), PF14629 (Orc4_C), PF14630 (Orc5_C),
968 PF05460 (Orc6), PF03801 (Ndc80_HEC), PF03800 (Nuf2), PF08234 (Spindle_Spc25) and PF08286
969 (Spc24). For Ndc80, Nuf2, Spc24 and Spc25 we also applied the HMMs models published in¹⁴.

970 **B. Supplementary results**

971 **B1. BUSCO completeness.**

972 A subset of 272 BUSCO proteins from the odb9 database was used for a phylogeny-guided search for
973 divergent orthologs. This revealed that: *i*) 27 out of 272 BUSCO (9.9%) proteins are absent in all
974 metamonads, *ii*) only 101 (~41%) of the remaining 245 proteins were shared by all metamonad
975 proteomes, and *iii*) up to 38% are absent in all Fornicata. Metamonad genomes only contained 60% to
976 91% of the BUSCO proteins (**Table 1, Supplementary Table 1**, note that the BUSCO presence-
977 absence patterns of the transcriptomic data from *Trepomonas sp.* PC1 are consistent with those of the

978 remaining diplomonads). These analyses demonstrate that the Metamonada have secondarily lost a
979 relatively large number of highly conserved eukaryotic proteins and, therefore, BUSCO analysis
980 cannot be used on its own to evaluate metamonad genome completeness.

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

982 Metamonad-specific HMM retrieved two candidates for Orc1/Cdc6 proteins from *C. frisia* (*i.e.*,
983 Cfrisia_2222, Cfrisia_2845) and one from *C. membranifera* (*i.e.*, c4603.t1), and one Orc4 candidate
984 from each *Carpediemonas* species (*i.e.*, Cfrisia_2559, ds58_16707). Further inspection of these hits
985 showed that only the AAA+ region shared similarity among all of these proteins, which is expected
986 as ORC and Cdc6 proteins belong to the ATPase superfamily. However, based on full protein
987 identity, full profile composition and domain architecture, the proteins retrieved with the Orc1/Cdc6
988 HMM were confidently annotated as Katanin P60 ATPase-containing subunit A1 (Cfrisia_2222),
989 Replication factor C subunits 1 (c4603.t1) and 5 (Cfrisia_2845), and proteins retrieved with Orc4
990 HMM were members of the Dynein heavy chain (Cfrisia_2559) and AAA-family ATPase families
991 (ds58_16707). The latter is a 744 aa protein that has a C-terminal region with no sequence similarity
992 or amino acid profile frequencies that resembles a Orc4_C Pfam domain from other metamonads or
993 model eukaryotes. All the additional search strategies yielded false positives in *Carpediemonas*
994 species, as these retrieved AAA-family members lacking sequence similarity to orc proteins, showed
995 completely different protein domain architecture than the expected one and were associated with
996 different functional annotation (data not shown). When reconstructing the domain architecture of
997 ORC and Cdc6 proteins in metamonads, we noted that Fornicata Orc1/Cdc6-like proteins are
998 remarkably smaller (*i.e.*, 1.5 to 3 times smaller) than Orc1 and Cdc6 from the model organisms and
999 other protists used later in phylogenetic reconstruction (**Supplementary Figure 1A and B,**
1000 **Supplementary Table 1**). In most cases, the small proteins lack protein domains rendering a
1001 different domain architecture with respect to their homologs in *S. cerevisiae*, *H. sapiens*, *A. thaliana*

1002 and *T. vaginalis* (**Supplementary Figure 1A, Supplementary Table 1**). For example, Orc1 and
1003 Cdc6 paralogs in Fornicata lack BAH, and AAA_lid10 and Cdc6_C domains. Protein alignments
1004 show that the conserved areas of these proteins correspond to AAA+ domain that have relatively
1005 conserved Walker domains A and B (except MONOS_13325 from *M. exilis*), with a few proteins
1006 lacking the arginine finger motif (R-finger) within the Walker B motif (**Supplementary Figure 1B**).
1007 The latter may negatively affect ATPase activity of the R-finger-less proteins. In an attempt to
1008 establish orthology, metamonad Orc1/Cdc6 candidates were used for phylogenetic reconstruction
1009 together with publicly available proteins that have reliable annotations for Orc1 and Cdc6, expected
1010 domain architecture and/or with experimental evidence of their functional activity in the replisome.
1011 Phylogenetic analysis shows that metamonad proteins form separate clades from the *bona fide* Orc1
1012 and Cdc6 sequences (**Supplementary Figure 1C**). One of these separate clades encompasses Orc1-b
1013 from *T. brucei* that has been shown to participate during DNA replication despite lacking the typical
1014 domain architecture¹⁵.

1015 **B3. DNA replication streamlining in nucleomorphs**

1016 The loss of ORC/Cdc6 accompanied by the partial retention of MCM, PCNA, Cdc45, RCF, GINS
1017 and the homologous recombination (HR) recombinase Rad51 was observed in cryptophyte and
1018 chlorarachniophyte nucleomorphs (**Supplementary table 1**). ORC and Cdc6 were found as single
1019 copies (except Orc2) in the nuclear DNA of these two groups; their predicted proteins lack obvious
1020 signal and targeting peptides which would likely prevent them from participating in a nucleus-
1021 coordinated nucleomorph replication. Hence, nucleomorph DNA replication likely occurs by HR
1022 without the assistance of ORC/Cdc6 origin-binding, but this replication might nonetheless be
1023 regulated at the transcriptional level by the nucleus as shown by¹⁶. Many of the remaining nuclear-
1024 encoded proteins involved in replication are present in more than one copy in those taxa, with several
1025 of them containing signal and transit peptides (*e.g.*, H2A, POLD, RCF1 and RFA1)^{16,17}.

1026 **B4. Acquisition of Endonuclease IV, RarA and RNase H1 by lateral gene transfer**

1027 The Endonuclease IV (Apn1 in yeast) and exonuclease III (Exo III) function in the removal of
1028 abasic sites in DNA via the BER pathway. Our analyses show that *C. frisia* and *C. membranifera*
1029 have Exo III and have a prokaryotic version of Endo IV (**Supplementary Fig 8**). Interestingly, none
1030 of the parabasalids and *Giardia* spp. have an Endo IV homolog, either eukaryotic or prokaryotic. *S.*
1031 *salmonicida* and *Trepomonas* sp. PC1, by contrast, appear to encode a typical eukaryotic Endo IV.

1032 The RarA (Replication-Associated Recombination protein A, also named MgsA) protein is
1033 ubiquitous in bacteria and eukaryotes (*e.g.*, homologs Msg1 in yeast and WRNIP1 in mammals) and
1034 acts in the context of collapsed replication forks^{18,19}. *Carpediemonas* possesses a prokaryotic-like
1035 version (**Supplementary Fig 9**) that lacks the ubiquitin-binding Zn finger N-terminal domain typical
1036 of eukaryotic homologs¹⁸. No canonical eukaryotic RarAs were detected in the remaining
1037 metamonads, but it appears that prokaryotic-like RarA proteins in *Giardia*, *S. salmonicida* and
1038 *Trepomonas* sp. PC1 were acquired in an independent event from that of *Carpediemonas*.

1039 Both *Carpediemonas* genomes have a eukaryotic RNase H2, lack eukaryotic RNase H1 but
1040 encode up to two copies of a prokaryotic-like RNase H1 (**Supplementary Fig. 10**) which do not
1041 have the typical eukaryotic HBD domain²⁰. The HBD domain is thought to be responsible for the
1042 higher affinity of this protein for DNA/RNA duplexes rather than for dsRNA^{21,22}. All prokaryotic-
1043 like RNase H1s in metamonads are highly divergent (**Supplementary Fig. 10**) and, in the case of *S.*
1044 *salmonicida* RNaseH1 proteins, these formed very long branches in all of our preliminary trees, that
1045 had to be removed for the final phylogenetic reconstruction. Remarkably, the phylogenetic
1046 reconstruction that includes other metamonad proteins suggests that *Giardia*, *Trepomonas* sp. PC1, *T.*
1047 *foetus* and *T. vaginalis*, also acquired bacterial RNaseH1. *Trepomonas* sp. PC1 and *Giardia*
1048 sequences cluster together but the *T. foetus* and *T. vaginalis* enzymes each emerge amidst different
1049 bacterial branches, suggesting that they have been acquired independently from the *Carpediemonas*

1050 homologs. It should, however, be noted that the support values are overall low, partly due to the fact
1051 that these sequences and their relatives are highly divergent from each other, from *Carpediemonas*
1052 bacterial-like sequences, and from typical eukaryotic RNaseH1.

1053 **C. Supplementary discussion**

1054 **C1. BUSCO incompleteness**

1055 Both eukaryote-wide and protist BUSCO analyses using the BUSCO methods underperformed in our
1056 analyses. Despite using a phylogeny-guided search with the Eukaryota database, a more
1057 comprehensive database than the protist BUSCO database, a remarkably large number of BUSCO
1058 proteins were inconsistently present in Metamonada. This is not surprising, as the clade harbors a very
1059 diverse group of taxa with varied lifestyles and many have undergone genome streamlining^{9,10,23-25},
1060 and the BUSCO databases are expected to be more accurate with greater taxonomic proximity to the
1061 studied genome^{1,26,27}. While it might be tempting to suggest the 101 BUSCO proteins that are shared
1062 by all metamonads be used to evaluate genome completion in the clade, the overwhelming evidence of
1063 differential genome streamlining strongly indicates that databases should be lineage specific (*e.g.*,
1064 *Carpediemonas*, *Giardia*, etc). Hence, our results highlight the need for constructing such databases
1065 including proteins that showcase the sequence diversity of the groups and genes that are truly single
1066 copy in each of these lineages. Regardless, using only standard BUSCO methods to capture genome
1067 completion will still fall short in such assessments as it will fail to evaluate the most difficult-to-
1068 assemble regions of the genome^{27,28}. For that reason, combined approaches such as the ones used here
1069 provide a more comprehensive global overview of genome completeness.

1070 **D. Supplementary references**

1071 1 Waterhouse, R. M. *et al.* BUSCO applications from quality assessments to gene prediction
1072 and phylogenomics. *Mol. Biol. Evol.* **35**, 543-548 (2018).

- 1073 2 Leger, M. M. *et al.* Organelles that illuminate the origins of *Trichomonas* hydrogenosomes
1074 and *Giardia* mitosomes. *Nat Ecol Evol* **1**, 0092 (2017).
- 1075 3 Mi, H. *et al.* PANTHER version 11: expanded annotation data from Gene Ontology and
1076 Reactome pathways, and data analysis tool enhancements. *Nucleic Acids Res.* **45**, D183-D189
1077 (2017).
- 1078 4 Eddy, S. R. Accelerated profile HMM searches. *PLoS Comp. Biol.* **7**, e1002195 (2011).
- 1079 5 Yamada, K. D., Tomii, K. & Katoh, K. Application of the MAFFT sequence alignment
1080 program to large data-reexamination of the usefulness of chained guide trees. *Bioinformatics*
1081 **32**, 3246-3251 (2016).
- 1082 6 Nguyen, L. T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective
1083 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**,
1084 268-274 (2015).
- 1085 7 Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and Visualization of
1086 Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635-1638 (2016).
- 1087 8 Hamann, E. *et al.* Syntrophic linkage between predatory *Carpodemonas* and specific
1088 prokaryotic populations. *ISME J* **11**, 1205-1217 (2017).
- 1089 9 Tanifuji, G. *et al.* The draft genome of *Kipferlia bialata* reveals reductive genome evolution
1090 in fornicate parasites. *PLoS One* **13**, e0194487 (2018).
- 1091 10 Xu, F. *et al.* On the reversibility of parasitism: adaptation to a free-living lifestyle via gene
1092 acquisitions in the diplomonad *Trepomonas sp.* PC1. *BMC Biol.* **14**, 62 (2016).
- 1093 11 Ying, H. *et al.* Comparative genomics reveals the distinct evolutionary trajectories of the
1094 robust and complex coral lineages. *Genome Biol.* **19**, 175-175 (2018).

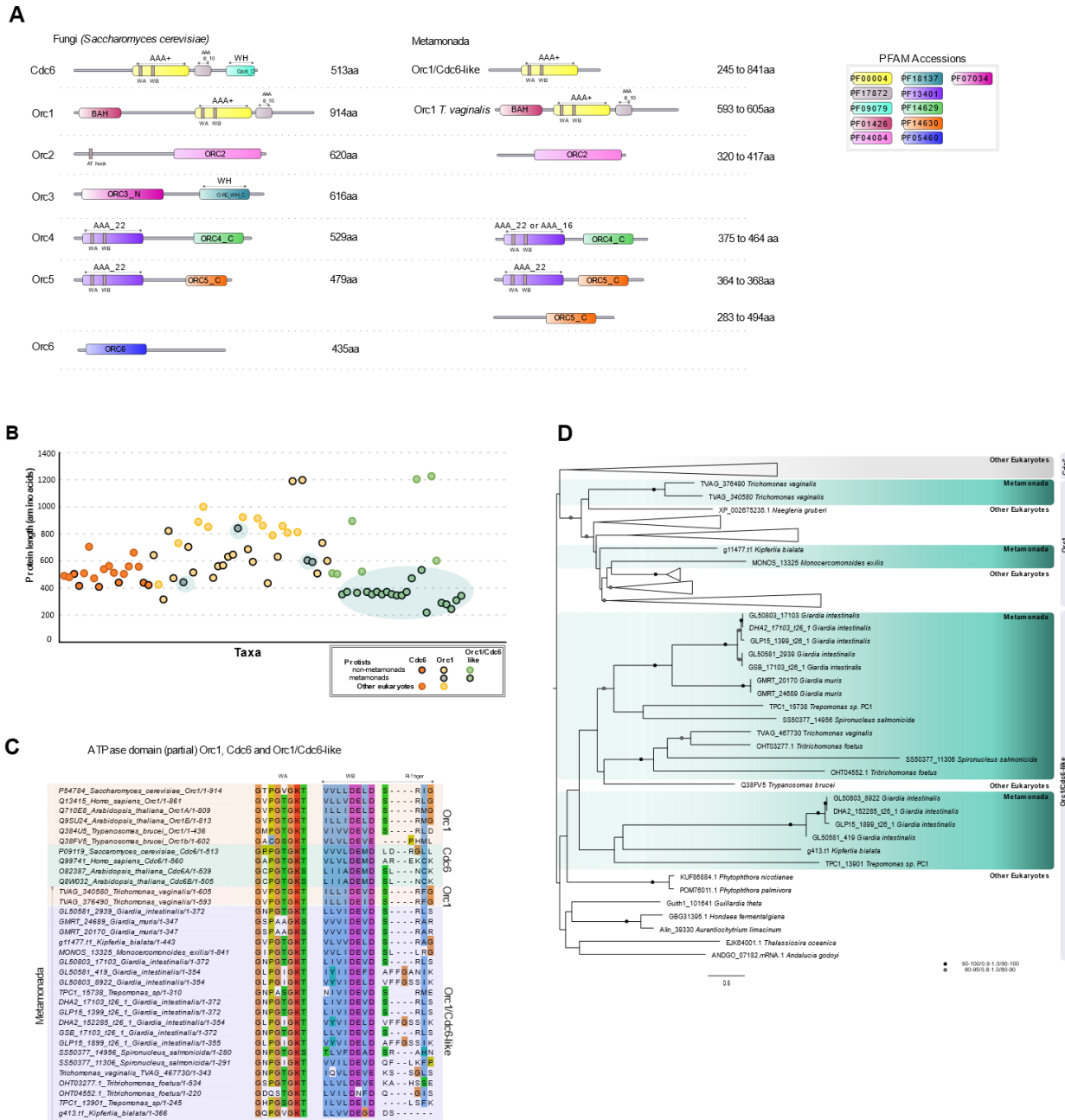
- 1095 12 Woolstra, C. *et al.* The ReFuGe 2020 Consortium—using “omics” approaches to explore the
1096 adaptability and resilience of coral holobionts to environmental change. *Front. Mar. Sci.* **2**
1097 (2015).
- 1098 13 Benchimol, M. *et al.* Draft genome sequence of *Tritrichomonas foetus* Strain K. *Genome*
1099 *Announc.* **5**, e00195-00117 (2017).
- 1100 14 Tromer, E., van Hooff, J., Kops, G. & Snel, B. Mosaic origin of the eukaryotic kinetochore.
1101 *Proc. Natl. Acad. Sci.*, 201821945 (2019).
- 1102 15 Dang, H. Q. & Li, Z. The Cdc45.Mcm2-7.GINS protein complex in trypanosomes regulates
1103 DNA replication and interacts with two Orc1-like proteins in the origin recognition complex.
1104 *J. Biol. Chem.* **286**, 32424-32435 (2011).
- 1105 16 Onuma, R., Mishra, N. & Miyagishima, S. Y. Regulation of chloroplast and nucleomorph
1106 replication by the cell cycle in the cryptophyte *Guillardia theta*. *Sci. Rep.* **7**, 2345 (2017).
- 1107 17 Suzuki, S., Ishida, K. & Hirakawa, Y. Diurnal transcriptional regulation of endosymbiotically
1108 derived genes in the Chlorarachniophyte *Bigeloviella natans*. *Genome Biol. Evol.* **8**, 2672-
1109 2682 (2016).
- 1110 18 Romero, H. *et al.* Single molecule tracking reveals functions for RarA at replication forks but
1111 also independently from replication during DNA repair in *Bacillus subtilis*. *Sci. Rep.* **9**, 1997
1112 (2019).
- 1113 19 Yoshimura, A., Seki, M. & Enomoto, T. The role of WRNIP1 in genome maintenance. *Cell*
1114 *Cycle* **16**, 515-521 (2017).
- 1115 20 Cerritelli, S. *et al.* Failure to produce mitochondrial DNA results in embryonic lethality in
1116 RNaseH1 null mice. *Mol. Cell* **11**, 807-815 (2003).
- 1117 21 Nowotny, M. *et al.* Specific recognition of RNA/DNA hybrid and enhancement of human
1118 RNase H1 activity by HBD. *EMBO J.* **27**, 1172-1181 (2008).

- 1119 22 Cerritelli, S. M. & Crouch, R. J. Ribonuclease H: the enzymes in eukaryotes. *FEBS J.* **276**,
1120 1494-1505 (2009).
- 1121 23 Morrison, H. G. *et al.* Genomic Minimalism in the Early Diverging Intestinal Parasite *Giardia*
1122 *lamblia*. *Science* **317**, 1921-1926 (2007).
- 1123 24 Xu, F. *et al.* The compact genome of *Giardia muris* reveals important steps in the evolution of
1124 intestinal protozoan parasites. *Microb. Genom.* (2020).
- 1125 25 Xu, F. *et al.* The genome of *Spironucleus salmonicida* highlights a fish pathogen adapted to
1126 fluctuating environments. *PLoS Genet.* **10**, e1004053 (2014).
- 1127 26 Saary, P., Mitchell, A. L. & Finn, R. D. Estimating the quality of eukaryotic genomes
1128 recovered from metagenomic analysis with EukCC. *Genome Biol.* **21**, 244 (2020).
- 1129 27 Hanschen, E., Hovde, B. & Starckenburg, S. An evaluation of methodology to determine algal
1130 genome completeness. *Algal Res.* **51**, 102019 (2020).
- 1131 28 Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality,
1132 completeness, and phasing assessment for genome assemblies. *Genome Biol.* **21**, 245 (2020).
- 1133 29 Waterhouse, A. M., Procter, J. B., Martin, D. M., Clamp, M. & Barton, G. J. Jalview Version
1134 2--a multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189-
1135 1191 (2009).
- 1136 30 Musacchio, A. The molecular biology of spindle assembly checkpoint signaling dynamics.
1137 *Curr. Biol.* **25**, R1002-R1018 (2015).
- 1138 31 Alfieri, C., Zhang, S. & Barford, D. Visualizing the complex functions and mechanisms of the
1139 anaphase promoting complex/cyclosome (APC/C). *Open Biol.* **7** (2017).
- 1140 32 Tromer, E., Bade, D., Snel, B. & Kops, G. J. Phylogenomics-guided discovery of a novel
1141 conserved cassette of short linear motifs in BubR1 essential for the spindle checkpoint. *Open*
1142 *Biol.* **6** (2016).

- 1143 33 Vleugel, M. *et al.* Arrayed BUB recruitment modules in the kinetochore scaffold KNL1
1144 promote accurate chromosome segregation. *J. Cell Biol.* **203**, 943-955 (2013).
- 1145 34 Shepperd, L. A. *et al.* Phosphodependent recruitment of Bub1 and Bub3 to Spc7/KNL1 by
1146 Mph1 kinase maintains the spindle checkpoint. *Curr. Biol.* **22**, 891-899 (2012).
- 1147 35 Tromer, E., Snel, B. & Kops, G. Widespread recurrent patterns of rapid repeat evolution in
1148 the kinetochore scaffold KNL1. *Genome Biol. Evol.* **7**, 2383-2393 (2015).
- 1149 36 Moyle, M. W. *et al.* A Bub1-Mad1 interaction targets the Mad1-Mad2 complex to unattached
1150 kinetochores to initiate the spindle checkpoint. *J. Cell Biol.* **204**, 647-657 (2014).
- 1151 37 Ji, Z., Gao, H., Jia, L., Li, B. & Yu, H. A sequential multi-target Mps1 phosphorylation
1152 cascade promotes spindle checkpoint signaling. *Elife* **6** (2017).
- 1153 38 Zhang, G. *et al.* Bub1 positions Mad1 close to KNL1 MELT repeats to promote checkpoint
1154 signalling. *Nat. Commun.* **8**, 15822 (2017).
- 1155 39 Faesen, A. C. *et al.* Basis of catalytic assembly of the mitotic checkpoint complex. *Nature*
1156 **542**, 498-502 (2017).
- 1157 40 Izawa, D. & Pines, J. The mitotic checkpoint complex binds a second CDC20 to inhibit active
1158 APC/C. *Nature* **517**, 631 (2014).
- 1159 41 Di Fiore, B., Wurzenberger, C., Davey, N. E. & Pines, J. The mitotic checkpoint complex
1160 requires an evolutionary conserved cassette to bind and inhibit active APC/C. *Mol. Cell* **64**,
1161 1144-1153 (2016).
- 1162 42 Burton, J. L. & Solomon, M. J. D box and KEN box motifs in budding yeast Hsl1p are
1163 required for APC-mediated degradation and direct binding to Cdc20p and Cdh1p. *Genes Dev.*
1164 **15**, 2381-2395 (2001).

1165

1166 **E. Supplementary figures** (Note: Any reference in Supplementary Figure Legends can be found in
 1167 Supplementary References)



1168

1169

1170

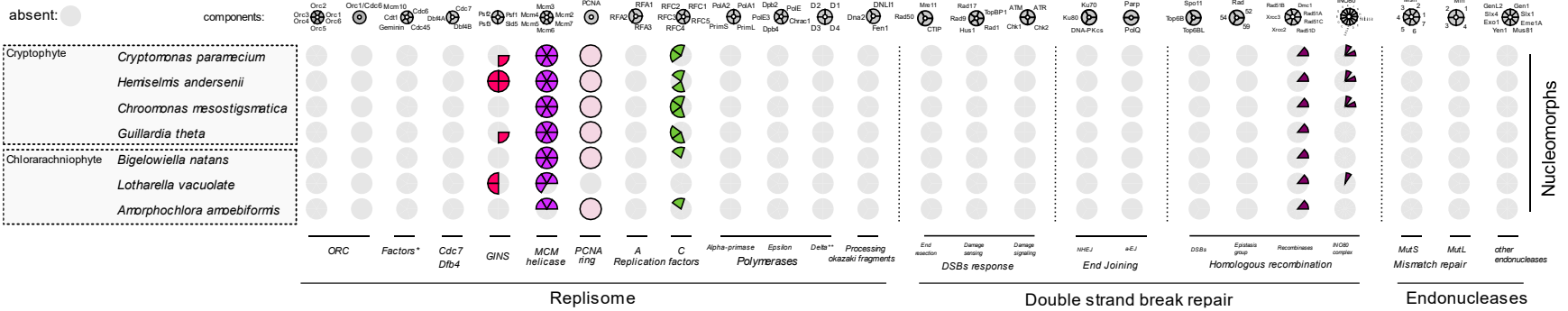
1171

Supplementary Fig 1 Orc1-6 and Cdc6 proteins. **A)** Left: typical domain architecture observed for Orc1-6 and Cdc6 in *Saccharomyces cerevisiae*, Right: representative domain architecture of metamonad proteins drawn to reflect the most common protein size. If no species name is given, then

1172 the depicted domain structure was found in all of the metamonads where present. Numbers on the right
1173 of each depiction correspond to the total protein length or its range in the case of metamonads
1174 (additional information in **Supplementary Table 1**). B) Comparison of Orc1, Cdc6 and Orc1/Cdc6-
1175 like protein lengths across 81 eukaryotes encompassing metamonads and non-metamonads protists
1176 (source information in **Supplementary Table 1**). Metamonad proteins are highlighted with green
1177 shaded bubbles in the background. C) Orc1/Cdc6 partial ATPase domain showing Walker A and
1178 Walker B motifs including R-finger. Reference species at the top. Multiple sequence alignment was
1179 visualized with Jalview²⁹ using the Clustal colouring scheme. D) Phylogenetic reconstruction of Orc1,
1180 Cdc6 and Orc1/Cdc6-like proteins inferred with IQ-TREE⁶ under the LG+ C10+F+ Γ model using
1181 1000 ultrafast bootstraps (bootstrap value ranges for branches are shown with black and grey dots).
1182 The alignment consists of 81 taxa with 367 sites after trimming. Orc1/Cdc6-like proteins do not form a
1183 clade with *bona fide* Orc1 and Cdc6 proteins making it impossible to definitively establish whether or
1184 not they are orthologs.

1185

DNA Replication & DNA Repair



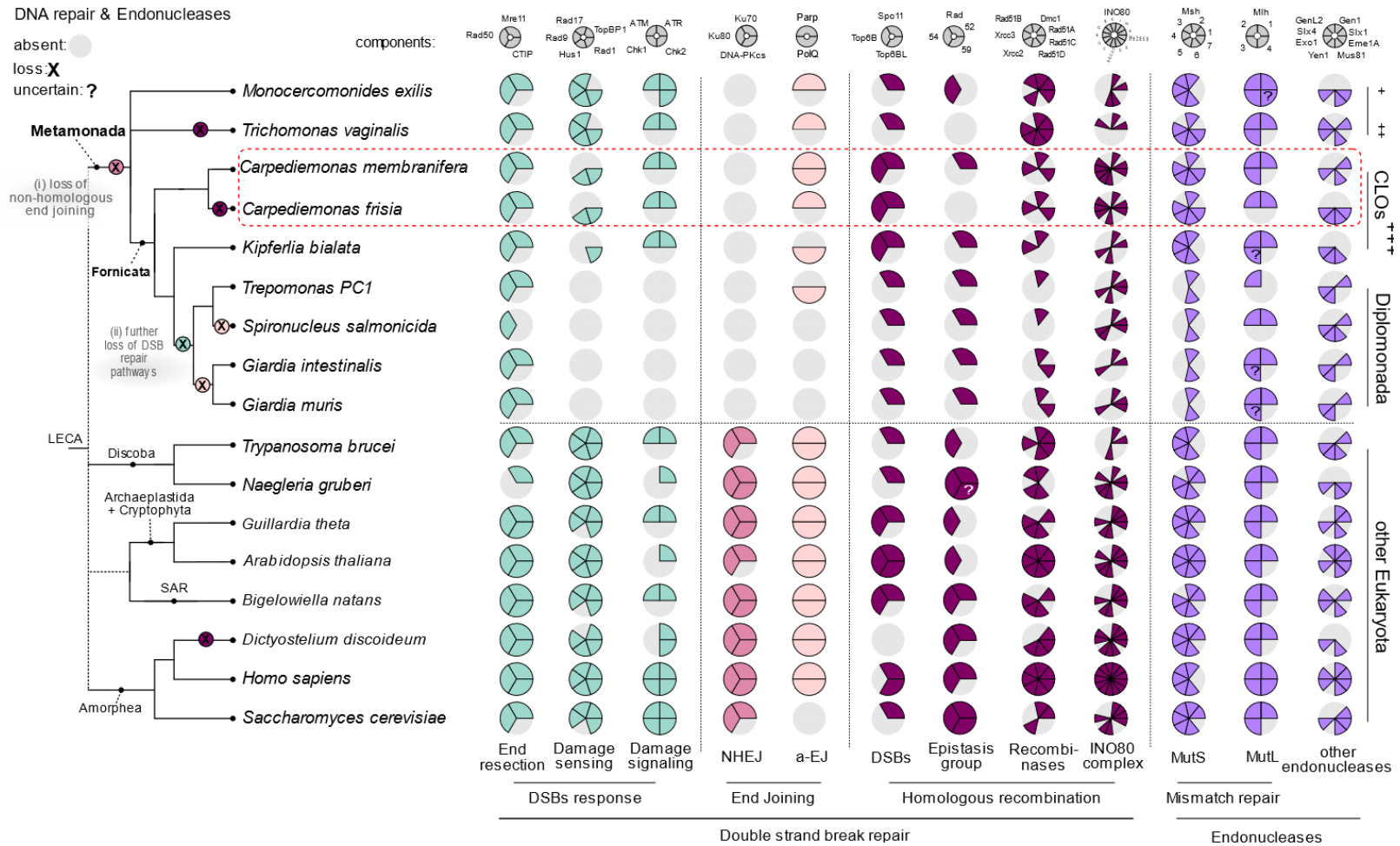
1186

1187

1188 **Supplementary Fig 2** The distribution of core molecular systems of the replisome, double strand break repair and endonucleases in

1189 nucleomorph genomes of cryptophyte and chlorarachniophytes.

1190

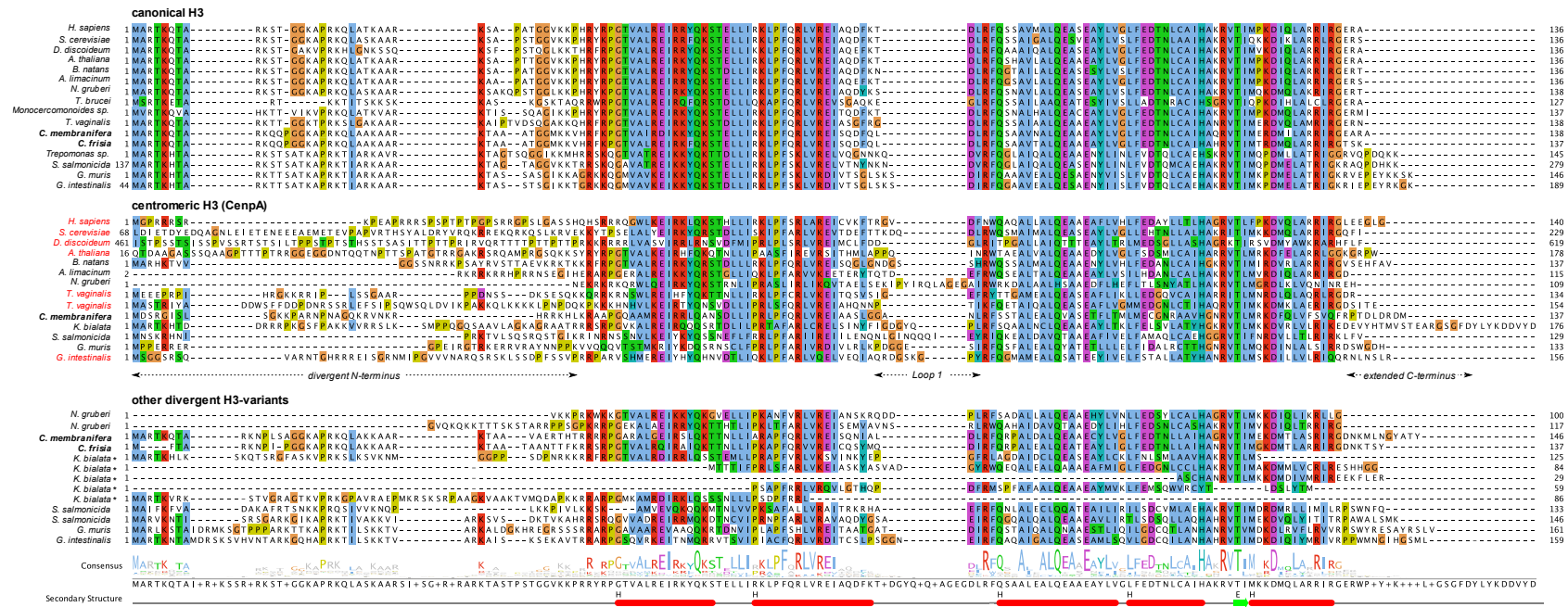


1191

1192 **Supplementary Fig 3** The distribution of core molecular systems of DNA repair across eukaryotic diversity. A schematic global
 1193 eukaryote phylogeny is shown on the left with classification of the major metamonad lineages indicated. Double strand break repair and
 1194 endonuclease sets. ****Carpediemonas*-Like Organisms. ‘?’ is used in cases where correct orthology was difficult to establish, so the
 1195 protein name appears with the suffix ‘-like’ in tables.

1201 respectively. Colour schemes correspond to the kinetochore overview figure on the right and to that used in Figure 1. Right: cartoon of
1202 the components of the kinetochore, SAC signalling, the APC/C and its substrates (Cyclin A/B) in LECA and Carpediemonas species to
1203 indicate the loss of components (light grey shading). Blue lines indicate the presence of proteins that are part of the MCC. Asterisk:
1204 Apc10 has three paralogs in *C. membranifera* and two in *C. frisia*. One is the canonical Apc10, the two others are fused to a BTB-Kelch
1205 protein of which its closest homologs is a likely adapter for the E3 ubiquitin ligase Cullin 3.

1206



1207

1208 **Supplementary Fig 5. *Carpediemonas* harbours three different types of Histone H3 proteins, a centromere-specific variant**

1209 **(CenpA).** Multiple sequence alignment of different Histone H3 variants in eukaryotes and metamonads, including the secondary structure

1210 of canonical H3 in humans (pdb: 6ESF_A). CenpA orthologs are characterized by extended amino and carboxy termini and a large L1

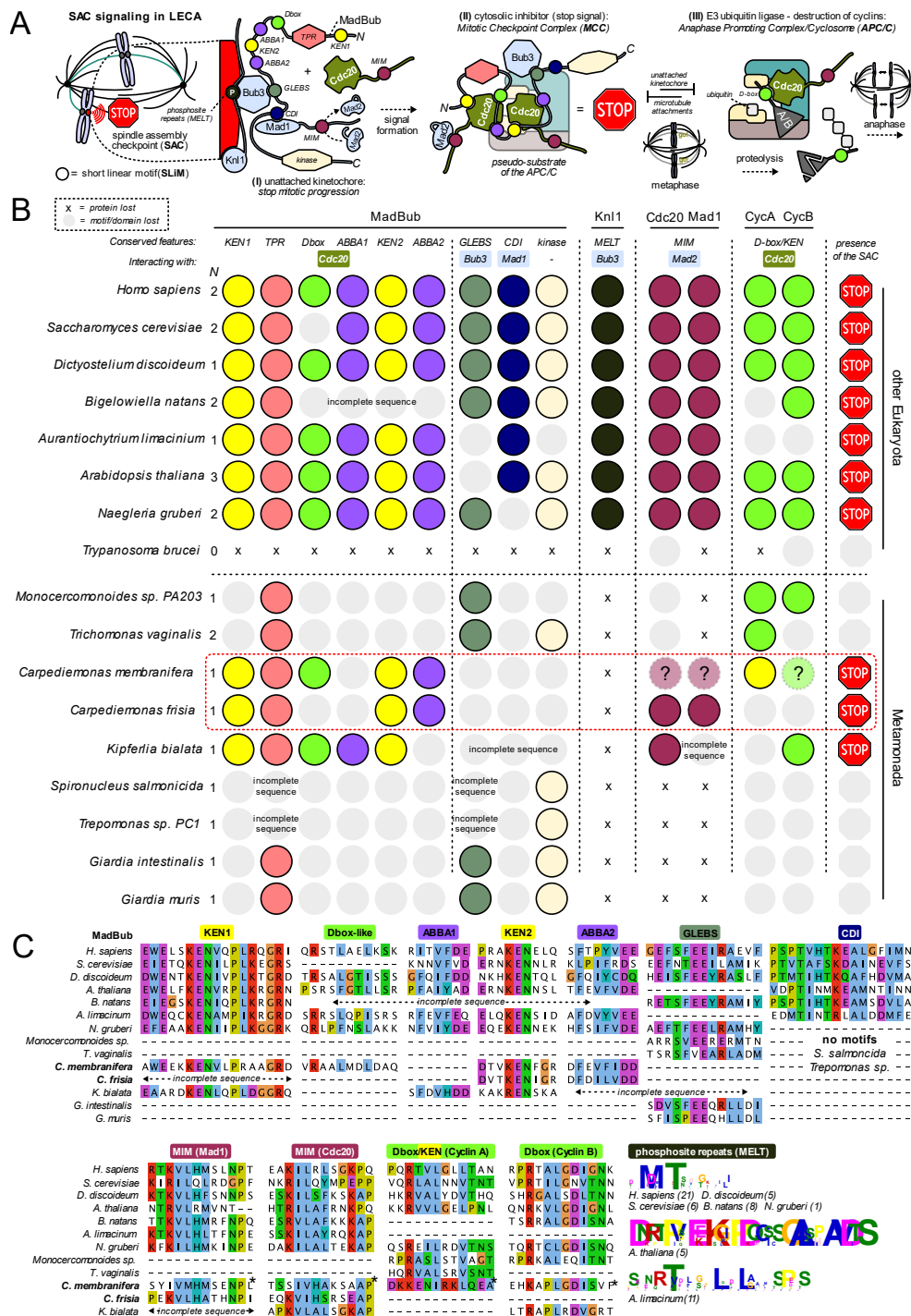
1211 loop. Red names in the CenpA panel indicate for which species centromere/kinetochore localization has been confirmed. In addition to

1212 CenpA and canonical Histone H3-variants, multiple eukaryotes, including *C. membranifera* and *C. frisia*, harbour other divergent H3

1213 variants. Such divergent variants make the annotation of Histone H3 homologs ambiguous (see Asterisks; incomplete sequences).

1214 Multiple sequence alignments were visualized with Jalview²⁹, using the Clustal colour scheme. Asterisks indicate two potential CenpA

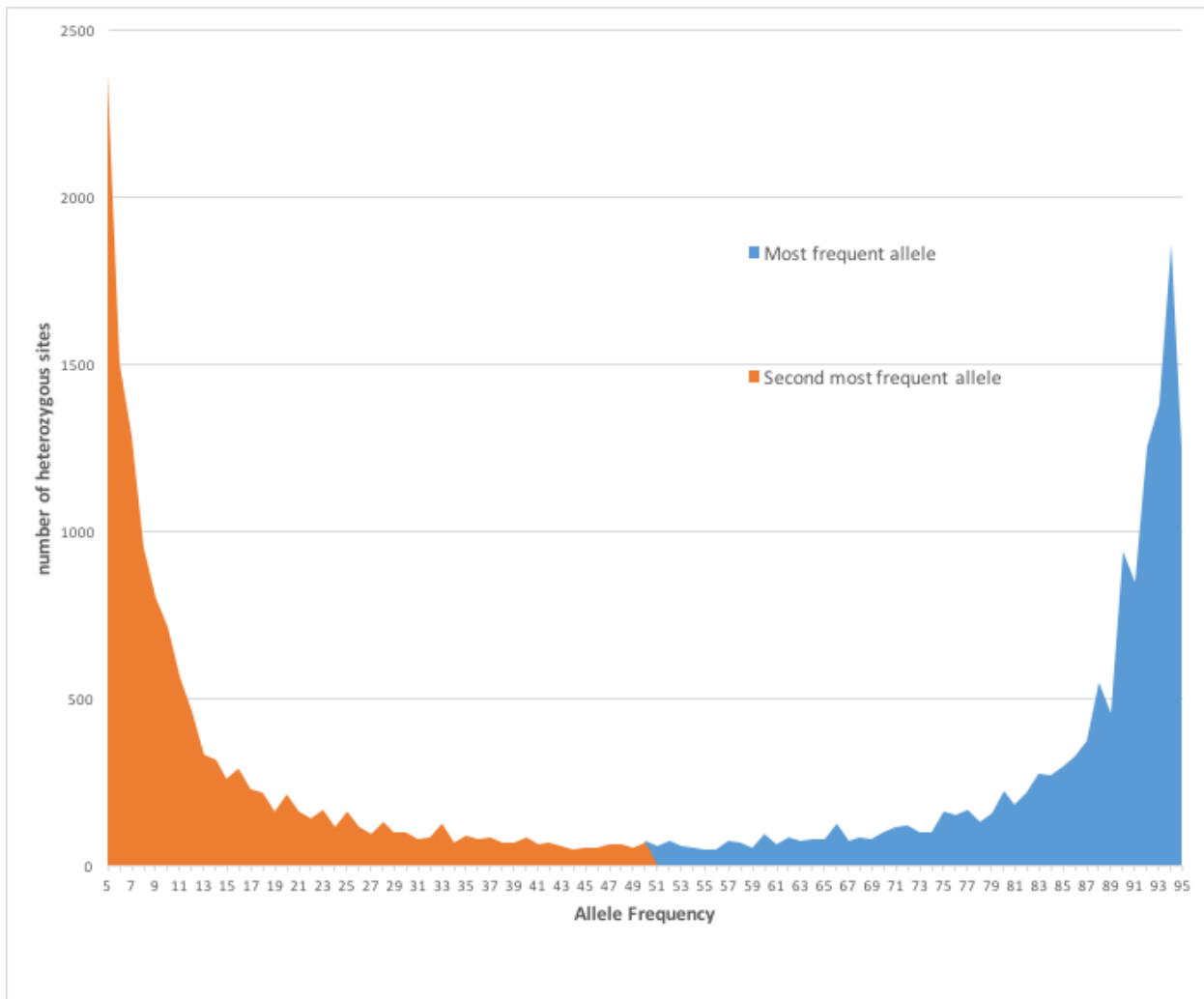
1215 candidates in *T. vaginalis*.



1221 which consist of many short linear motifs (SLiMs) that mediate the interaction of SAC components
1222 and the APC/C (light blue)^{31,32}. MadBub itself is recruited to the kinetochore through interaction with
1223 Bub3 (GLEBS), which on its turn binds repeated phosphomotifs in Knl1³³⁻³⁵. The CDI or CMI motif
1224 aids to recruit Mad1³⁶⁻³⁸, which has a Mad2-interaction Motif (MIM) that mediated the kinetochore-
1225 dependent conversion of open-Mad2 to Mad2 in a closed conformation³⁹. **(II)** Mad2, MadBub, Bub3
1226 and 2x Cdc20 (APC/C co-activator) form the mitotic checkpoint complex (MCC) and block the
1227 APC/C^{32,40,41}. MadBub contains 3 different APC/C degrons (D-box, KEN-box and ABBA motif)³¹ that
1228 direct its interaction with 2x Cdc20s and effectively make the MCC a pseudo substrate of the APC/C.
1229 **(III)** Increasing amounts of kinetochore-microtubule attachments silence the production of the MCC at
1230 kinetochores and the APC/C is released. Cdc20 now presents its substrates Cyclin A and Cyclin B
1231 (some eukaryotes have other substrates as well, but they are not universally conserved) for
1232 ubiquitination and subsequent degradation through recognition of a Dbox motif⁴². Chromosome
1233 segregation will now be initiated (anaphase). **B)** Presence/absence matrix of motifs involved in SAC
1234 signalling in a selection of Eukaryotes and Metamonads, including *C. membranifera* and *C. frisia*.
1235 Colours correspond to the motifs in panel A, light grey indicates motif loss. *N* signifies the number of
1236 MadBub homologs that are present in each species. ‘Incomplete’ points to sequences that were found
1237 to be incomplete due to gaps in the genome assembly. Question marks indicate the uncertainty in the
1238 presence of that particular motif. Although Metamonads have all four MCC components (Mad2, Bub3,
1239 MadBub and Cdc20), most homologs do not contain the motifs to elicit a canonical SAC signalling
1240 and it is therefore likely that they do not have a SAC response. Exceptions are *C membranifera*, *C.*
1241 *frisia* and *Kipferlia bialata*. They retained the N-terminal KEN-boxes and one ABBA motif, which are
1242 involved in the binding of two Cdc20s and a Mad2-interaction motif (MIM) in Mad1 and Cdc20. **C)**
1243 Multiple sequence alignments of the motifs from panel A and B. Coloured motif boxes correspond to

1244 panel A and B. Multiple sequence alignments were visualized with Jalview²⁹, using the Clustal
1245 colouring scheme. Asterisks indicate ambiguous motifs in *Carpodemonas membranifera*.

1246

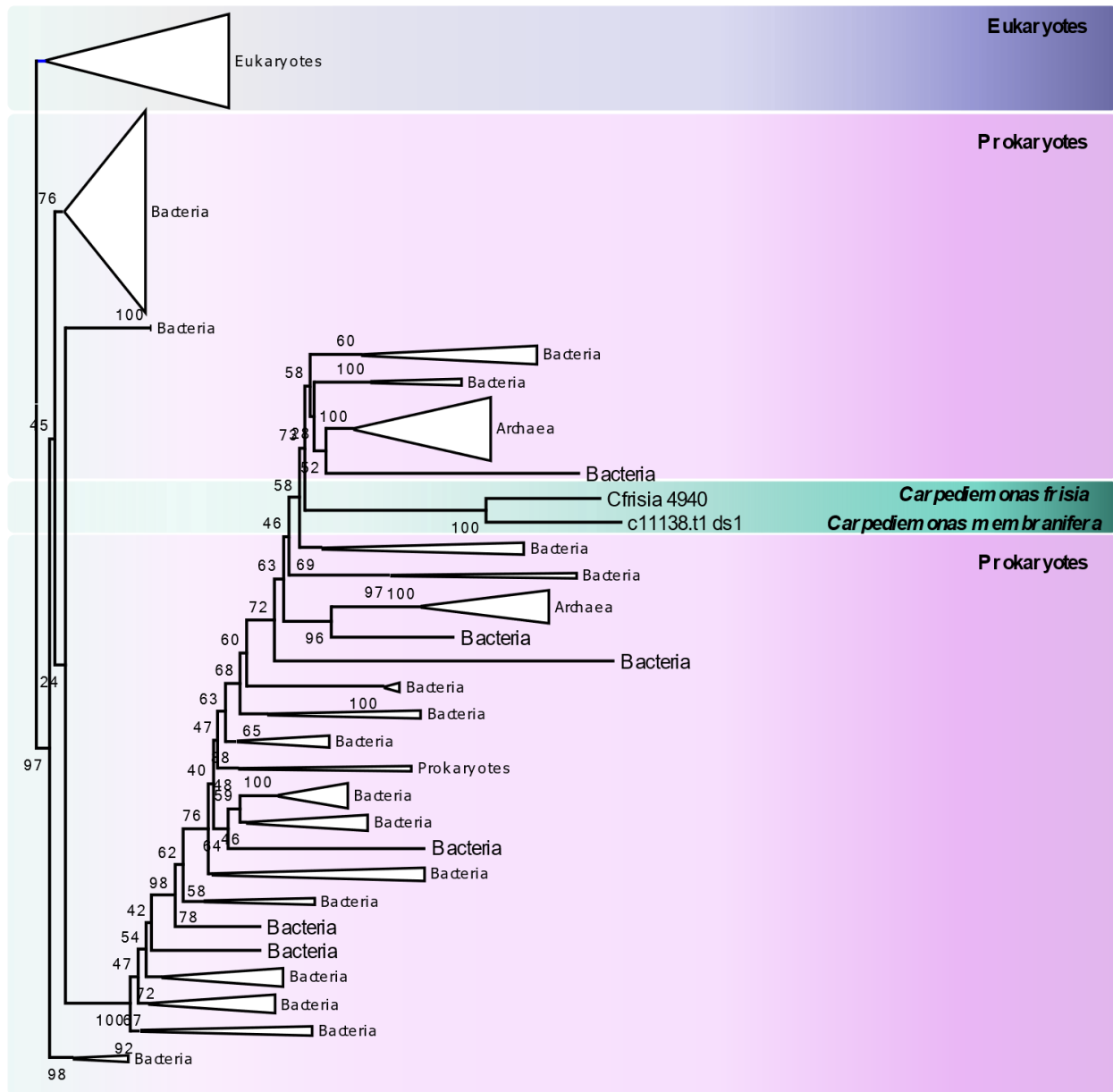


1247

1248 **Supplementary Fig 7 Histogram showing the frequency distribution of single nucleotide variants**

1249 **in the genome of *C. membranifera*.** Diagram showing the typical distribution of a haploid genome.

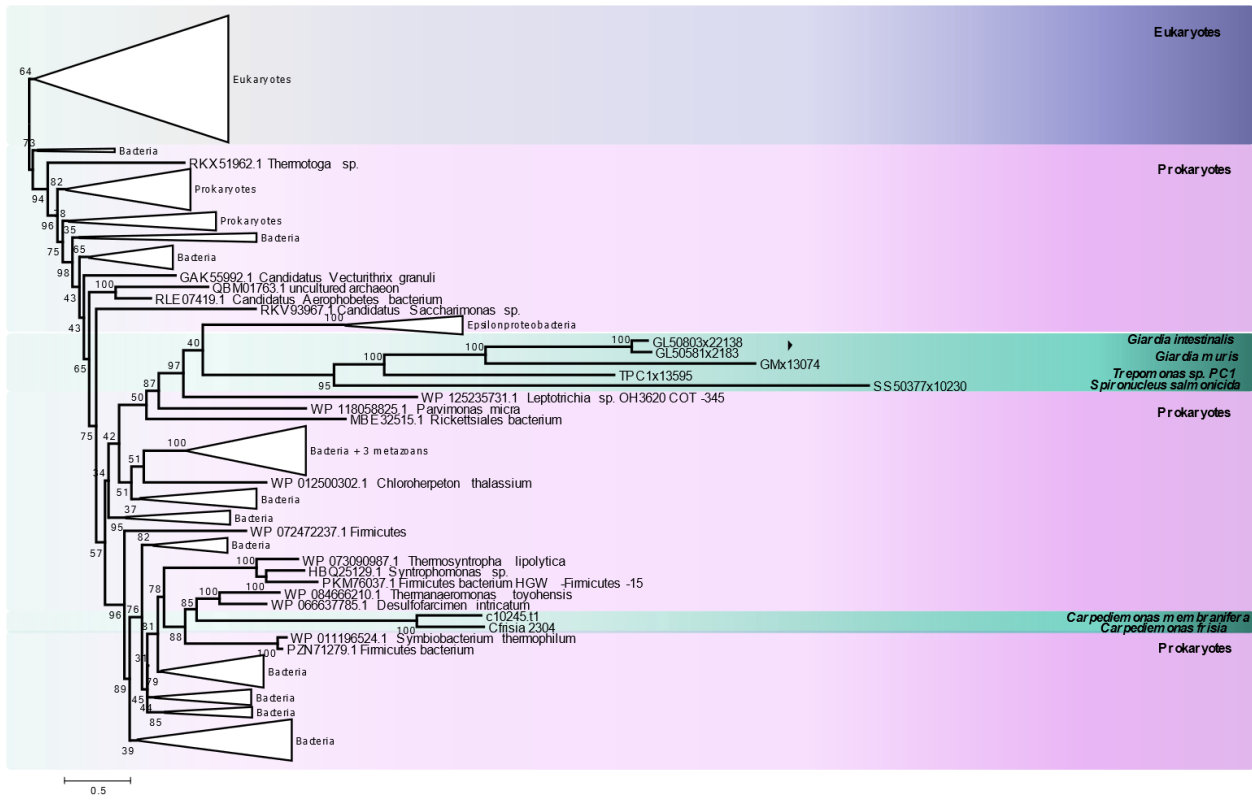
1250



1251

0.5

1252 **Supplementary Fig 8 Maximum likelihood reconstruction of Endo IV.** The unrooted tree contains
1253 eukaryotic and prokaryotic Endo IV sequences, showing *Carpediemonas* sequences emerging within
1254 bacterial proteins. The tree was inferred with IQ-TREE under the LG+I+C20 model with 1000
1255 ultrafast bootstraps; alignment length was 276. Scale bar shows the inferred number of amino acid
1256 substitutions per site.



1267 proteins have been acquired in different events. The tree was inferred with IQ-TREE under the
1268 LG+I+G+C20 model with 1000 ultrafast bootstraps; alignment length was 149. Scale bar shows the
1269 inferred number of amino acid substitutions per site.

1270 **F. Supplementary tables**

1271 Secure download link: http://perun.biochem.dal.ca/downloads/dsalas/Supplementary_Table1.zip

1272 **Supplementary Table 1:**

1273 **Supplementary Table 1A** BUSCO proteins found in Metamonada based on searches for 245 proteins
1274 present in at least one taxon

1275 **Supplementary Table 1B** DNA replication and repair orthologs in 18 diverse eukaryotic genomes

1276 **Supplementary Table 1C** Spindle assembly, kinetochore and APC/C orthologs in 18 diverse
1277 eukaryotic genomes

1278 **Supplementary Table 1D** Additional genomes queried during the searches for ORC, Cdc6 and Ndc80
1279 proteins

1280 **Supplementary Table 1E** Lengths of Orc1-6, Cdc6 and Orc1/Cdc6-like proteins and domain
1281 architecture comparisons between metamonads and other eukaryotes.

1282 **Supplementary Table 1F** Orc1, Cdc6 and Orc1/Cdc6-like proteins. Information used in
1283 Supplementary Figure 1 panels B and D