

25 **Title: Unexpected mitochondrial genome diversity revealed by targeted single-**  
26 **cell genomics of heterotrophic flagellated protists**

27 **Short title:** Single-cell mito-genomics of heterotrophic flagellates

28  
29 Jeremy G. Wideman<sup>a,b,c,d,1,\*</sup>, Adam Monier<sup>a,1</sup>, Raquel Rodríguez-Martínez<sup>a,e,1</sup>, Guy Leonard<sup>a</sup>, Emily Cook<sup>a</sup>,  
30 Camille Poirier<sup>f,g</sup>, Finlay Maguire<sup>a,h</sup>, David Milner<sup>a</sup>, Nicholas A. T. Irwin<sup>i</sup>, Karen Moore<sup>a</sup>, Alyson E. Santoro<sup>j</sup>,  
31 Patrick J. Keeling<sup>j</sup>, Alexandra Z. Worden<sup>f,g</sup>, and Thomas A. Richards<sup>a,\*</sup>

32  
33 <sup>a</sup>Living Systems Institute, University of Exeter, Stocker Road, Exeter EX4 4QD, United Kingdom.

34  
35 <sup>b</sup>Wissenschaftskolleg zu Berlin, Wallotstrasse 19, 14193, Berlin, Germany.

36  
37 <sup>c</sup>Department of Biochemistry & Molecular Biology, Dalhousie University, Halifax, Nova Scotia, B3H 4R2  
38 Canada.

39  
40 <sup>d</sup>Center for Mechanisms of Evolution, Biodesign Institute, School of Life Sciences, Arizona State  
41 University, Tempe, Arizona, 85287 USA.

42  
43 <sup>e</sup>Laboratorio de Complejidad Microbiana y Ecología Funcional, Instituto Antofagasta, Universidad de  
44 Antofagasta, Antofagasta, Chile.

45  
46 <sup>f</sup>Monterey Bay Aquarium Research Institute, 7700 Sandholdt Road, Moss Landing, CA 95039, USA.

47  
48 <sup>g</sup>Ocean EcoSystems Biology Unit, Division of Marine Ecology, GEOMAR Helmholtz Centre for Ocean  
49 Research Kiel, Kiel, Germany.

50  
51 <sup>h</sup>Faculty of Computer Science, Dalhousie University, 1459 Lemarchant Street, Halifax, NS B3H 3P8,  
52 Canada.

53  
54 <sup>i</sup>Department of Botany, University of British Columbia, 3529-6270 University Boulevard, Vancouver, BC  
55 V6T 1Z4, Canada.

56  
57 <sup>j</sup>Department of Ecology, Evolution and Marine Biology, University of California, Santa Barbara, CA 93106,  
58 USA.

59  
60 <sup>1</sup>These authors contributed equally to this work.

61  
62 \*Corresponding authors: Jeremy G. Wideman, Center for Mechanisms of Evolution, Biodesign Institute,  
63 School of Life Sciences, Arizona State University, Tempe, Arizona, 85287 USA. E-mail:  
64 Jeremy.Wideman@asu.edu and Thomas A. Richards, Living Systems Institute, University of Exeter, Stocker  
65 Road, Exeter EX4 4QD, United Kingdom. E-mail: T.A.Richards@exeter.ac.uk

66 **Abstract:**

67

68 Most eukaryotic microbial diversity is uncultivated, under-studied, and lacks nuclear genome data.

69 Mitochondrial genome sampling is more comprehensive, yet many phylogenetically important groups

70 remain unsampled. Using a single-cell sorting approach combining tubulin-specific labelling with

71 photopigment exclusion, we sorted flagellated heterotrophic unicellular eukaryotes from Pacific Ocean

72 samples. We recovered 206 single amplified genomes (SAGs) predominantly from under-represented

73 branches on the tree of life. Seventy SAGs contained unique mitochondrial contigs including 21 complete,

74 or near-complete, mitochondrial genomes from formerly under-sampled phylogenetic branches including

75 telonemids, katablepharids, cercozoans, and marine stramenopiles (MASTs), more than doubling the

76 available sampling of heterotrophic flagellate mitochondrial genomes. Collectively, these data identify a

77 dynamic history of mitochondrial genome evolution including intron gain/loss, extensive patterns of

78 genetic code variation, and complex patterns of gene loss. Surprisingly, we found that stramenopile

79 mitochondrial content is highly plastic, resembling patterns of variation previously only observed in

80 plants.

81 Mitochondria originate from an alphaproteobacteria-like endosymbiont<sup>1</sup>, often contain their own  
82 genomes, and make ATP via oxidative phosphorylation. Most of the 900-1100 different mitochondrial  
83 proteins are nuclear encoded<sup>2</sup>. The progenitor endosymbiont encoded many more genes than extant  
84 mitochondrial genomes (mtDNA), with numerous genes lost or transferred to the nucleus<sup>3</sup>. Mitochondria-  
85 encoded genes vary, but include those essential for mitochondrial transcription, translation, and the  
86 electron transport chain (ETC)<sup>4</sup>. Understanding the dynamics of mitochondrial gene loss and gene transfer  
87 to the nucleus is, however, limited by poor sampling from diverse lineages, especially heterotrophic  
88 flagellates<sup>5</sup>.

89  
90 Microbial eukaryotes, including heterotrophic flagellates, are important constituents of trophic  
91 networks and global biogeochemical cycles<sup>5</sup>, but most remain uncultured. In the absence of cultures,  
92 researchers have used single-cell or targeted metagenome approaches to acquire genomic samples. Three  
93 studies have analyzed partial nuclear or plastid genomes from photosynthetic marine cells<sup>6,7</sup>, and  
94 considerable information exists for cultured phytoplankton<sup>8</sup>. Recent studies have tried to fill gaps, relying  
95 on hand-picking cells of interest<sup>9,10</sup> or fluorescence-activated cell sorting (FACS). Among the latter, a few  
96 have attempted genome sequencing and assembly<sup>11-15</sup> while others have analyzed SSU-rRNA genes from  
97 PCR-amplicons<sup>16,17</sup>. These FACS-based studies have used LysoTracker, to stain acidic compartments such  
98 as food vacuoles<sup>18</sup>, or the permissive DNA stain SYBRGreen in preserved cells<sup>15,17</sup> combined with  
99 chlorophyll exclusion to enrich for putatively phagotrophic cells. Where genome sequencing has been  
100 attempted, insight into the genome sequences of a few eukaryotes has been provided; however, the  
101 highly fragmented incomplete nature of single amplified genomes (SAGs) has restricted their use for  
102 comparative genomics.

103  
104 Here, we hypothesize that mtDNAs will be sampled in SAGs at a tractable frequency, allowing for  
105 comparative analysis. We developed a cell sorting pipeline to select for presence of tubulin, combined  
106 with chlorophyll exclusion, to target heterotrophic flagellates for single cell isolation. Using these samples,  
107 we conducted whole-genome amplification and sequencing, recovering numerous and diverse mtDNAs.  
108 Using these data we investigated mitochondrial gene content evolution, confirming a dynamic pattern of  
109 gene loss, and hitherto unexplored patterns of genetic code variation and intron acquisition.

## 110 111 **Results**

### 112 113 ***Single-cell sampling of marine flagellates***

114 Heterotrophic protists employ diverse lifestyles important for ecosystem function. Since heterotrophs are  
115 poorly sampled it is important that diverse methods are developed to recover diverse forms. Many feed  
116 by phagotrophy, employing acidic vacuoles to digest engulfed prey. Previous studies have used FACS  
117 combined with LysoTracker, which stains acidic vacuoles, to target actively feeding cells for genomic  
118 investigation. However, many heterotrophic flagellates (e.g., obligate osmotrophs<sup>19</sup>), do not phagocytose;  
119 and furthermore, acidic vacuoles can be deployed for a diversity of alternative cellular processes<sup>20,21</sup>.  
120 Therefore, such approaches can yield false positives<sup>22</sup>. To develop alternative ways of recovering  
121 heterotrophic flagellates, while limiting the recovery of false-positives (e.g., prokaryotic cells and detrital  
122 particles), we developed an approach combining flow cytometry with tubulin-specific fluorescence  
123 staining, following the logic that many protists, especially in the marine water column, use their flagella  
124 to find food, hunt prey, and in some cases infect hosts.

125  
126 We sorted small tubulin-positive photopigment-lacking cells from the sub-surface chlorophyll  
127 maximum (SCM, 30 m), isolated DNA and performed multiple displacement amplification (MDA)<sup>23</sup>. V9 PCR  
128 combined with Sanger sequencing identified 206 SAGs containing eukaryotic nuclear Small Subunit (nSSU)  
129 rRNA genes. Our strategy did not include sub-cloning of the SSU-rDNA amplified-template. The Sanger  
130 chromatographs did not show evidence of mixed amplicons suggesting that the V9 sequences recovered  
131 were the predominant rDNA signal from each SAG. These were mapped to a universal eukaryotic  
132 reference tree, revealing a diversity of nSSU sequences that cluster with heterotrophic flagellates (Figs. 1  
133 and 2). Of these, 189 (92%) branched closely to marine heterotrophic flagellates (e.g.<sup>24</sup>) demonstrating  
134 the efficacy of our approach (Figs. 1 and 2 and Supplementary Table 1). Six grouped with taxa containing  
135 photosynthetic/heterotrophic forms (e.g., haptophytes and ochrophytes), and eleven were derived from  
136 non-flagellated fungi previously sampled from marine environments<sup>25</sup> (Fig. 2, panel 1).

137  
138 A recent TARA Oceans-related project presented a broad diversity of heterotrophic nSSUs from  
139 sorted cryopreserved SAGs using SYBR green and chlorophyll exclusion and enriched for different taxa  
140 compared to our analysis<sup>17</sup>. The majority of TARA heterotrophic flagellates recovered were MASTs (362,  
141 71%); whereas, while our protocol recovered some MASTs (12, 6%), we predominantly recovered  
142 cercozoans (53, 26%), Marine ALveolates (MALVs) and dinoflagellates (51, 25%), choanoflagellates (22,  
143 11%), telonemids (13, 6%) and euglenozoans (20, 10%). Although from different geographic sites, the  
144 differences in taxa recovered highlights the importance of developing approaches that target specific  
145 cellular attributes.

146 A rank abundance analysis on nSSU-V9 diversity tag sequences was performed using DNA  
147 isolated from parallel seawater samples from the same depth and 10 m above. We searched these  
148 community profiles for representation of the 206 SAGs and found that our SAGs, were among the rarer  
149 taxa identified (Extended Data 1 and Supplementary Table 2). This was expected as the vast majority of  
150 eukaryotes at the SCM are photosynthetic. These data also show that many abundant heterotrophs were  
151 not recovered in our cell-sampling. This could be a product of bias arising from size exclusion or due to  
152 the limitation of sampling hundreds of cells from a community of millions. We conclude that our sorting  
153 method was effective in targeting heterotrophic flagellates while excluding phototrophs and non-target  
154 cells/particles and can be applied to various environments (e.g. freshwater and potentially, with  
155 modification, in soils).

### 156 157 ***Genome sequencing of single-cell samples***

158  
159 Based on the phylogenetic affiliation of the 206 SAGs, we chose 99 cells from under-sampled lineages for  
160 DNA sequencing (Fig. 2). We generated 204 Gbps with a mean (median) sequencing depth of ~1.61 (1.35)  
161 Gbp/SAG. The resulting reads were assembled generating a mean assembly size of 14.5 Mbp (SD = 13.8)  
162 and N50 of 3.4 kbp (SD = 2.4)/SAG. Full length nSSU rRNAs recovered from these assemblies were used to  
163 confirm the V9 phylogenetic position of the SAGS sampled by BLAST<sup>26</sup> discussed above. In all cases only a  
164 single full-length eukaryotic SSU-type was recovered from each SAG, suggesting that co-sampling of  
165 multiple eukaryotic cells was minimal. After database curation, three of the nSSU-V9-types previously  
166 mapped to a tree were determined to be artefactual: As1 and As2, which mapped as ascomycete fungi  
167 (likely due to long branch attraction), were actually shown to represent a picozoan and a rhizarian,  
168 respectively. Furthermore, the T8-SAG assembly contained a complete telonemid SSU; thus, the V9  
169 amplicon sequence was judged to have mapped erroneously as a dinoflagellate (Supplementary Table 3).

170  
171 To estimate genome completion, we implemented the Core-Eukaryotic-Genes-Mapping-  
172 Approach (CEGMA)<sup>13,14</sup> demonstrating recovery of 0.81-48% of CEGMA genes (mean/median 11.4%/6.5%)  
173 (Tables S3-4), comparable to 2-45% recovery in other studies<sup>11,13,14</sup>. However, this approach to genome  
174 completion estimation is subject to a range of artefacts stemming from: i) sampling wells occupied by  
175 more than one cell, and ii) underestimated completeness due to biases in the CEGMA reference taxa. In  
176 some cases, we know that our assemblies are derived from a mixture of eukaryotic, prokaryotic and viral  
177 signatures (Supplemental Data S1 DOI: 10.6084/m9.figshare.8859014); however, the lack of multiple SSUs  
178 in individual SAGs suggests that eukaryote-eukaryote contamination was minimal.

179 ***Biased recovery of mtDNAs from SAGs***

180  
181 Mitochondrial genome contigs were recovered in 70 of 99 SAGs (Supplementary Table 4). In the 53 SAGs  
182 that demonstrate > 50% predicted mitochondrial completion, the relative coverage of mtDNAs was higher  
183 and more variable (M = 17.0x SD = 17.2) compared to the SAG assemblies (M = 4.9x SD = 2.5)  
184 (Supplementary Table 5) consistent with their derivation from organellar genomes that are often present  
185 in higher copy numbers than nuclear genomes. Interestingly, we observe three distinct groups of SAGs  
186 (Fig. 3), those with 'high' nuclear CEGMA completion, those with high mitochondrial coverage, and those  
187 with both low/intermediate nuclear completion and mitochondrial coverage (Hotelling's T2-test<sup>27</sup>  $p =$   
188  $9.07e-13$ ), but no SAGs with both high nuclear and mtDNA recovery (Fig. 3). The mutually exclusive  
189 recovery of mtDNAs or higher CEGMA score could be due to several factors: mtDNA could be abundant in  
190 some cells, mtDNAs could be preferentially amplified by the SAG methodology (as a product of biased  
191 MDA of circular, or AT-rich genomes), or alternatively, nuclear DNA sampling and amplification may be  
192 retarded relative to mtDNAs due to chromatin wrapping or the complex secondary structures of nuclear  
193 DNA. Regardless of the explanation, our data demonstrate that when mitochondrial DNA is preferentially  
194 recovered from SAG genomes, nuclear gene sampling is limited. The differences between mitochondrial  
195 and nuclear genome coverage, the lack of intervening stop codons in open reading frames, and the  
196 absence of bordering nuclear sequence in mitochondrial contigs, all suggest that we have sequenced bona  
197 fide mitochondrial genomes and not mitochondrial insertions into nuclear genomes.

198  
199 A total of 10 unique, complete circular-mapping mtDNAs were assembled from individual SAGs.  
200 These include: two telonemids (T1 NCBI Accession: MK188946, T12 MN082145), a katablepharid (K4  
201 MK188945), an unknown alveolate (see below As1 MK188935), two MAST3s (S11 MK188941, S18  
202 MK188943), a MAST1 (S17 MK188942), a haptophyte (H2 MK188944), and two choanoflagellates (C14  
203 MK188937, C15 MK188938) (Fig. 4- bold). Two cercozoan mtDNAs were assembled, judged linear, and  
204 likely to be complete based on protein repertoires (R17 MK188936, R32 MN082144, bold in Fig. 4). A  
205 further nine unique near-complete (~75-95% complete, see methods) mtDNAs were identified, but could  
206 not be completed by additional assembly approaches or by PCR. In some cases, these incomplete mtDNAs  
207 provide additional samples validating the provenance of the mitochondrial sampling (Fig. 4). From publicly  
208 available datasets<sup>14,28</sup>, we assembled three additional complete mtDNAs: *Incisomonas marina*, a MAST4a,  
209 and a MAST4e (Fig. 4 asterisks). Additionally, we identified a likely complete MAST4a mtDNA  
210 (EU795181.1) misannotated as a bacterial fosmid in the NCBI database. A near-complete mtDNA from a

211 MAST4d SAG<sup>13</sup> was also assembled (Fig. 5). In total, this effort provided 26 complete or near-complete  
212 unique mtDNAs from poorly sampled eukaryotic branches.

213  
214 To confirm that the mtDNAs belong to the expected taxa, we used our complete and near-  
215 complete mitochondrial assemblies as BLAST queries into the NCBI non-redundant database  
216 (Supplementary Table 6). The choanoflagellate (C14, C15), katablepharid (K4), and haptophyte (H2)  
217 mtDNAs hit related mtDNAs (Supplementary Table 6). Surprisingly, the top hits for As1 were all alveolate  
218 dinoflagellates indicating conflict between the mitochondrial and nuclear signal (see below). All mtDNAs  
219 from stramenopiles (S2, 4, 6, 11, 14, 16, 18), except S17, retrieved other stramenopiles as best hits. Since  
220 the S17 mtDNA did not retrieve sequenced stramenopiles, the Cox1 protein sequence was extracted and  
221 used as a BLAST query retrieving only stramenopile sequences (Supplementary Table 6). Unexpectedly,  
222 the cercozoan mtDNAs and translated Cox1 sequences retrieved stramenopiles and other eukaryotes as  
223 top hits, but not sequenced cercozoans (Supplementary Table 6). We therefore reconstructed a multi-  
224 gene phylogeny using stramenopile and cercozoan mtDNAs (Fig. 6). Our cercozoans bifurcated with  
225 *Bigelowiella* and *Paracercomonas* and not stramenopiles with full support, confirming their likely-identity  
226 as rhizarians. These results lead us to conclude that all assembled mtDNAs with the exception of As1 have  
227 the same taxonomic affiliation as the nSSUs present in each respective sample.

228  
229 The 'As1 SAG' contained a single assembled nSSU sequence 94% identical to the nSSU from  
230 picozoan MS5584-11<sup>11</sup> and a single circular mtDNA. The mtDNA encodes no tRNAs and only five putative  
231 genes including barely identifiable, fragmented, mitochondrial small and large ribosomal RNA genes, *cob*,  
232 *cox1*, and an unidentified open reading frame, but based on the predicted transmembrane architecture  
233 of the protein, is likely a divergent *cox3*<sup>29</sup> (Fig. 4). This repertoire is the same as myzozoan alveolates  
234 differing considerably from the picozoan MS5584-11 mtDNA<sup>11,30</sup>. Consistent with the BLAST results  
235 reported above, phylogenetic reconstruction using Cox1 demonstrated that the As1-derived protein  
236 branches within myzozoans (Extended Data 2). In contrast, the MS5584-11 Cox1 protein did not branch  
237 strongly with any eukaryotic group, as expected for orphan lineages. Given the phylogenetic position of  
238 Cox1 and the myzozoan-like coding content and ribosomal fragmentation, we conclude that the mtDNA  
239 assembled from As1 is derived from a myzozoan not a picozoan, a result potentially arising from sampling  
240 a cryptic cell-cell interaction (predator-prey or host-parasite).

241

242 ***Evolutionary diversity of mitochondrial gene repertoires***

243 The data reported here allowed us to sample a wide diversity of eukaryotic lineages and compare  
244 repertoires of mitochondrial genes (Figs. 4-5). Several gene families thought to be encoded in a small  
245 subset of eukaryotic mtDNAs were shown to be discontinuously distributed across a diversity of lineages  
246 (Fig. 5- red squares). For example, the telonemids possess 40 mitochondrial genes including: *rps1*, *rpl10*,  
247 *rpl18* (Extended Data 3), *rpl31*, *rpl32*, and *tatC* thought to be rare. Whereas the katablepharids contain a  
248 single discontinuously distributed gene (*nad8*). Within the katablepharid mtDNAs we also identified  
249 thirteen additional open reading frames with no similarity to ancestral mitochondrial proteins (Fig. 4).  
250 Some of these genes are similar to LAGLIDADG and GIY-YIG homing endonucleases, but some may  
251 represent undescribed selfish elements or mitochondrial proteins with lineage-specific functions requiring  
252 further investigation. MAST mtDNAs encode additional discontinuously distributed gene families  
253 including *tatA* and *tatC* in MAST1c, MAST3g, *Incisomonas marina*, MAST4, and MAST8 mtDNAs but are  
254 absent in closely related lineages (MAST3i and MAST3e). Prior studies have noted *tatC* in  
255 labyrinthulomycete mtDNAs (KU183024.1 and AF288091.2<sup>31,32</sup>), which is absent in our thraustochytrid-  
256 related cells (S2 MK188939, S4 MK188940). We also identified the RNA component (*rnpB*) of RNase P  
257 encoded by MAST3e and MAST4e, *rps1* by MAST1c, and *rpl31* by MAST1c, MAST4, and MAST8. The  
258 variable nature of stramenopile mtDNA repertoires reveals unexpected dynamics of gene loss and  
259 endosymbiotic transfer within this lineage.

#### 260 261 ***Introns in diverse protist mtDNAs***

262  
263 In addition to the standard bacterial-derived mitochondrial gene repertoire, mtDNAs sporadically contain  
264 Group I and Group II self-splicing introns<sup>33</sup>. Using mfanot ([http://megasun.bch.umontreal.ca/cgi-](http://megasun.bch.umontreal.ca/cgi-bin/mfanot/mfanotInterface.pl)  
265 [bin/mfanot/mfanotInterface.pl](http://megasun.bch.umontreal.ca/cgi-bin/mfanot/mfanotInterface.pl)), we identified introns in cercozoan, choanoflagellate, and  
266 katablepharid mtDNAs (Fig. 4 dark grey lines). Interestingly, the two choanoflagellate mtDNAs recovered  
267 have 97% identity but contain a different number of introns in the *cox1* gene (C14 = 4, C15 = 2, *M.*  
268 *brevicollis* = 3) (Fig. 4). The two encoded homing endonucleases in C15 are similar to two in C14 (89% and  
269 98% amino acid identity), but none are similar to those in *M. brevicollis cox1* (AF538053.1), suggesting a  
270 complex pattern of replacement or rapid intron diversification<sup>34</sup>.

271  
272 Similarly, in the cercozoan mtDNAs, while no introns can be detected in the R1 and R2 mtDNAs,  
273 the cercozoan mtDNAs M9, As2, R32, and R16/17 contain 23, 8, 9, and 8 introns, respectively. Even among  
274 mtDNAs from closely related cercozoans (e.g., R17 and R32 with 97% nSSU rRNA nucleotide identity, Fig.  
275 4), the differences between the number of introns and the different positions of the introns (e.g., R32 has



276 4 large introns in *cox1* whereas R17 has no introns in *cox1*) suggests that most of the introns have been  
277 acquired recently or the genomes sampled have undergone repeated invasion by related introns coupled  
278 with differential loss of intron variants (e.g.,<sup>34</sup>).

279  
280 While the *P. bilix* and some cryptophyte mtDNAs contain no, or very few introns<sup>35-37</sup>, the  
281 katablepharid K4 mtDNA contains seven introns (dark grey in Fig. 4). The published *L. marina* partial  
282 mtDNA sequence contains homing endonuclease-encoding group I introns in the *cob* and *cox1* genes at  
283 identical locations as introns identified in the katablepharid mtDNAs sampled here (49% and 73% amino  
284 acid identity, respectively). Our data confirm that multiple mitochondrial evolutionary lineages undergo a  
285 high turnover of self-splicing introns, whereas other lineages appear free from intron colonisation.

286  
287 ***Stramenopile mitochondrial phylogeny identifies organelle to nucleus transfers, and variations in the***  
288 ***mitochondrial genetic code***

289  
290 Using our MAST mtDNAs and sampling from public databases, we sought to calculate a stramenopile  
291 mtDNA phylogeny. Using sixteen conserved ETC proteins, we reconstructed a 4442-site concatenated  
292 protein phylogeny using members of cercozoans as an outgroup (Fig. 6). The phylogeny recovered  
293 previously established phylogenetic groups including Ochrophyta, Labyrinthulomycota, and  
294 Pseudofungi<sup>28,38</sup>. Similar to other mitochondrial phylogenies<sup>39</sup>, and in contradiction to phylogenies based  
295 on nuclear proteins, we could not recover Ochrophyta-Pseudofungi sisterhood<sup>28,40</sup> suggesting there is  
296 either conflicting phylogenetic signal in mtDNA compared to nuclear markers, or some systematic  
297 phylogenetic artefact is present, discussed previously<sup>39</sup>. Our phylogeny recovered some support for the  
298 placement of MAST clades previously proposed from nSSU rRNA phylogenies<sup>41</sup> and partially corroborated  
299 in a recent multi-gene phylogeny of nuclear encoded genes<sup>28</sup>. These relationships include: an opalozoa  
300 group that includes diverse MAST3s (although *Cafeteria roenbergensis* and MAST12 fall outside this group),  
301 and a sagenistan group containing MAST4s, MAST8, unexpectedly MAST1c, and labyrinthulomycetes (Fig.  
302 6). Given previous evidence of contradictory relationships identified in stramenopile mitochondrial and  
303 nuclear gene phylogenies<sup>39</sup>, the branching order presented here should be treated with caution. As such,  
304 additional sampling of stramenopile lineages is required to understand the conflict observed between  
305 mitochondrial and nuclear phylogenies.

306  
307 Using the mitochondrial phylogeny, we sought to polarise mitochondrial traits onto the  
308 stramenopile tree. So far, recent and frequent functional mitochondria to nuclear gene transfers have  
309 been reported only in Archaeplastida<sup>42</sup> (i.e., green plants). Identification of closely related lineages

310 containing different mitochondrial genes (i.e. MAST4s, MAST1, and MAST8) suggests that genes have  
311 relatively recently been transferred to the nucleus in stramenopile lineages. Indeed, there are numerous  
312 transfers of *atp1* and also partial transfers of *nad11* in multiple stramenopile lineages (Fig. 6 and<sup>40,43</sup>). The  
313 mtDNA of MAST1c lacks *nad7* and MAST12 encodes only the N-terminal half of *nad11*, whereas MAST4s  
314 lack *nad7*, *nad9*, and *nad11*, which are encoded in mtDNAs of most other stramenopiles. We therefore  
315 searched for nuclear-encoded versions of these genes in the MAST1c, MAST12 and MAST4 assemblies<sup>14</sup>.  
316 In MAST1c we identified a short contig encoding the C-terminal region of *nad7* adjacent to sequence with  
317 no similarity to known proteins or genomic DNA. In MAST12 we identified a contig with a C-terminal  
318 domain of *nad11*, which appears to contain spliceosomal introns. Finally, we also identified a contig in a  
319 MAST4 assembly encoding *nad9* adjacent to the U4/U6 small nuclear ribonucleoprotein Prp4 along with  
320 a number of unidentified proteins (Complex I contigs: 10.6084/m9.figshare.7314692). These results  
321 suggest that these essential genes have been relocated to the nucleus in these lineages.

322  
323 Our results demonstrate that stramenopile mtDNA repertoires are extremely diverse compared  
324 to other major lineages like animals and fungi and resemble more closely the dynamic repertoires in the  
325 plant lineage<sup>42</sup>. Interestingly, the patterns of variation identified (Fig. 5) generally correspond to a complex  
326 pattern of losses previously proposed as ‘predictable’ in which ‘non-core’ components of complexes (e.g.,  
327 Complex I components *nad7-11*) are more readily transferred to the nucleus than core (defined as  
328 energetically central) components (e.g., Complex I components *nad1-6*)<sup>44</sup>. These results further support  
329 the hypothesis that the evolutionary diversification of the mitochondrial lineage, deep within the  
330 eukaryotic radiation, was typified by a pattern of early conservation of a wider gene repertoire, followed  
331 by numerous independent gene losses<sup>30</sup>.

332  
333 Lastly, our stramenopile and cercozoan mtDNAs allowed us to trace the evolutionary history of  
334 three genetic code changes. Several mitochondrial code changes have been documented<sup>45</sup>, the most  
335 common being TGA recoded from a stop codon to tryptophan. This simple change has occurred  
336 independently in several lineages including holozoans, fungi, haptophytes, some diatoms, *C. roebergensis*,  
337 cercozoans, picozoan MS584-11, ciliates, and some red and green algae (e.g., see<sup>4</sup>). We show that the  
338 TGA-tryptophan genetic code change observed in *C. roebergensis* is shared with MAST12 and can be  
339 traced to their common ancestor. Likewise, since all cercozoans, including sequences presented here,  
340 encode TGA as tryptophan, it is likely that the code change occurred very early in this lineage. More  
341 strikingly, we identified a genetic code present in our thraustochytrid mtDNAs (two near-complete and  
342 three fragmented, S2, 4 and S1, 3, 15, respectively). In these mtDNAs, TGA and TTA (normally encodes

343 leucine) serve as the only termination codons, and TAG and TAA (normally termination codons) have been  
344 recoded to tyrosine (Extended Data 4). This finding is supported by the identification of a UUA anticodon  
345 tRNA encoded in the SAG mtDNAs (Extended Data 5). It is known that TTA was recoded as a stop codon  
346 in *Thraustochytrium aureum* (AF288091.2)<sup>31</sup> thus we can trace stepwise changes in the mtDNA code in  
347 this lineage (Fig. 6). These data demonstrate a complex pattern of genetic code variation across  
348 stramenopile mitochondria.

349

## 350 **Discussion**

351

352 We demonstrate that mtDNAs are readily recovered from heterotrophic flagellates using tubulin-targeted  
353 single-cell sorting with chlorophyll exclusion followed by whole genome amplification and sequencing.  
354 This represents a method for recovering mtDNAs from diverse uncultured eukaryotes that can be applied,  
355 with minor protocol variations, to investigate a range of environments. Such an approach will allow for  
356 higher resolution studies of protist population structures and for effective sampling of multiple genes with  
357 different rates of sequence variation useful for phylogenetic analyses. The data reported here have  
358 substantially increased publicly available heterotrophic flagellate mtDNAs. NCBI reports 9520 complete  
359 mtDNAs, 8685 from animals, 406 from photosynthetic algae and plants, 334 from fungi, and 50 from  
360 animal/plant parasites (apicomplexans and oomycetes). Of the remaining 44 genomes of heterotrophic  
361 protists, only 17 are heterotrophic flagellates spread across the eukaryote tree of life. Our data more than  
362 doubles this representation, adding complete or near-complete genomes from 5 un- or under-  
363 represented groups (Telonemida (0 + 2), Katablepharida (~30% of 1 genome + 1), heterotrophic flagellated  
364 stramenopiles (2 + 11), Rhizaria (6 + 5) and Choanozoa (1 + 2)). Further investigation in diverse  
365 environments will expand our sampling of heterotrophic protist mtDNAs from across the eukaryotic tree.

366

## 367 **Methods**

368

### 369 ***Sample collection and preparation***

370

371 Seawater was collected in Monterey Bay at 36.6893°N; 122.384°W (Monterey Bay Aquarium  
372 Research Institute timeseries station M2, 56 km from shore) on 7 October 2014 using a Niskin  
373 rosette. Water was collected at 20 m and 30 m (sub-surface chlorophyll maximum as determined  
374 by *in vivo* chlorophyll fluorescence). For general community diversity analyses 500 mL of water  
375 was filtered on to a 0.2 µm pore size Supor filter (Pall cat# 60301) and extracted using a

376 modification of the DNeasy kit (Qiagen) including the addition of a mechanical lysis by bead-  
377 beating<sup>46</sup>. For single-cell sorting, the 30 m water sample was pre-filtered through a 30 µm mesh,  
378 then concentrated by gravity ~70-100 times onto a 0.8 µm filter and stained with Paclitaxel,  
379 Oregon Green® 488 Conjugate (ThermoFisher, 100 ug/mL stock made in DMSO) at 10 µM  
380 (targeting tubulin from cytoskeleton). Cells were washed twice with sterile artificial sea water to  
381 remove unbound dye, then stained with Hoechst 33342 (targeting DNA) at 2 µg/ml. Stained  
382 samples were diluted into sterile artificial sea water in preparation for flow cytometry.

383

### 384 ***Cell sorting of marine heterotrophic flagellates***

385

386 Cells were analyzed and sorted on a BD Influx flow cytometer equipped with a 488 nm and a 355  
387 nm laser and using sterile nuclease-free PBS pH 7.4 as sheath fluid (ThermoFisher cat# AM9625).  
388 A combination of sort windows was applied to select the cells that showed green and blue  
389 fluorescence (captured by a 520/35 nm and a 460/50 nm bandpass filter for Oregon Green  
390 [tubulin] and Hoechst 33342 [Blue-DNA], respectively) as compared to unstained control  
391 samples, and baseline red fluorescence (692/40 nm bandpass filter) indicating the absence of  
392 chlorophyll, allowing us to exclude the majority of photosynthetic cells (See Extended Data 6).  
393 Eighteen SAGs with recovered mitochondrial genomes were obtained following this strategy and  
394 originated from sort 34 and sort 36 (Supplementary Table 2). A majority of SAGs (52) were  
395 recovered from sort 35 where cells were targeted based on Oregon Green fluorescence only and  
396 regardless of Hoechst fluorescence, however sort windows were refined using the forward angle  
397 light scatter (used as a proxy for cell size) to select cells larger than cyanobacterial cells present  
398 in the sample (i.e., *Synechococcus*, recognizable by the orange fluorescence of the phycoerythrin  
399 present in the cells detected in a 572/27 nm bandpass filter).

400

401 Targeted cells were sorted into 96-well plates so that all wells received one individual cell  
402 (Single-Cell sorting mode implemented in the BD FACS 'Sortware' sorter software v1.0.0.650),  
403 except for the outer column of wells which were left empty for negative controls. Duplicate plates  
404 were obtained for sort 34 and 36 and triplicate plates for sort 35. The plates were illuminated by  
405 UV radiation inside the sort chamber for 2 min prior to the sort, covered with foil and placed at -  
406 80°C immediately after the sort. The sort quality and correct drop delay was regularly checked

407 by sorting a known number of polystyrene beads (Polysciences, cat# 17153-10) on a slide and  
408 counting them on an epifluorescence microscope.

409  
410 ***Single-cell genome amplification and sequencing***

411  
412 Samples (sorted cells and negative controls) were lysed for 10 min at 65 °C using alkaline solution  
413 from the Repli-g Single Cell Kit (Qiagen) according to manufacturer's instructions for amplification  
414 of genomic DNA from single cells. After neutralization, samples were amplified using the Repli-g  
415 reagents for a final volume of 50 µl. The MDA reactions were run in a thermal cycler for 8 h at 30  
416 °C. All materials used during MDA procedures were UV-treated in a HL-2000 HybriLinker, UV  
417 Crosslinker (UVP) for 30 to 90 min. Single-cell MDA products were screened using Sanger  
418 sequencing of the V9 region of the nuclear small subunit (nSSU) rRNA gene amplicons derived  
419 from each MDA product. An aliquot of each MDA product was diluted 100-fold in water and 2 µl  
420 of this dilution served as the template for each PCR reaction in 25 µl final volume. PCR  
421 amplification was carried out using the primers: Forward 1389F (5'- TTGTACACACCGCCC-3') and  
422 reverse 1510R (5'-CCTTCYGCAGGTTACCTAC-3') as in<sup>47</sup>. PCR products were run on 1% agarose  
423 gel stained with GelGreen. Bands were cut using a Visi-Blue Plate (in a UVP transilluminator) to  
424 ensure that DNA was not damaged. Amplicons were purified with GeneJet gel extraction kit  
425 (Thermo Scientific), quantified with a Qubit fluorometer using the dsDNA BR kit (Invitrogen) and  
426 sent for Sanger sequencing (Eurofins).

427  
428 For Illumina library preparation an aliquot of each chosen MDA sample (including 6  
429 negative controls) was purified with AMPureXP magnetic beads (Beckmann) following the  
430 manufacturer's instructions, quantified with a Qubit and diluted in 10mM TrisCl (pH 8.0) to a final  
431 volume of 130 µL and a concentration of 7.7 ng/µL. DNA was fragmented using focused acoustic  
432 waves (Covaris E220), concentrated, and libraries made with Nextflex Rapid DNA library  
433 preparation kit and indexes (BIOO Scientific) without PCR amplification. For a subset of samples,  
434 3 µL of each was pooled and concentrated for 450-650 bp size selection using a Blue Pippin 1.5%  
435 agarose cassette with R2 marker. The average size of the recovered libraries was 420 bp (with  
436 295 bp inserts). For a second subset, libraries were prepared similarly but used bead-based size  
437 selection (420-620 bp), rather than Blue Pippin, quantified by qPCR and equimolar pooled at 2

439 nM. Library pools were denatured, diluted and 250 paired-end sequenced across two lanes on a  
440 HiSeq 2500 using Rapid Run SBS v2 reagents (Illumina). Nine repeated samples which were  
441 sequenced more deeply on an additional HiSeq 2500 lane in order to obtain better coverage of  
442 these genomes (Supplementary Table 4).

443  
444 For environmental census of nSSU amplicon libraries, 10 ng environmental DNA was  
445 amplified in a two-step protocol following the Illumina amplicon library preparation strategy.  
446 Sequencing primers comprised Illumina Nextera pad sequence, a 12 base unique molecular  
447 identifier, a spacer sequence, and 1389F or 1510R sequences described above. Two cycles of PCR  
448 were performed using these primers in four 25  $\mu$ L PCR reactions with 2.5 ng DNA in each.  
449 Reactions were pooled and purified using AmpureXP beads before adding NexteraXT indexes in  
450 a second PCR reaction (21 cycles) to complete the library preparations. Triplicate samples were  
451 prepared, pooled in equimolar amounts, and quantified by qPCR before 125 bp PE Illumina  
452 sequencing.

453  
454 ***Single-cell genomic assembly***

455  
456 All SAG sample libraries were assembled using the automatic workflow available at  
457 <https://zenodo.org/record/192677> (DOI: 10.5281/zenodo.192677) or  
458 [https://github.com/guyleonard/single\\_cell\\_workflow](https://github.com/guyleonard/single_cell_workflow). All Illumina read library samples were  
459 uploaded to an Amazon EC2 instance (m4.10xlarge) of Ubuntu Linux. The 150 bp PE read libraries  
460 were then overlapped using the program PEAR<sup>48</sup> in order to create "long" reads, the resulting  
461 long reads and the pairs that did not overlap were subsequently quality and adaptor trimmed  
462 using the program Trim Galore!  
463 ([https://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)). The resulting libraries  
464 were then assembled with SPAdes 3.7.1<sup>49</sup> using the single-cell mode, the careful option and with  
465 a combination of k-mers (21, 33, 55). Quality assessment of the resulting scaffolds was computed  
466 with the analysis software QUASt<sup>50</sup> and completeness profiles were made using CEGMA<sup>51</sup>. A set  
467 of "blobtools" charts are also made with a combination of scaffolds, read mapping and  
468 megaBLAST hits to the NCBI 'nt' database<sup>52</sup>. Additional analyses including KRONA taxonomy  
469 charts and QUALIMAP reports of mapping were computed from these data<sup>53</sup>.

470 ***Universal nSSU tree, V9 mapping, taxon identification from rDNA assemblies***

471  
472 For the SAG taxonomic classification, nSSU V9 sequences (primers 1389F-1510R<sup>54</sup>, a single  
473 sequence from each sample) corresponding to each of our 206 SAG were phylogenetically  
474 mapped onto an nSSU reference phylogenetic tree (see Supplementary Table 1) reconstructed  
475 from a processed version of the nSSU Protist Ribosomal Reference (PR2) database v4.4<sup>55</sup>; built  
476 from GenBank release 203). We first processed the PR2 database by removing short sequences  
477 (< 400 bp) and/or sequences not spanning the V9 region. In addition, sequences from metazoan  
478 organisms (based on PR2/GenBank taxonomic data) were also discarded. To remove sequence  
479 redundancy, the PR2 database was then clustered using CD-HIT v4.6<sup>56</sup> at 90% sequence identity  
480 for sequences classified as Opisthokonta (resulting in 2,694 clusters) and at 98% for non-  
481 Opisthokonta sequences (18,245 clusters). This final processed PR2 database, used for  
482 subsequent phylogenetic analysis, was composed of 20,939 nSSU clusters, representing a total  
483 of 132,235 nSSU sequences.

484  
485 Cluster representatives, along with SAG V9 sequences, were then aligned with PyNAST  
486 v1.2<sup>57</sup> using the nSSU seed alignment from Silva release v123<sup>58</sup> as a template alignment. The  
487 resulting alignment was then edited and trimmed using Trimal v1.4<sup>59</sup> to remove sites with gaps  
488 in more than 25% of the sequences, but conserving at least half of the original alignment (i.e., -  
489 gt 0.25 -cons 50 parameters); the final alignment was composed of 1,750 sites. Aligned SAG V9  
490 sequences were removed from the alignment and the PR2-based maximum-likelihood (ML) tree  
491 was reconstructed using RAxML v8.2 (multithreaded version; PTHREADS-SSE3)<sup>60</sup> under the GTR  
492 model with CAT approximation. SAG V9 sequences were mapped onto the PR2 reference ML tree  
493 using the RAxML Evolutionary Placement Algorithm (EPA<sup>61</sup>) under GTR-CAT. To evaluate local  
494 node supports, a Shimodaira-Hasegawa (SH)-like test<sup>62</sup> was run using FastTree v2.1 (double  
495 precision build<sup>63</sup> in 'accurate' mode (-mlacc 2 -slowlni parameters) and under GTR-CAT.  
496 Subsequently to the phylogenetic mapping, and for tree display purposes, taxa with long  
497 branches were pruned from the phylogenetic tree; specifically, branches were pruned if the  
498 length of the inner node's parent branch was longer than 0.2 substitutions per site or if the  
499 terminal branch (i.e., linking a leaf to a node) was longer than 3 substitutions per site. These long

500 branches were identified and removed using the Newick utilities package<sup>64</sup>; note that no SAG V9  
501 sequences were mapped onto these long branches. The figures corresponding to the full, circular  
502 PR2 phylogenetic tree with SAG V9 mapping (Fig. 1) and clade-specific trees (Fig. 2) were  
503 rendered using the R package ggtree<sup>65</sup>.

504  
505 Contigs from assemblies containing rRNA gene sequences were extracted and used as  
506 queries in BLAST searches to confirm V9 mapping results (Supplementary Table 2). Out of 99  
507 sequenced SAGs, 86 V9 placements corresponded closely with the respective assembled nSSU  
508 BLAST hits, whereas 3 did not corroborate the V9 mapping results, including both sequences that  
509 mapped to ascomycetes, and one sequence that mapped to dinophyte. In these cases the nSSU  
510 assembly data clearly indicate that the V9 regions were misplaced during mapping, the first two  
511 due to long-branch attraction, and the third due to poor V9 sequence quality. The negative  
512 controls contained predominantly very small fragments of contigs most similar to bacterial SSU  
513 sequences possibly due to contamination. However, two of the six total negative control samples  
514 subjected to sequencing contained low-coverage contigs most similar to the nSSU sequence of  
515 *Cryothecomonas aestivalis* (97-99% identity). Since these controls were taken from different 96-  
516 well plates than our samples related to *C. aestivalis*, it is extremely unlikely that these control  
517 wells were contaminated either biologically or during library preparation. Instead, it is much  
518 more likely that the large signal from the 25 SAG samples that contained contigs with extremely  
519 high coverage (sometimes in the thousands) most similar (97-99% identity) to nSSU sequences  
520 of *C. aestivalis* interfered with the detector during the sequencing run. The abundance and over-  
521 representation of these sequences in our SAG samples is a plausible source of the apparent  
522 technical contamination (i.e., instrument-derived) of these two negative controls, as well as some  
523 other samples (see Supplementary Table 1).

524

### 525 ***Monterey Bay V9 tag sequencing diversity census of whole seawater samples***

526

527 Primers and other technical sequences were trimmed from demultiplexed paired end reads using  
528 cutadapt v1.14<sup>66</sup>. To identify artefactual sequences, reads were searched against a V9 reference  
529 database (a V9-trimmed version of PR2, clustered at 80% sequence identity using CD-HIT) using  
530 BLASTn<sup>67</sup>; reads with no significant hit (E-value < 1e-5) against the reference database were



531 discarded. Reads were then processed using DADA2 v1.4<sup>68</sup>. Based on quality profiles, forward  
532 reads were truncated at 150 bp, reverse reads at 100 bp and reads with more than two expected  
533 errors were filtered out. Forward and reverse reads were then independently corrected using  
534 run-specific error rate modelling and dereplicated. Amplicon sequence variants (ASVs; i.e.,  
535 unique sequences) were inferred from these merged reads; Chimeric ASVs were identified and  
536 discarded from the datasets. Next, ASVs were assigned a taxonomy using the RDP naïve Bayesian  
537 classifier<sup>69</sup>, as implemented in DADA2, and using PR2 as a reference database. ASVs classified as  
538 Bacteria, Archaea, Organelle, Metazoa or with no eukaryotic supergroup classification (i.e.,  
539 classified only as “Eukaryota”) were discarded. The final Monterey Bay V9 census dataset was  
540 comprised of 1,073 ASVs representing a total of 89,376 quality controlled, merged sequences  
541 (Supplementary Table 6). Comparison between V9 sequences from Monterey Bay SAGs and  
542 environmental census, in terms of sequence identity (Supplementary Table 5), were conducted  
543 using EMBOSS Water pairwise sequence alignment<sup>70</sup>. Subsequent V9 analyses were conducted  
544 using the R package Phyloseq v1.20<sup>71</sup>.

545  
546 ***Mitochondrial genome contig identification, re-assembly, annotation, and confirmation***

547  
548 In 70 of 99 (70%) SAG assemblies, contigs encoding multiple mitochondrial-like genes were  
549 identified from the assembly. To ensure that no contaminating DNAs were included in our  
550 analysis we removed any contigs >90% identical to known bacterial, chloroplast, or  
551 contaminating (e.g. fungal) mitochondrial DNAs. To obtain better mitochondrial genome  
552 assemblies, reads mapping to each of the identified mitochondrial scaffolds for each SAG were  
553 extracted (using BWA<sup>72</sup>, SAMTOOLS<sup>73</sup>, and BAMTOOLS<sup>74</sup>) and reassembled with SPAdes 3.7.1<sup>49</sup>  
554 in assembly-only mode. The best assemblies were chosen for further analysis, manual  
555 adjustment, and annotation. Mitochondrial genes, including introns, were annotated using  
556 mfannot (<http://megasun.bch.umontreal.ca/cgi-bin/mfannot/mfannotInterface.pl>) with manual  
557 correction as needed. Myzozoan ribosomal fragments in the As1 mitochondrial genome were  
558 identified by nhmmer<sup>75</sup> searches with HMMER v3.1 using Hidden Markov Models generated from  
559 alignments of known fragments (evalue < 1e<sup>-5</sup>)<sup>76</sup>. Complete or near-complete contigs (see below)  
560 were used as queries to identify shorter (i.e. encoding only single mitochondrial proteins or RNA

561 genes) *bona fide* mitochondrial contigs in assemblies from closely related cells. Mitochondrial  
562 genome completion percentages were estimated by comparing incomplete mitochondrial  
563 genomes to complete (100% circular) or near complete genomes (arbitrarily designated at 95%  
564 when no or nearly no coding sequence is missing based on comparisons with closely related taxa).

565  
566 Samples T1, T12, K4, As1 (mapped near ascomycetes but top BLAST is picozoan), H2, C14,  
567 C15, S11, S17, S18, as well as *Incisomonas marina* and two MAST4 mitochondrial genomes from  
568 previous studies assembled into complete circular genomes. T11 could be assembled into a single  
569 contig but could not be circularized. R32 reassembled into a single linear contig with repeats at  
570 5' and 3' ends and was used to identify contigs in similar SAGs. R16 and R17 have identical nSSU  
571 rRNA and nearly identical mitochondrial sequences (>99.9%) and were used to infer a likely-  
572 complete linear mitochondrial genome molecule with repeats at both 5' and 3' ends similar to  
573 R32. R1 and R2 also have identical SSU and nearly identical mitochondrial sequences (>99.5%).  
574 Overlapping contigs from R1 and R2 were joined to form two large contigs that could not be  
575 confidently joined further. As2 (mapped on V9-SSU rDNA phylogenetic trees near ascomycetes  
576 but was actually a rhizarian cell) contained a *Mataza*-like SSU and assembled into seven contigs  
577 that could not be joined but contained nearly all of the predicted genes present in R17 and R32.  
578 M6 and M7 mitochondrial genomes were near identical (>99.6%) and were used to infer a near-  
579 complete *Mataza* mitochondrial genome consisting of two non-overlapping contigs. Two SAGs  
580 related to thraustochytrids, S2 and S4, contained single large mitochondrial genome contigs that  
581 could not be circularized by PCR. However, based on synteny, the missing stretches of DNA could  
582 be inferred since the sequences lacking were present in the reciprocal SAG (shaded and labelled  
583 'inferred' in Fig. 4). S16 (MAST3g) assembled into a single contig and appears to be complete in  
584 terms of coding content, however, a repeat region was assembled in the 3' region of the contig  
585 which appears to contain fragments of *cox2* which could indicate the presence of an inverted  
586 repeat. We could not verify this as we did not recover any other MAST3g SAGs. Similarly, S6  
587 (MAST12) and S14 (MAST8b) were assembled into 2 and 3 contigs respectively. S14 appears  
588 complete with respect to coding content, although the contigs could not be joined. S6 was  
589 incomplete, but when compared to the coding content of its closest sequenced relative *Cafeteria*  
590 *roebergensis* (which also contains a TGA-W code change) it lacked only 7 of 32 genes and

591 therefore was estimated at 78% complete. Complete and near-complete mitochondrial genomes  
592 were visualized using the CGview server<sup>77</sup> and manually edited for figure construction. Closely  
593 related mitochondrial genome molecules were manually examined for synteny (Fig. 4 inner  
594 coloured circles within boxed mitochondrial genomes).

595  
596 Since mitochondrial genomes were well represented in SAG assemblies, we calculated  
597 the relative coverage of mitochondrial genomes compared to the total SAG assembly. We  
598 defined relative coverage as the minimum read coverage over 80% of the representative genome  
599 as defined by BamQC in BAMTOOLS output reports (Supplementary Tables 3 and 4). 51x coverage  
600 was the maximum coverage in the output of this tool. The relationship between relative  
601 mitochondrial genome coverage was compared with that of the nuclear coverage (as estimated  
602 by CEGMA%) using the 'ggplot2' (v2.2.1)<sup>78</sup> and 'DescTools' (v0.99.23)<sup>79</sup> packages in the R (v3.4.3)  
603 programming language<sup>80</sup>. A two-sided Hotelling's T2 test<sup>27</sup> (df1=2, df2=30, T.2=44.942, p=9.07e-  
604 13) was used to test whether the groupings of SAGs showing high mitochondrial coverage (n=17)  
605 and those with high nuclear coverage (n=16) were sampled from populations showing distinct  
606 template profiles. This was performed under the assumptions that they were independently  
607 sampled from multivariate normal distributions with approximately equal covariance matrices.

### 608 609 ***Identification of an alternative genetic code in thraustochytrids***

610  
611 The recovered thraustochytrid mitochondrial genomes (S1-4 and 15) use TTA as a stop codon and  
612 contain in-frame TAG and TAA codons that align with conserved tyrosine residues when  
613 compared to homologues in other thraustochytrids (Extended Data 4), suggesting that these stop  
614 codons have been reassigned to code for tyrosine. *Cob* genes with internal stop codons were  
615 identified in mitochondrial contigs from each SAG and translated using the standard genetic  
616 code. These genes were aligned using MUSCLE<sup>81</sup> with publicly available *cob* genes from  
617 thraustochytrid mitochondrial genomes (KU183024.1 and AF288091.2) (Extended Data 4). The  
618 lack of a tRNA containing the UAA anticodon and the presence of a tRNA with an AAU anticodon  
619 corroborates this hypothesis (Extended Data 5). Since *Thraustochytrium aureum* is known to have  
620 reassigned TTA to a stop codon (GenBank: AF288091.2), these findings support the sister  
621 relationship of thraustochytrids and the phylogenetically related SAGs sampled here (Fig. 6).

622 ***Phylogenetic analysis of representative stramenopiles from concatenated mitochondria-***  
623 ***encoded ETC proteins***

624  
625 Since mitochondrial ribosomes and ribosomal proteins are fast evolving and have a greater  
626 propensity to be lost or relocated to the nucleus, we chose to reconstruct a phylogeny of the  
627 stramenopiles using 16 conserved mitochondria-encoded ETC proteins. These included Nad1, 2,  
628 3, 4, 4L, 5, 6, 7, 9, Cob, Cox1-3, and Atp6, 8, 9. After alignment and manual trimming using  
629 Mesquite v2.75, this resulted in a concatenated alignment with 4442 sites. IQ-Tree<sup>82</sup> was used  
630 for model testing resulting in LG as the highest scoring model by BIC. Phylogenetic tree  
631 reconstructions were performed using MrBayes v3.2.6 for Bayesian analysis<sup>83</sup>. MrBayes analyses  
632 were run with the following parameters prset aamodelpr = fixed (WAG); mcmcngen = 1,000,000;  
633 samplefreq = 1000; nchains = 4; startingtree = random; sumt burnin = 250. Split frequencies were  
634 checked to ensure convergence. Maximum likelihood bootstrap values (100 pseudoreplicates)  
635 were obtained using RAxML v8.2.10<sup>84</sup> under the LG model<sup>85</sup>.

636  
637 ***Phylogenetic analysis of Cox1 proteins from diverse eukaryotes***

638  
639 Cox1 proteins were collected from representative eukaryote groups from the NCBI non-  
640 redundant protein database using BLAST<sup>26</sup>. Resulting sequences were aligned using MUSCLE<sup>81</sup>,  
641 and manually trimmed to a resulting 402 sites. A phylogenetic reconstruction was conducted  
642 using RAxML v8.2.10<sup>84</sup> (100 bootstrap pseudoreplicates) under the LG model<sup>85</sup>.

643  
644 **Data Availability:**

645  
646 Complete mtDNAs assembled from this study are found at NCBI under the accessions MK188935-47,  
647 MN082144-5. Sequencing data can be found under NCBI BioProject: PRJNA379597. Reads are deposited  
648 at NCBI SRA: SRP102236. Partial mtDNA contigs and other important contigs mentioned in text are  
649 available at DOI: 10.6084/m9.figshare.7314728. Nuclear SAG assemblies are available at DOI:  
650 10.6084/m9.figshare.7352966. A public method can be accessed at:  
651 [dx.doi.org/10.17504/protocols.io.ywpxdn](https://dx.doi.org/10.17504/protocols.io.ywpxdn)

652  
653 **Code availability:**

654 Bioinformatic workflow is published at DOI: 10.5281/zenodo.192677. List of programs used in this study:  
655 BD FACS 'Software' sorter software v1.0.0.650, PEAR 0.9.8, Trim Galore!  
656 www.bioinformatics.babraham.ac.uk/projects/trim\_galore/, SPAdes 3.7.1, QUASt 5.0.2, CEGMA v2,  
657 Blobtools v1.0, Qualimap v2.2.1, BWA 0.7.17, SAMTOOLS 1.9, BAMTOOLS 2.4.0, PyNAST v1.2, Trimal v1.4,  
658 RAxML v8.2 , FastTree v2.1 , DADA2 v1.4, Phyloseq v1.20, mfannot [http://megasun.bch.umontreal.ca/cgi-](http://megasun.bch.umontreal.ca/cgi-bin/mfannot/mfannotInterface.pl)  
659 [bin/mfannot/mfannotInterface.pl](http://megasun.bch.umontreal.ca/cgi-bin/mfannot/mfannotInterface.pl), HMMER v3.1 , ggplot2 v2.2.1, DescTools v0.99.23, MUSCLE  
660 <https://www.ebi.ac.uk/Tools/msa/muscle/>, MrBayes v3.2.6.

661  
662 **Author contributions**

663  
664 JGW performed bioinformatic and phylogenetic analyses, and wrote the manuscript. RR-M performed  
665 molecular biological analyses. AM performed bioinformatic and phylogenetic analyses and GL performed  
666 bioinformatic analyses. EC and CP collected the samples and performed flow cytometry. FM performed  
667 statistical and bioinformatic analyses. DM performed molecular biological experiments and generated  
668 biochemical reagents. KM performed the genome sequencing. NATI analysed genomic data. TAR devised  
669 the project. JGW, AES, PJK, AZW, and TAR supervised the project and wrote the manuscript. All authors  
670 contributed to the editing of the final manuscript.

671  
672 **Competing interests**

673  
674 The authors declare no competing interests.

675  
676 **Corresponding authors**

677  
678 Jeremy G. Wideman: [Jeremy.Wideman@asu.edu](mailto:Jeremy.Wideman@asu.edu) and Thomas A. Richards: [T.A.Richards@exeter.ac.uk](mailto:T.A.Richards@exeter.ac.uk).

679  
680 **Acknowledgements**

681  
682 We would like to thank Franz Lang and Natacha Beck for annotation assistance and access to an  
683 unreleased version of mfannot, Dana Price for assistance with picozoan SAG data, and Cory Dunn for  
684 fruitful discussions and encouragement.

685  
686 This project was supported by a Gordon and Betty Moore foundation grant (GBMF3307) to TAR,  
687 AES, AZW and PJK, and a Philip Leverhulme Award (PLP-2014-147) to TAR. Field sampling was supported  
688 by the David and Lucile Packard Foundation and GBMF3788 to AZW. TAR and AM are supported by Royal  
689 Society University Research Fellowships. JGW was supported by the European Molecular Biology

690 Organization Long-term Fellowship (ALTF 761-2014) co-funded by European Commission  
691 (EMBOCOFUND2012, GA-2012-600394) support from Marie Curie Actions and a College for Life Sciences  
692 Fellowship at the Wissenschaftskolleg zu Berlin. RRM is supported by CONICYT FONDECYT 11170748. FM  
693 is supported by Genome Canada.

694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742

## References:

1. Martijn, J., Vosseberg, J., Guy, L., Offre, P. & Ettema, T. J. G. Deep mitochondrial origin outside the sampled alphaproteobacteria. *Nature* **557**, 101–105 (2018).
2. Roger, A. J., Muñoz-Gómez, S. A. & Kamikawa, R. The Origin and Diversification of Mitochondria. *Curr. Biol.* **27**, R1177–R1192 (2017).
3. Martin & Herrmann. Gene transfer from organelles to the nucleus: how much, what happens, and Why? *Plant Physiol.* **118**, 9–17 (1998).
4. Gray, M. W. *et al.* Genome structure and gene content in protist mitochondrial DNAs. *Nucleic Acids Res.* **26**, 865–78 (1998).
5. Worden, A. Z. *et al.* Rethinking the marine carbon cycle: Factoring in the multifarious lifestyles of microbes. *Science.* **347**, (2015).
6. Cuvelier, M. L. *et al.* Targeted metagenomics and ecology of globally important uncultured eukaryotic phytoplankton. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 14679–84 (2010).
7. Worden, A. Z. *et al.* Global distribution of a wild alga revealed by targeted metagenomics. *Curr. Biol.* **22**, R675–R677 (2012).
8. Keeling, P. J. *et al.* The Marine Microbial Eukaryote Transcriptome Sequencing Project (MMETSP): Illuminating the Functional Diversity of Eukaryotic Life in the Oceans through Transcriptome Sequencing. *PLoS Biol.* **12**, e1001889 (2014).
9. Gawryluk, R. M. R. *et al.* Morphological Identification and Single-Cell Genomics of Marine Diplonemids. *Curr. Biol.* **26**, 3053–3059 (2016).
10. Strassert, J. F. H. *et al.* Single cell genomics of uncultured marine alveolates shows paraphyly of basal dinoflagellates. *ISME J.* **12**, 304–308 (2018).
11. Yoon, H. S. *et al.* Single-Cell Genomics Reveals Organismal Interactions in Uncultivated Marine Protists. *Science.* **332**, 714–717 (2011).
12. Bhattacharya, D. *et al.* Single cell genome analysis supports a link between phagotrophy and primary plastid endosymbiosis. *Sci. Rep.* **2**, 356 (2012).
13. Roy, R. S. *et al.* Single cell genome analysis of an uncultured heterotrophic stramenopile. *Sci. Rep.* **4**, 4780 (2014).
14. Mangot, J.-F. *et al.* Accessing the genomic information of unculturable oceanic picoeukaryotes by combining multiple single cells. *Sci. Rep.* **7**, 41498 (2017).
15. Seeleuthner, Y. *et al.* Single-cell genomics of multiple uncultured stramenopiles reveals underestimated functional diversity across oceans. *Nat. Commun.* **9**, 310 (2018).
16. Martinez-Garcia, M. *et al.* Unveiling in situ interactions between marine protists and bacteria through single cell sequencing. *ISME J.* **6**, 703–707 (2012).

- 743 17. Sieracki, M. E. *et al.* Single cell genomics yields a wide diversity of small planktonic protists across  
744 major ocean ecosystems. *Sci. Rep.* **9**, 6025 (2019).  
745
- 746 18. Rose, J., Caron, D., Sieracki, M. & Poulton, N. Counting heterotrophic nanoplanktonic protists in  
747 cultures and aquatic communities by flow cytometry. *Aquat. Microb. Ecol.* **34**, 263–277 (2004).  
748
- 749 19. Richards, T. A. & Talbot, N. J. Horizontal gene transfer in osmotrophs: playing with public goods.  
750 *Nat. Rev. Microbiol.* **11**, 720–727 (2013).  
751
- 752 20. Vrieling, E. G., Gieskes, W. W. C. & Beelen, T. P. M. SILICON DEPOSITION IN DIATOMS: CONTROL  
753 BY THE pH INSIDE THE SILICON DEPOSITION VESICLE. *J. Phycol.* **35**, 548–559 (1999).  
754
- 755 21. Kawai, A., Uchiyama, H., Takano, S., Nakamura, N. & Ohkuma, S. Autophagosome-Lysosome  
756 Fusion Depends on the pH in Acidic Compartments in CHO Cells. *Autophagy* **3**, 154–157 (2007).  
757
- 758 22. Wilken, S. *et al.* The need to account for cell biology in characterizing predatory mixotrophs in  
759 aquatic environments. *Philos. Trans. B* (2019). doi:10.1098/rstb.2019.0090  
760
- 761 23. Dean, F. B. *et al.* Comprehensive human genome amplification using multiple displacement  
762 amplification. *Proc. Natl. Acad. Sci.* **99**, 5261–5266 (2002).  
763
- 764 24. Richards, T. A. & Bass, D. Molecular screening of free-living microbial eukaryotes: diversity and  
765 distribution using a meta-analysis. *Curr. Opin. Microbiol.* **8**, 240–252 (2005).  
766
- 767 25. Richards, T. A., Jones, M. D. M., Leonard, G. & Bass, D. Marine Fungi: Their Ecology and Molecular  
768 Diversity. (2011). doi:10.1146/annurev-marine-120710-100802  
769
- 770 26. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database search  
771 programs. *Nucleic Acids Res.* **25**, 3389–402 (1997).  
772
- 773 27. Hotelling, H. The Generalization of Student's Ratio. *Ann. Math. Stat.* **2**, 360–378 (1931).  
774
- 775 28. Derelle, R., López-García, P., Timpano, H. & Moreira, D. A Phylogenomic Framework to Study the  
776 Diversity and Evolution of Stramenopiles (=Heterokonts). *Mol. Biol. Evol.* **33**, 2890–2898 (2016).  
777
- 778 29. Flegontov, P. *et al.* Divergent mitochondrial respiratory chains in phototrophic relatives of  
779 apicomplexan parasites. *Mol. Biol. Evol.* **32**, 1115–1131 (2015).  
780
- 781 30. Janouškovec, J. *et al.* A New Lineage of Eukaryotes Illuminates Early Mitochondrial Genome  
782 Reduction. *Curr. Biol.* **27**, 3717–3724.e5. (2017). doi:10.1016/j.cub.2017.10.051  
783
- 784 31. Gray, M. W., Lang, B. F. & Burger, G. Mitochondria of Protists. *Annu. Rev. Genet.* **38**, 477–524  
785 (2004).  
786
- 787 32. Wang, Z. *et al.* Complete mitochondrial genome of a DHA-rich protist *Schizochytrium* sp.  
788 TIO1101. *Mitochondrial DNA Part B* **1**, 126–127 (2016).  
789
- 790 33. Saldanha, R., Mohr, G., Belfort, M. & Lambowitz, A. M. Group I and group II introns. *FASEB J.* **7**,  
791 15–24 (1993).



- 792 34. Goddard, M. R. & Burt, A. Recurrent invasion and extinction of a selfish gene. *Proc. Natl. Acad. Sci. U. S. A.* **96**, 13880–5 (1999).  
793  
794
- 795 35. Hauth, A. M., Maier, U. G., Lang, B. F. & Burger, G. The *Rhodomonas salina* mitochondrial  
796 genome: bacteria-like operons, compact gene arrangement and complex repeat region. *Nucleic  
797 Acids Res.* **33**, 4433–4442 (2005).  
798
- 799 36. Kim, E. *et al.* Complete Sequence and Analysis of the Mitochondrial Genome of *Hemiselmis  
800 andersenii* CCMP644 (Cryptophyceae). *BMC Genomics* **9**, 215 (2008).  
801
- 802 37. Nishimura, Y. *et al.* Mitochondrial Genome of *Palpitomonas bilix*: Derived Genome Structure and  
803 Ancestral System for Cytochrome c Maturation. *Genome Biol. Evol.* **8**, 3090–3098 (2016).  
804
- 805 38. Riisberg, I. *et al.* Seven Gene Phylogeny of Heterokonts. *Protist* **160**, 191–204 (2009).  
806
- 807 39. Oudot-Le Secq, M.-P., Loiseaux-de Goër, S., Stam, W. T. & Olsen, J. L. Complete mitochondrial  
808 genomes of the three brown algae (Heterokonta: Phaeophyceae) *Dictyota dichotoma*, *Fucus  
809 vesiculosus* and *Desmarestia viridis*. *Curr. Genet.* **49**, 47–58 (2006).  
810
- 811 40. Leonard, G. *et al.* Comparative genomic analysis of the ‘pseudofungus’ *Hyphochytrium  
812 catenoides*. *Open Biol.* **8**, 170184 (2018).  
813
- 814 41. Massana, R., del Campo, J., Sieracki, M. E., Audic, S. & Logares, R. Exploring the uncultured  
815 microeukaryote majority in the oceans: reevaluation of ribogroups within stramenopiles. *ISME J.*  
816 **8**, 854–866 (2014).  
817
- 818 42. Kannan, S., Rogozin, I. B. & Koonin, E. V. MitoCOGs: clusters of orthologous genes from  
819 mitochondria and implications for the evolution of eukaryotes. *BMC Evol. Biol.* **14**, 237 (2014).  
820
- 821 43. Ševčíková, T. *et al.* A Comparative Analysis of Mitochondrial Genomes in Eustigmatophyte Algae.  
822 *Genome Biol. Evol.* **8**, 705–722 (2016).  
823
- 824 44. Johnston, I. G. & Williams, B. P. Evolutionary Inference across Eukaryotes Identifies Specific  
825 Pressures Favoring Mitochondrial Gene Retention. *Cell Syst.* **2**, 101–111 (2016).  
826
- 827 45. Keeling, P. J. Genomics: Evolution of the Genetic Code. *Curr. Biol.* **26**, R851–R853 (2016).  
828
- 829 46. Demir-Hilton, E. *et al.* Global distribution patterns of distinct clades of the photosynthetic  
830 picoeukaryote *Ostreococcus*. *ISME J.* **5**, 1095–1107 (2011).  
831
- 832 47. Logares, R. *et al.* Patterns of rare and abundant marine microbial eukaryotes. *Curr. Biol.* **24**, 813–  
833 21 (2014).  
834
- 835 48. Zhang, J., Kobert, K., Flouri, T. & Stamatakis, A. PEAR: a fast and accurate Illumina Paired-End  
836 reAd mergeR. *Bioinformatics* **30**, 614–620 (2014).  
837
- 838 49. Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-  
839 Cell Sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).

- 840 50. Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUASt: quality assessment tool for genome  
841 assemblies. *Bioinformatics* **29**, 1072–1075 (2013).  
842
- 843 51. Parra, G., Bradnam, K. & Korf, I. CEGMA: a pipeline to accurately annotate core genes in  
844 eukaryotic genomes. *Bioinformatics* **23**, 1061–1067 (2007).  
845
- 846 52. Kumar, S., Jones, M., Koutsovoulos, G., Clarke, M. & Blaxter, M. Blobology: exploring raw genome  
847 data for contaminants, symbionts and parasites using taxon-annotated GC-coverage plots. *Front.*  
848 *Genet.* **4**, 237 (2013).  
849
- 850 53. Okonechnikov, K., Conesa, A. & García-Alcalde, F. Qualimap 2: advanced multi-sample quality  
851 control for high-throughput sequencing data. *Bioinformatics* **32**, btv566 (2015).  
852
- 853 54. Amaral-Zettler, L. A., McCliment, E. A., Ducklow, H. W. & Huse, S. M. A Method for Studying  
854 Protistan Diversity Using Massively Parallel Sequencing of V9 Hypervariable Regions of Small-  
855 Subunit Ribosomal RNA Genes. *PLoS One* **4**, e6372 (2009).  
856
- 857 55. Guillou, L. *et al.* The Protist Ribosomal Reference database (PR2): a catalog of unicellular  
858 eukaryote small sub-unit rRNA sequences with curated taxonomy. *Nucleic Acids Res.* **41**, D597-  
859 604 (2013).  
860
- 861 56. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation  
862 sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).  
863
- 864 57. Caporaso, J. G. *et al.* PyNAST: a flexible tool for aligning sequences to a template alignment.  
865 *Bioinformatics* **26**, 266–267 (2010).  
866
- 867 58. Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data processing and  
868 web-based tools. *Nucleic Acids Res.* **41**, D590–D596 (2012).  
869
- 870 59. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated alignment  
871 trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–3 (2009).  
872
- 873 60. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
874 phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).  
875
- 876 61. Berger, S. A., Krompass, D. & Stamatakis, A. Performance, Accuracy, and Web Server for  
877 Evolutionary Placement of Short Sequence Reads under Maximum Likelihood. *Syst. Biol.* **60**, 291–  
878 302 (2011).  
879
- 880 62. Shimodaira, H. & Hasegawa, M. Multiple Comparisons of Log-Likelihoods with Applications to  
881 Phylogenetic Inference. *Mol. Biol. Evol.* **16**, 1114–1116 (1999).  
882
- 883 63. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 – Approximately Maximum-Likelihood Trees for  
884 Large Alignments. *PLoS One* **5**, e9490 (2010).  
885
- 886 64. Junier, T. & Zdobnov, E. M. The Newick utilities: high-throughput phylogenetic tree processing in  
887 the UNIX shell. *Bioinformatics* **26**, 1669–1670 (2010).

- 888 65. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T.-Y. ggtree : an r package for visualization and  
889 annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol.*  
890 *Evol.* **8**, 28–36 (2017).  
891
- 892 66. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
893 *EMBnet.journal* **17**, 10 (2011).  
894
- 895 67. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).  
896
- 897 68. Callahan, B. J. *et al.* DADA2: High-resolution sample inference from Illumina amplicon data. *Nat.*  
898 *Methods* **13**, 581–583 (2016).  
899
- 900 69. Wang, Q., Garrity, G. M., Tiedje, J. M. & Cole, J. R. Naive Bayesian classifier for rapid assignment  
901 of rRNA sequences into the new bacterial taxonomy. *Appl. Environ. Microbiol.* **73**, 5261–7 (2007).  
902
- 903 70. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open Software Suite.  
904 *Trends Genet.* **16**, 276–7 (2000).  
905
- 906 71. McMurdie, P. J. & Holmes, S. phyloseq: An R Package for Reproducible Interactive Analysis and  
907 Graphics of Microbiome Census Data. *PLoS One* **8**, e61217 (2013).  
908
- 909 72. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
910 *Bioinformatics* **25**, 1754–1760 (2009).  
911
- 912 73. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079  
913 (2009).  
914
- 915 74. Barnett, D. W., Garrison, E. K., Quinlan, A. R., Stromberg, M. P. & Marth, G. T. BamTools: a C++  
916 API and toolkit for analyzing and managing BAM files. *Bioinformatics* **27**, 1691–1692 (2011).  
917
- 918 75. Wheeler, T. J. & Eddy, S. R. nhmmer: DNA homology search with profile HMMs. *Bioinformatics*  
919 **29**, 2487–2489 (2013).  
920
- 921 76. Jackson, C. J. *et al.* Broad genomic and transcriptional analysis reveals a highly derived genome in  
922 dinoflagellate mitochondria. *BMC Biol.* **5**, 41 (2007).  
923
- 924 77. Grant, J. R. & Stothard, P. The CGView Server: a comparative genomics tool for circular genomes.  
925 *Nucleic Acids Res.* **36**, W181–W184 (2008).  
926
- 927 78. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New York, 2009).  
928
- 929 79. Signorell, A. DescTools: Tools for descriptive statistics. R package version 0.99.23. (2017).  
930
- 931 80. R Core Team. R: A language and environment for statistical computing. (2013). Available at:  
932 <http://www.r-project.org/>. (Accessed: 19th July 2018)  
933
- 934 81. Edgar, R. C. MUSCLE: a multiple sequence alignment method with reduced time and space  
935 complexity. *BMC Bioinformatics* **5**, 113 (2004).  
936
- 937 82. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and Effective

- 938 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol. Biol. Evol.* **32**, 268–  
939 274 (2015).  
940  
941 83. Ronquist, F. & Huelsenbeck, J. P. MrBayes 3: Bayesian phylogenetic inference under mixed  
942 models. *Bioinformatics* **19**, 1572–4 (2003).  
943  
944 84. Stamatakis, A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands  
945 of taxa and mixed models. *Bioinformatics* **22**, 2688–90 (2006).  
946  
947 85. Le, S. Q. & Gascuel, O. An improved general amino acid replacement matrix. *Mol. Biol. Evol.* **25**,  
948 1307–20 (2008).  
949

## 950 Figure Legends

951  
952 **Fig. 1. V9-nSSU phylogenetic mapping of Monterey Bay single amplified genomes.** Maximum-likelihood  
953 (ML) phylogenetic tree of reference nSSU sequences, retrieved and curated from the PR2 reference  
954 database onto which SAG nSSU-V9 sequences (from 206 flagellum-targeted flow-cytometry sorted single  
955 cells from eastern North Pacific waters) were phylogenetically mapped (red circles). The ML tree was  
956 inferred under the GTR-CAT model, based on a multiple sequence alignment of 20,939 PR2 representative  
957 sequences and totalling 1,750 sites. Major eukaryotic clades are labelled (See Fig. 2). Groups with  
958 representative SAGs are shaded in blue. Numbers in brackets beside taxon names indicate the number of  
959 SAGs that were obtained from each taxonomic group.

960  
961 **Fig. 2. Clade-specific ML subtrees.** Six distinct eukaryotic clades from which numerous SAGs were  
962 recovered. Shown here are SAG nSSU V9 sequences mapped to a full-length reference tree that  
963 incorporated PR2 reference sequences. For each subtree, specific lineages that attracted SAG V9  
964 sequences are highlighted in pink frames with the lineage name provided and in parentheses the number  
965 of SAG V9 that mapped onto the lineage. SAGs with mitochondrial contigs present are labelled: complete  
966 mtDNAs, white font on black circle; near complete, bold; partial genome, italics. SH-like local node  
967 supports are shown for > 0.9 (black circles). taxonomic colour legend: Alveolata, orange; Apusozoa,  
968 yellow; Euglenozoa, white; Opisthokonta, grey; Stramenopiles, blue; Amoebozoa, pink; Archaeplastida,  
969 green; Hacrobia, turquoise; Rhizaria, purple. Scale bars represent the number of estimated substitutions  
970 per site.

971  
972 **Fig. 3. Distribution and groupings of mitochondrial sequence coverage relative to estimated nuclear**  
973 **genome completeness.** Sequenced genomes showed either high nuclear completion (CEGMA%) (green,  
974 n=17 biologically independent mitochondrial genomes), high mitochondrial coverage (X-fold coverage of  
975 >80% of mtDNA) (red, n=16 biologically independent mitochondrial genomes) or simultaneously low

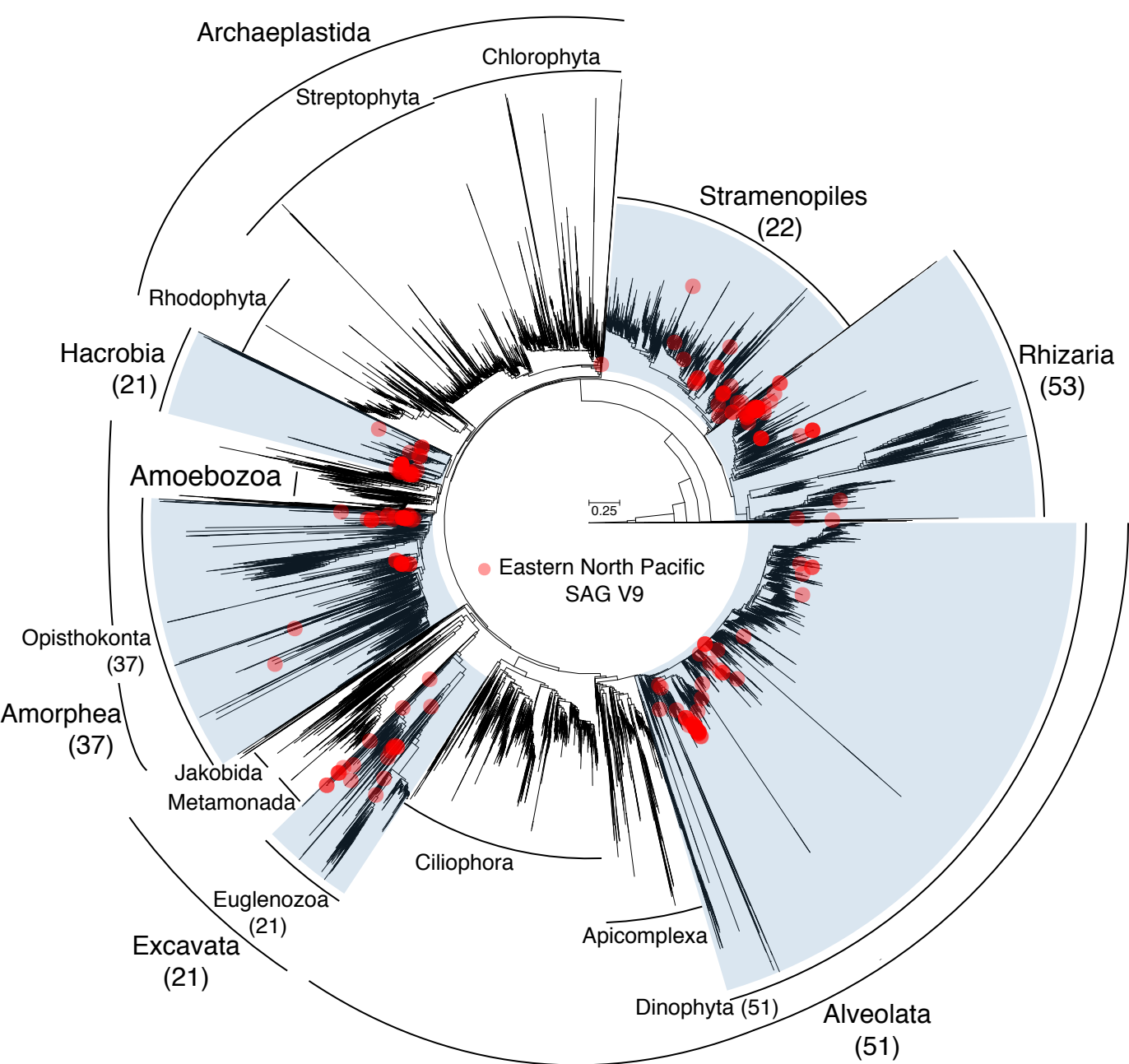
976 nuclear and mitochondrial coverage (blue). The blue density contours were plotted using 'ggplot2'. X-fold  
977 coverage of >80% of each sequenced mtDNA was calculated using BamQC in BAMTOOLS for each SAG  
978 with >50% estimated completion of the mtDNA. This score was plotted against the estimated CEGMA  
979 completion scores (%). The result of a Hotelling's T2-test to assess whether the SAGs with high  
980 mitochondrial coverage and those with high nuclear coverage supports rejection of the null hypothesis  
981 that these are sampled from the same population. The rejection of the null hypothesis suggests that the  
982 there is a fundamental difference between these two SAG sub populations.

983  
984 **Fig. 4. Uncharacterized mtDNAs from underrepresented eukaryotic groups.** Complete and near-  
985 complete mtDNAs assembled from heterotrophic marine flagellate SAGs. Mitochondrial contigs were  
986 annotated using mfanot with manual corrections as needed  
987 (<http://megasun.bch.umontreal.ca/RNAweasel/>). MtDNAs are represented as circular diagrams (bold  
988 central font) or broken circles if contigs could not be joined (regular central font). Complete genomes  
989 assembled herein using publicly available metagenomes and previously published SAG datasets are  
990 marked with an asterisk in the centre of the genome map. Genomes from *Cryothecomonas*-like cells did  
991 not map as circular. Where present, coloured central circles correspond to syntenic regions shared  
992 between closely related genomes (within boxes). Some mtDNAs were inferred from multiple cells with  
993 identical nSSU sequences containing nearly identical stretches of mtDNA sequences that could be stitched  
994 together (see methods). Colour coded genes: blue, protein coding; purple, rRNA; red, tRNA, dark grey,  
995 putative introns.

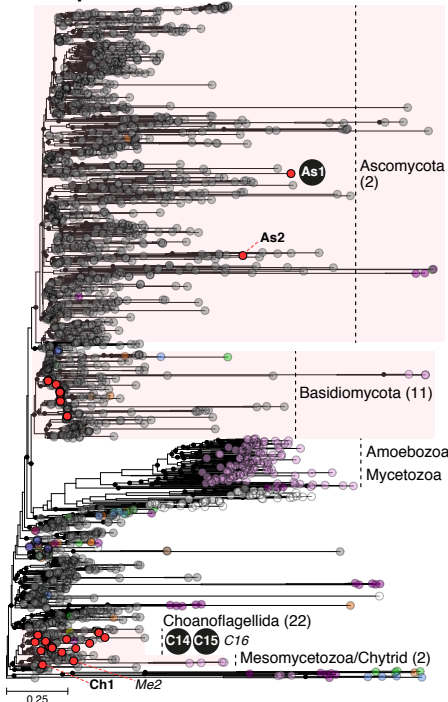
996  
997 **Fig. 5. Comparison between mtDNA gene repertoires.** Mitochondrial genomes newly assembled in this  
998 study (bold font), previously sequenced mtDNAs (regular font), and ancestral reconstructions (L-Dia-CA,  
999 Last Diaphoretickes Common Ancestor; L-Amo-CA, Last Amorphean Common Ancestor - including  
1000 malawimonads and collodictyonids); L-Jak-CA, Last Jakobid Common Ancestor; LECA, Last Eukaryote  
1001 Common Ancestor) are shown. Black square, present; empty square, absent; red square, rare protein. #  
1002 symbols indicate incomplete mtDNA. Asterisks indicate genomes assembled from publicly available  
1003 datasets.

1004  
1005 **Fig. 6. Phylogenetic reconstruction of representative stramenopiles using concatenated conserved**  
1006 **mitochondria-encoded electron transport chain proteins.** Electron transport chain proteins encoded in  
1007 publicly available mtDNAs and our newly sequenced mtDNAs of stramenopiles and rhizarians were  
1008 collected, aligned, masked, and concatenated, resulting in a 16-protein 4442-site alignment. We excluded

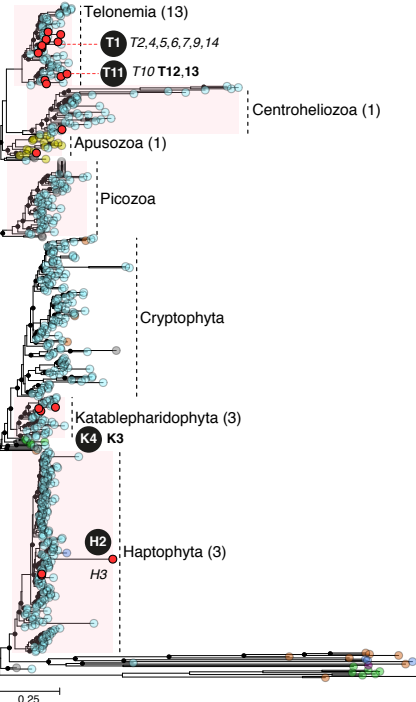
1009 alveolate mtDNAs from this analysis because most of these datasets encode very few (e.g. dinoflagellates  
1010 and apicomplexans) and/or highly divergent proteins (e.g. ciliates). Phylogenies were reconstructed and  
1011 node support values were calculated using MrBayes v3.2.6 for posterior probability<sup>83</sup> and RAxML v8.2.10  
1012 for maximum likelihood<sup>84</sup> and presented as inset (MrBayes/RAxML). The MrBayes tree topology is shown.  
1013 Changes in genetic code are mapped to nodes as indicated. Genes encoding electron transport chain  
1014 components (*atp1*, *nad7*, *nad9*, *nad11*) that have putatively moved to the nucleus are bolded and mapped  
1015 to nodes as indicated. The *atp1* gene has been lost within the opalozoans and is indicated with a  
1016 strikethrough. N-*nad11* indicates that the N-terminal domain of *nad11* is encoded in the nucleus while C-  
1017 *nad11* indicates the C-terminal domain of *nad11* is encoded in the nucleus. MtDNAs presented in this  
1018 study are indicated in bold. Percentages indicate the estimated completeness of each mtDNA presented  
1019 in this study.



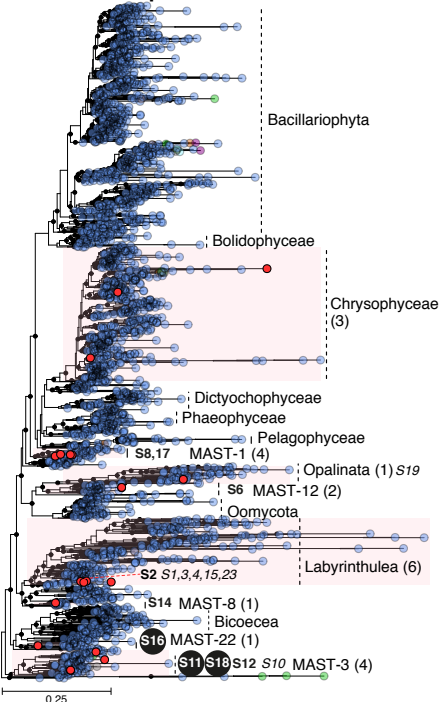
### 1 - Opisthokonta



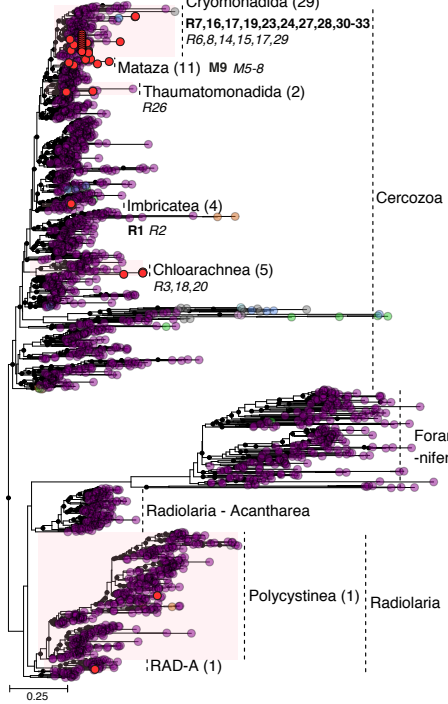
### 2 - Hacrobia



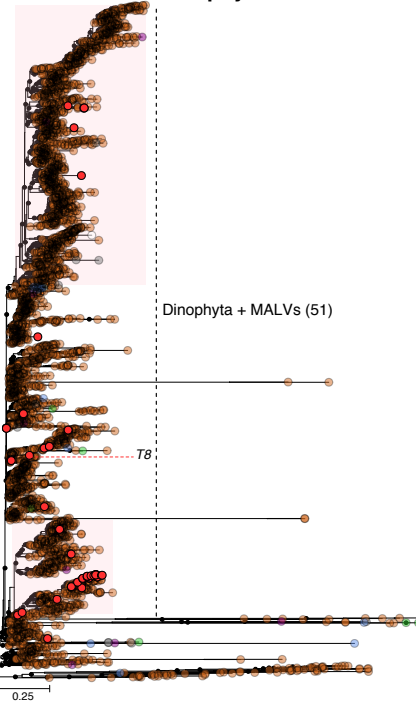
### 3 - Stramenopiles



### 4 - Rhizaria



### 5 - Alveolata / Dinophyta



### 6 - Excavata / Euglenozoa

