

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

Single cell genomics reveals plastid-lacking Picozoa are close relatives of red algae

Max E. Schön^{1,2}, Vasily V. Zlatogursky¹, Rohan P. Singh³, Camille Poirier^{4,5}, Susanne Wilken⁵, Varsha Mathur⁶, Jürgen F. H. Strasser¹, Jarone Pinhassi⁷, Alexandra Z. Worden^{4,5}, Patrick J. Keeling⁶, Thijs J. G. Ettema⁸, Jeremy G. Wideman³, Fabien Burki^{1,9*}

¹Department of Organismal Biology, Program in Systematic Biology, Uppsala University, Uppsala, Sweden

²Department of Cell and Molecular Biology, Program in Molecular Evolution, Uppsala University, Uppsala, Sweden

³Biodesign Center for Mechanisms of Evolution, School of Life Sciences, Arizona State University, Tempe, AZ, USA

⁴Ocean EcoSystems Biology, RD3, GEOMAR Helmholtz Centre for Ocean Research Kiel, Germany

⁵Monterey Bay Aquarium Research Institute, Moss Landing, CA, USA

⁶Department of Botany, University of British Columbia, Vancouver, British Columbia, Canada

⁷Centre for Ecology and Evolution in Microbial Model Systems – EEMiS, Linnaeus University, Kalmar, Sweden

⁸Laboratory of Microbiology, Wageningen University and Research, NL-6708 WE Wageningen, The Netherlands

⁹Science for Life Laboratory, Uppsala University, Uppsala, Sweden

Present address (VVZ): Department of Invertebrate Zoology, Faculty of Biology, St. Petersburg State University, Russia

Present address (JFHS): Department of Ecosystem Research, Leibniz Institute of Freshwater Ecology and Inland Fisheries, Berlin, Germany

Present address (CP): Department of Zoology, University of Oxford, 11a Mansfield Road, Oxford OX1 3SZ, UK

Present address (SW): Institute for Biodiversity and Ecosystem Dynamics, University of Amsterdam, NL

*Corresponding author: fabien.burki@ebc.uu.se

28 **Abstract**

29 The endosymbiotic origin of plastids from cyanobacteria gave eukaryotes photosynthetic capabilities and
30 launched the diversification of countless forms of algae. These primary plastids are found in members of the
31 eukaryotic supergroup Archaeplastida. All known archaeplastids still retain some form of primary plastids,
32 which are widely assumed to have a single origin. Here, we used single-cell genomics from natural samples
33 combined with phylogenomics to infer the evolutionary origin of the phylum Picozoa, a globally distributed but
34 seemingly rare group of marine microbial heterotrophic eukaryotes. Strikingly, the analysis of 43 single-cell
35 genomes shows that Picozoa belong to Archaeplastida, specifically related to red algae and the phagotrophic
36 rhodelphids. These picozoan genomes support the hypothesis that Picozoa lack a plastid, and further reveal no
37 evidence of an early cryptic endosymbiosis with cyanobacteria. These findings change our understanding of
38 plastid evolution as they either represent the first complete plastid loss in a free-living taxon, or indicate that red
39 algae and rhodelphids obtained their plastids independently of other archaeplastids.

40

41 **Introduction**

42 The origin of plastids by endosymbiosis between a eukaryotic host and cyanobacteria was a fundamental
43 transition in eukaryotic evolution, giving rise to the first photosynthetic eukaryotes. These ancient primary
44 plastids, estimated to have originated >1.8 billion years ago¹, are found in Rhodophyta (red algae),
45 Chloroplastida (green algae, including land plants), and Glaucophyta (glaucophytes)—together forming the
46 eukaryotic supergroup Archaeplastida². Unravelling the complex sequence of events leading to the establishment
47 of the cyanobacterial endosymbiont in Archaeplastida is complicated by antiquity, and by the current lack of
48 modern descendants of early-diverging relatives of the main archaeplastidan groups in culture collections or
49 sequence databases. Indeed, the only other known example of primary endosymbiosis are the chromatophores in
50 one unrelated genus of amoeba (*Paulinella*), which originated about a billion years later^{3,4}. Recently, two newly
51 described phyla (Prasinodermophyta and Rhodelphidia) were found to branch as sister to green and red algae,
52 respectively^{5,6}. Most transformative was the discovery that rhodelphids are obligate phagotrophs that maintain
53 cryptic non-photosynthetic plastids, implying that the ancestor of red algae was mixotrophic, a finding that
54 greatly alters our perspectives on early archaeplastid evolution⁵.

55 While there is substantial evidence that Archaeplastida is a group descended from a photosynthetic ancestor,
56 non-photosynthetic and plastid-lacking lineages have been found to branch near the base or even within
57 archaeplastids in phylogenomic trees. For example, Cryptista (which includes plastid-lacking and secondary
58 plastid-containing species), have been inferred to be sister to either green algae and glaucophytes⁷, or red
59 algae^{5,8}, although other phylogenomic analyses have recovered the monophyly of Archaeplastida to the
60 exclusion of the cryptists^{5,9,10}. Another non-photosynthetic group that recently showed affinities to red algae
61 based on phylogenomics is Picozoa^{5,9,10}. But as for cryptists, the position of Picozoa has lacked consistent

62 support, mostly because there is no member of Picozoa available in continuous culture, and genomic data are
63 currently restricted to a few, incomplete, single amplified genomes (SAGs)¹¹. Thus, the origin of Picozoa
64 remains unclear.

65 Picozoa (previously known as picobiliphytes) were first described in 2007 in marine environmental clone
66 libraries of the 18S ribosomal RNA (rRNA) gene and observed by epifluorescence microscopy in temperate
67 waters¹². Due to orange autofluorescence reminiscent of the photosynthetic pigment phycobiliprotein and
68 emanating from an organelle-like structure, picozoans were initially described as likely containing a plastid.
69 Orange fluorescence was also observed in association with these uncultured cells in subtropical waters¹³.
70 However, the hypothesis that the cells were photosynthetic was challenged by the characterization of SAG data
71 from three picozoan cells isolated by fluorescence-activated cell sorting (FACS)¹¹. The analysis of these SAGs
72 revealed neither plastid DNA nor nuclear-encoded plastid-targeted proteins, but the scope of these conclusions is
73 limited due to the small number of analyzed cells and their highly fragmented and incomplete genomes¹¹. Most
74 interestingly, a transient culture was later established, enabling the formal description of the first (and so far
75 only) picozoan species—*Picomonas judraskeda*—as well as ultrastructural observations with electron
76 microscopy¹⁴. These observations revealed an unusual structural feature in two body parts, a feeding strategy by
77 endocytosis of nano-sized colloid particles, and confirmed the absence of plastids¹⁴. Only the 18S rRNA gene
78 sequence of *P. judraskeda* is available as the transient culture was lost before genomic data could be generated.
79 Here, we present an analysis of genomic data from 43 picozoan single-cell genomes sorted with FACS from the
80 Pacific Ocean off the California coast and from the Baltic Sea. Using a gene and taxon-rich phylogenomic
81 dataset, these new data allowed us to robustly infer Picozoa as a lineage of archaeplastids, branching with red
82 algae and rhodophids. With this expanded genomic dataset, we confirm Picozoa as the first archaeplastid
83 lineage lacking a plastid. We discuss the important implications that these results have on our understanding of
84 the origin of plastids.

85

86 **Results**

87 *Single cell assembled genomes representative of Picozoa diversity*

88 We isolated 43 picozoan cells (40 from the eastern North Pacific off the coast of California, 3 from the Baltic
89 Sea) using FACS and performed whole genome amplification by multiple displacement amplification (MDA).
90 The taxonomic affiliation of the SAGs was determined either by PCR with Picozoa-specific primers¹⁴ or 18S
91 rRNA gene sequencing using general eukaryotic primers, followed by Illumina sequencing of the MDA products
92 (see methods). The sequencing reads were assembled into genomic contigs, with a total assembly size ranging
93 from 350 kbp to 66 Mbp (Fig 1a, Table S1). From these contigs, the 18S rRNA gene was found in 37 out of the
94 43 SAGs, which we used to build a phylogenetic tree with reference sequences from the protist ribosomal
95 reference PR2 database (Fig S1). Based on this tree, we identified 6 groups representing 32 SAGs that possessed

96 nearly identical 18S rRNA gene sequences. These SAGs with identical ribotype were reassembled by pooling all
97 reads in order to obtain longer, more complete co-assemblies (CO-SAGs). The genome size of the CO-SAGs
98 ranged from 32 Mbp to 109 Mbp (Fig 1a, Tab S1), an increase of 5% to 45% over individual SAGs. The genome
99 completeness of the SAGs and CO-SAGs was estimated based on two datasets: (i) a set of 255 eukaryotic
100 marker genes available in BUSCO¹⁵, and (ii) a set of 317 conserved marker genes derived from a previous pan-
101 eukaryote phylogenomic dataset¹ that we used here as starting point (Fig 1b). These comparisons showed that
102 while most SAGs were highly incomplete (Fig 1a and b), the CO-SAGs were generally more complete (up to
103 60%). When taken together, 90% of the BUSCO markers and 88% of the phylogenomic markers were present in
104 at least one assembly, suggesting that while the single-cell genome assemblies are fragmentary, they together
105 represent a much more complete Picozoa meta-assembly.

106 The final 17 assemblies (11 SAGs and 6 CO-SAGs) were mainly placed within the three proposed groups of
107 Picozoa BP1-3 (Fig 1c), sensu Cuvelier et al. (2008)¹³. One SAG (SAG11) was placed outside of these groups.
108 The deep-branching picozoan lineages identified by Moreira and López-García (2014)¹⁶, as well as other
109 possibly early-diverging lineages were not represented in our data (Fig 1c). Interestingly, one CO-SAG
110 (COSAG03) was closely related (18S rRNA gene 100% identical) to the only described species, *Picomonas*
111 *judraskeda*, for which no genomic data are available. Using our assemblies and reference sequences from PR2 as
112 queries, we identified by sequence identity 362 OTUs related to Picozoa ($\geq 90\%$) in the data provided by the
113 *Tara* Oceans project¹⁷. Picozoa were found in all major oceanic regions, but had generally low relative
114 abundance in V9 18S rRNA gene amplicon data (less than 1% of the eukaryotic fractions in most cases, Fig S2).
115 An exception was the Southern Ocean between South America and Antarctica, where the Picozoa-related OTUs
116 in one sample represented up to 30% of the V9 18S rRNA gene amplicons. Thus, Picozoa seems widespread in
117 the oceans but generally low in abundance based on available sampling, although they can reach higher relative
118 abundances in at least circumpolar waters.

119 120 *Phylogenomic dataset construction*

121 To infer the evolutionary origin of Picozoa, we expanded on a phylogenomic dataset that contained a broad
122 sampling of eukaryotes and a large number of genes that was recently used to study deep nodes in the eukaryotic
123 tree¹. Homologues from the SAGs and CO-SAGs as well as a number of newly sequenced key eukaryotes were
124 added to each single gene (see Table S2 for a list of taxa). After careful examination of the single genes for
125 contamination and orthology based on individual phylogenies (see material and methods), we retained all six
126 CO-SAGs and four individual SAGs together with the available SAG MS584-11 from a previous study¹¹. The
127 rest of the SAGs were excluded due to poor data coverage (less than 5 markers present) and, in one case
128 (SAG33), because it was heavily contaminated with sequences from a cryptophyte (see Data availability for
129 access to the gene trees). In total, our phylogenomic dataset contained 317 protein-coding genes, with
130 orthologues from Picozoa included in 279 genes (88%), and 794 taxa (Fig 1b). This represents an increase in

131 gene coverage from 18% to 88% compared to the previously available genomic data for Picozoa. The most
132 complete assembly was CO-SAG01, from which we identified orthologues for 163 (51%) of the markers.

133 *Picozoa group with Rhodophyta and Rhodelphidia*

134 Concatenated protein alignments of the curated 317 genes were used to infer the phylogenetic placement of
135 Picozoa in the eukaryotic Tree of Life. Initially, a maximum likelihood (ML) tree was reconstructed from the
136 complete 794-taxa dataset using the site-homogeneous model LG+F+G and ultrafast bootstrap support with
137 1,000 replicates (Fig S3). This analysis placed Picozoa together with a clade comprising red algae and
138 rhodelphids with strong support (100% UFBoot2), but the monophyly of Archaeplastida was not recovered due
139 to the internal placement of the cryptists. To further investigate the position of Picozoa, we applied better-fitting
140 site-heterogeneous models to a reduced dataset of 67 taxa, since these models are computationally much more
141 demanding. The process of taxon reduction was driven by the requirement of maintaining representation from all
142 major groups, while focusing sampling on the part of the tree where Picozoa most likely belong to, i.e.
143 Archaeplastida, TSAR, Haptista and Cryptista. We also merged several closely related lineages into OTUs based
144 on the initial ML tree in order to reduce missing data (Table S3). This 67-taxa dataset was used in ML and
145 Bayesian analyses with the best-fitting site-heterogeneous models LG+C60+F+G+PMSF (with non-parametric
146 bootstrapping) and CAT+GTR+G, respectively. Both ML and Bayesian analyses produced highly similar trees,
147 and received maximal support for the majority of relationships, including deep divergences (Fig 2). Most
148 interestingly, both analyses recovered the monophyly of Archaeplastida (BS=93%; PP=1), with cryptists as sister
149 lineage (BS=100%; PP=1). Consistent with the initial ML tree (Fig S3), red algae and rhodelphids branched
150 together (BS=95%; PP=1), with Picozoa as their sister with full support (BS=100%; PP=1). This grouping was
151 robust to fast-evolving sites removal analysis (Fig S4), trimming of the 25% and 50% compositionally most
152 biased sites (Fig S5), and was also recovered in a supertree method (ASTRAL-III) consistent with the multi-
153 species coalescent model (Fig S6). Although this group is robust, we observed one variation in the branching
154 order between Picozoa, rhodelphids and red algae when trimming the 50% most heterogenous sites (Fig S7) and
155 after removing genes with less than two picozoan sequences (Fig S8). In these analyses, Picozoa and red algae
156 were most closely related, although this relationship was never significantly supported. An approximately
157 unbiased (AU) test rejected all tested topologies except in the two cases where Picozoa branched as the closest
158 sister to red algae ($p=0.237$) and the topology of Fig 2 ($p=0.822$; Table S4). Finally, we identified in Picozoa and
159 rhodelphids a two amino acids replacement signature in the eukaryotic translation elongation factor 2 protein
160 (SA instead of the ancestral GS residues, see Supplementary Data S1) that was previously shown to unite red
161 and green algae (and land plants), haptophytes and some cryptists¹⁸. The presence of SA in Picozoa supports
162 their affiliation with red algae and rhodelphids.

163

164 *Picozoa SAGs show no evidence of a plastid*

165 Since there have been conflicting conclusions about the occurrence of plastids in picozoans, we extensively
166 searched our genomic data for evidence of cryptic plastids. First, we searched the SAG and CO-SAG assemblies
167 for plastidial contigs as evidence of a plastid genome. While there were some contigs that initially showed
168 similarities to reference plastid genomes, these were all rejected as bacterial (non-cyanobacterial) contamination
169 upon closer inspection. In contrast, mitochondrial contigs were readily identified in 26 of 43 SAGs (Table S5).
170 Although mitochondrial contigs remained fragmented in most SAGs, four complete or near-complete
171 mitochondrial genomes were recovered with coding content near-identical to the published mitochondrial
172 genome from picozoa MS5584-11¹⁹ (Fig S9). The ability to assemble complete mitochondrial genomes from the
173 SAGs suggests that the partial nature of the data does not specifically hinder organelle genome recovery if
174 present, at least in the case of mitochondria²⁰.

175 Second, we investigated the possibility that the plastid genome was lost while the organelle itself has been
176 retained—as is the case for *Rhodolphis*⁵. For this, we reconstructed phylogenetic trees for several essential
177 nuclear-encoded biochemical plastid pathways derived by endosymbiotic gene transfer (EGT) that were shown
178 to be at least partially retained even in cryptic plastids^{5,21,22}. These included genes involved in the biosynthesis of
179 isoprenoids (ispD,E,F,G,H, dxr, dxs), fatty acids (fabD,F,G,H,I,Z, ACC), heme (hemB,D,E,F,H,Y, ALAS), and
180 iron-sulfur clusters (sufB,C,D,E,S, NifU, iscA; see also Table S6). In all cases, the picozoan homologues
181 grouped either with bacteria—but not cyanobacteria, suggesting contamination—or the mitochondrial/nuclear
182 copies of host origin. Furthermore, none of the picozoan homologues contained predicted N-terminal plastid
183 transit peptides. We also searched for picozoan homologues of all additional proteins (n=62) that were predicted
184 to be targeted to the cryptic plastid in rhodelphids⁵. This search resulted in one protein (Arogenate
185 dehydrogenase, OG0000831) with picozoan homologues that were closely related to red algae and belonged to a
186 larger clade with host-derived plastid targeted plant sequences, but neither the picozoan nor the red algal
187 sequences displayed predicted transit peptides. Finally, to eliminate the possibility of missing sequences because
188 of errors during the assembly and gene prediction, we additionally searched the raw read sequences for the same
189 plastid-targeted or plastid transport machinery genes, which revealed no obvious candidates. In contrast, we
190 readily identified mitochondrial genes (e.g. homologues of the mitochondrial import machinery from the
191 TIM17/TIM22 family), which further strengthened our inference that the single-cell data are in principle
192 adequate to identify organellar components, when they are present.

193 The lack of cryptic plastids in diverse modern-day picozoans does not preclude photosynthetic ancestry if the
194 plastid was lost early in the evolution of the group. To assess this possibility, we searched more widely for
195 evidence of a cyanobacterial footprint on the nuclear genome that would rise above a background of horizontal
196 gene transfers for proteins functioning in cellular compartments other than the plastids. The presence of a
197 significant number of such proteins may be evidence for a plastid-bearing ancestor. We clustered proteins from

198 419 genomes, including all major eukaryotic groups as well as a selection of bacteria into orthologous groups
199 (OGs) (Table S7). We built phylogenies for OGs that contained at least cyanobacterial and algal sequences, as
200 well as a sequence from one of 33 focal taxa, including Picozoa, a range of photosynthetic taxa, but also non-
201 photosynthetic plastid-containing, and plastid-lacking taxa to be used as controls. Putative gene transfers from
202 cyanobacteria (EGT) were identified as a group of plastid-bearing eukaryotes that included sequences from the
203 focal taxa and branched sister to a clade of cyanobacteria. We allowed up to 10% of sequences from groups with
204 no plastid ancestry. This approach identified 16 putative EGTs for Picozoa where at least 2 different SAGs/CO-
205 SAGs grouped together, compared to between 89–313 EGTs for photosynthetic species, and up to 59 EGTs for
206 species with non-photosynthetic plastids (Fig 3a). At the other end of the spectrum for species with non-
207 photosynthetic plastids, we observed that the number of inferred cyanobacterial genes for e.g. rhodelphids (14)
208 or *Paraphysomonas* (12) was comparable to Picozoa (16) or other, plastid-lacking taxa such as *Telonema* (15) or
209 *Goniomonas* (18). In order to differentiate these putative endosymbiotic transfers from a background of bacterial
210 transfers (or bacterial contamination), we next attempted to normalise the EGT signal by estimating an extended
211 bacterial signal (indicative of putative HGT: horizontal gene transfers) using the same tree sorting procedure
212 (Fig S10). When comparing the number of inferred EGT with that of inferred HGT, we found a marked
213 difference between plastid-containing (including non-photosynthetic) and plastid-lacking lineages. While all
214 plastid-containing taxa—with the notable exception of *Rhodelphis*—showed a ratio of EGT to HGT above 1, all
215 species without plastid ancestry and *Hematodinium*, one of the few taxa with reported plastid loss, as well as
216 *Rhodelphis* and Picozoa showed a much higher number of inferred HGT than EGT.

217

218 Discussion

219 The 17 SAGs and CO-SAGs of Picozoa obtained in this study provide robust data for phylogenomic analyses of
220 this important phylum of eukaryotes. With this data, we are able to firmly place Picozoa within the supergroup
221 Archaeplastida, most likely as a sister lineage to red algae and rhodelphids. Archaeplastids contain all known
222 lineages with primary plastids (with the exception of *Paulinella*), which are widely viewed to be derived from a
223 single primary endosymbiosis with a cyanobacterium. This notion of a common origin of primary plastids is
224 supported by cellular and genomic data (see^{23,24} and references therein for review), as well as plastid
225 phylogenetics^{25,26}. The phylogenetic support for Archaeplastida based on host (nuclear) data has been less
226 certain^{7,8,27}, but our analysis is consistent with recent reports that have also recovered a monophyletic origin—
227 here including Picozoa—when using gene and taxon-rich phylogenomic datasets^{1,9,10}. This position has
228 important implications for our understanding of plastid origins because, in contrast to all other archaeplastids
229 known to date, our results indicate that Picozoa lack plastids and plastid-associated EGTs. The lack of plastid in
230 Picozoa was also inferred based on smaller initial SAG data¹¹ as well as ultrastructural observation of *P.*

231 *judraskeda*¹⁴. Two main possible hypotheses exist to explain the lack of plastids in Picozoa: that this group was
232 never photosynthetic, or complete plastid loss occurred early in their evolution.

233 To suggest that Picozoa was never photosynthetic requires that the current distribution of primary plastids is due
234 to multiple independent endosymbioses, specifically that red algae arose from a separate primary endosymbiosis
235 from that leading to green algae and glaucophytes. This scenario would have involved the endosymbioses of
236 closely related cyanobacterial lineages in closely related hosts to explain the many similarities between primary
237 plastids²⁴. Although this may sound unlikely, there is accumulating evidence that similar plastids were derived
238 independently from similar endosymbionts in closely related hosts in dinoflagellates with tertiary plastids^{28–30},
239 and has been argued before for primary plastids^{31–34}. However, the current bulk of cell and molecular evidence
240 suggests that multiple independent origins of primary plastids are unlikely, including several features of plastid
241 biology that are not present in cyanobacteria (e.g., protein targeting systems, light-harvesting complex proteins,
242 or plastid genome architecture)^{23,24,35}. A related explanation could involve a secondary endosymbiosis where the
243 plastid in red algae, for example, was secondarily acquired from a green alga³⁶. This latter scenario would be
244 made unlikely by the identification of host-derived plastid components shared between all archaeplastid lineages.

245 The second hypothesis implies that a common ancestor of Picozoa entirely lost its primary plastid. The
246 possibility of plastid loss in a free-living lineage like Picozoa would be unprecedented because to date, the only
247 known unambiguous cases of total plastid loss all come from parasitic lineages (all in myxozoan alveolates: in
248 *Cryptosporidium*³⁷, certain gregarines^{22,38}, and the dinoflagellate *Hematodinium*³⁹). To evaluate this possibility,
249 we searched our data for a cyanobacterial footprint in the nuclear genome that would result from an ancestral
250 endosymbiosis. The transfer of genes from endosymbiont to host nucleus via EGT, and the targeting of the
251 product of some or all of these genes back to the plastids, are recognised as a hallmark of organelle
252 integration^{40,41}. EGT has occurred in all algae, although its impact on nuclear genomes can vary and the
253 inference of EGT versus other horizontally acquired genes (HGT) can be difficult to decipher for ancient
254 endosymbioses^{42–46}. Our analysis of the normalised cyanobacterial signal in Picozoa, which we used as a proxy
255 for quantifying EGT, provides no clear evidence for the existence of a plastid-bearing ancestor. However, it
256 should be noted that evaluating the possibility of plastid loss in groups where a photosynthetic ancestry is not
257 confirmed—such as Picozoa—is complicated because there is no baseline for the surviving footprint of
258 endosymbiosis following plastid loss. Notably, we found no significant difference in the number of inferred
259 EGTs in Picozoa compared to lineages with demonstrated plastid loss (e.g. *Hematodinium* with 10 inferred
260 EGT), lineages with non-photosynthetic plastids (e.g. *Rhodolphis*: 14 inferred EGT), or with no photosynthetic
261 ancestry (e.g. *Telonema*: 15 inferred EGT).

262 The lack of a genomic baseline to assess plastid loss in Picozoa is further complicated by limitations of our data
263 and methods. The partial nature of eukaryotic SAGs makes it possible that EGTs are absent from our data, even
264 with >90% of inferred genomic completeness. Additionally, the possibility exists that the number of EGT might

265 have always been low during the evolution of the group, even if a plastid was once present. Recent
266 endosymbioses where EGT can be pinpointed with precision showed a relatively low frequency. For example,
267 they represent at most a few percent of the chromatophore proteome in *Paulinella*⁴⁷, or as few as 9 genes in
268 tertiary endosymbiosis in dinoflagellates⁴⁸. Thus, it is possible that the much higher number of EGT inferred in
269 red algae (e.g. 168 in *Galdieria*) occurred after the divergence of Picozoa, and that Picozoa quickly lost its
270 plastid before more EGT occurred. An observation that supports this hypothesis is the low number of putative
271 EGTs found in *Rhodolphis* (14), suggesting that the bulk of endosymbiotic transfers in red algae may have
272 happen after their divergence from rhodelphids.

273

274 **Conclusion**

275 In this study, we used single-cell genomics to demonstrate that Picozoa are a plastid lacking major lineage of
276 archaeplastids. This is the first example of an archaeplastid lineage without plastids, which can be interpreted as
277 either plastid loss, or evidence of independent endosymbiosis in the ancestor of red algae and rhodelphids. Under
278 the most widely accepted scenario of a single plastid origin in Archaeplastida, Picozoa would represent the first
279 known case of plastid loss in this group, but also more generally in any free-living species. In order to
280 discriminate plastid loss from multiple plastid gains in the early archaeplastid evolution, and more generally
281 during the evolution of secondary or tertiary plastids, a better understanding of the early steps of plastid
282 integration is required. In the recently evolved primary plastid-like chromatophores of *Paulinella*, the transfer of
283 endosymbiotic genes at the onset of the integration was shown to be minimal⁴. Similar examples of integrated
284 plastid endosymbionts but with apparently very few EGTs are known in dinoflagellates^{48,49}. Therefore, new
285 important clues to decipher the origin of plastids will likely come from a better understanding of the role of the
286 host in driving these endosymbioses, and crucially the establishment of a more complete framework for
287 archaeplastid evolution with the search and characterization of novel diversity of lineages without plastids. The
288 fact that this lineage has never been successfully maintained in culture, with just one study achieving transient
289 culture¹⁴, might indicate a lifestyle involving close association with other organisms (such as symbiosis) and
290 further underscores the enigma of picozoan biology, the lack of information on which hinders our interpretation
291 of their evolution.

292

293 **Material and Methods**

294 *Cell Isolation, Identification, and genome amplification*

295 *Baltic Sea.* Surface (depth: up to 2 m) marine water was collected from the Linnaeus microbial Observatory
296 (LMO) in the Baltic Sea located at 56°N 55.85' and 17°E 03.64' on two occasions: 2 May 2018 (6.1°C and 6.8
297 ppt salinity) and 3 April 2018 (2.4°C and 6.7 ppt salinity). The samples were transported to the laboratory and

298 filter-fractionized (see Table 1 for detailed sample information). The size fractions larger than 2 μm were
299 discarded whereas the fraction collected on 0.2 μm filters was resuspended in 2 mL of the filtrate. The obtained
300 samples were used for fluorescence-activated cell sorting (FACS). Aliquots of 4 μL of 1 mM Mitotracker Green
301 FM (ThermoFisher) stock solution were added to the samples and were kept in the dark at 15°C for 15–20
302 minutes. Then the cells were sorted into empty 96-well plates using MoFlo Astrios EQ cell sorter (Beckman
303 Coulter). Gates were set mainly based on Mitotracker intensity and the dye was detected by a 488 nm and 640
304 nm laser for excitation, 100 μm nozzle, sheath pressure of 25 psi and 0.1 μm sterile filtered 1 x PBS as sheath
305 fluid. The region with the highest green fluorescence and forward scatter contained the target group and was
306 thereafter used alongside with exclusion of red autofluorescence.

307 The SAGs were generated in each well with REPLI-g® Single Cell kit (Qiagen) following the manufacturer's
308 recommendations but scaled down to 5 μl reactions. Since the cells were sorted in dry plates, 400 nl of 1xPBS
309 was added prior to 300 nl of lysis buffer D2 for 10 min at 65°C and 10 min on ice, followed by 300 nl stop
310 solution. The PBS, reagent D2, stop solution, water, and reagent tubes were UV-treated at 2 Joules before use. A
311 final concentration of 0.5 μM SYTO 13 (Invitrogen) was added to the MDA mastermix. The reaction was run at
312 30°C for 6 h followed by inactivation at 65°C for 5 min and was monitored by detection of SYTO13
313 fluorescence every 15 minutes using a FLUOstar® Omega plate reader (BMG Labtech, Germany). The single
314 amplified genome (SAG) DNA was stored at -20°C until further PCR screening. The obtained products were
315 PCR-screened using Pico-PCR approach, as described in¹⁴ and the wells showing signal for Picozoa were
316 selected for sequencing.

317

318 *Eastern North Pacific*. Seawater was collected and sorted using a BD InFlux Fluorescently Activated Cell Sorter
319 (FACS) on three independent cruises in the eastern North Pacific. The instrument was equipped with a 100 mW
320 488 nm laser and a 100 mW 355 nm laser and run using sterile nuclease-free 1× PBS as sheath fluid. The
321 stations where sorting occurred were located at 36.748°N, 122.013°W (Station M1; 20 m, 2 April 2014 and 10
322 m, 5 May 2014); 36.695°N, 122.357°W (Station M2, 10 m, 5 May 2014); and 36.126°N, 123.49°W (Station 67-
323 70, 20 m 15 October 2013). Water was collected using Niskin bottles mounted on a CTD rosette. Prior to sorting
324 samples were concentrated by gravity over a 0.8 μm Supor filter. Two different stains were used: LysoSensor (2
325 April 2014, M1) and LysoTracker (5 May 2014, M1; 15 October 2013, 67-70), or both together (5 May 2014,
326 M2). Selection of eukaryotic cells stained with LysoTracker Green DND-26 (Life Technologies; final
327 concentration, 25 nM) was based on scatter parameters, positive green fluorescence (520/35 nm bandpass), as
328 compared to unstained samples, and exclusion of known phytoplankton populations, as discriminated by their
329 forward angle light scatter and red (chlorophyll-derived) autofluorescence (i.e., 692/40 nm bandpass) under 488
330 nm excitation, similar to methods in⁵⁰. Likewise, selection of cells stained with LysoSensor Blue DND-167 (Life
331 Technologies; final concentration, 1 μM), a ratiometric probe sensitive to intracellular *pH* levels, e.g. in
332 lysosomes, was based on scatter parameters, positive blue fluorescence (435/40 nm bandpass), as compared to
333 unstained samples, and exclusion of known phytoplankton populations, as discriminated by their forward angle

334 light scatter and red (chlorophyll-derived) autofluorescence (i.e., 692/40 nm bandpass filter) under 355 nm
335 excitation. For sorts using both stains, all of the above criteria, and excitation with both lasers (with emissions
336 collected through different pinholes and filter sets), were applied to select cells. Before each sort was initiated,
337 the respective plate was illuminated with UV irradiation for 2 min. Cells were sorted into 96- or 384-well plates
338 using the Single-Cell sorting mode from the BD FACS Software v1.0.0.650. A subset of wells was left empty or
339 received 20 cells for negative and positive controls, respectively. After sorting, the plates were covered with
340 sterile, nuclease free foil and frozen at -80 °C immediately after completion.

341 Whole genome amplification of individual sorted cells followed methods outlined in⁵⁰. For initial screening,
342 18S rRNA gene amplicons were amplified from each well using the Illumina adapted TAREuk454FWD1 and
343 TAREukREV3 primers targeting the V4 hypervariable region. PCR reactions contained 10 ng of template
344 DNA and 1X 5PRIME HotMasterMix (Quanta Biosciences) as well as 0.4 mg ml⁻¹ BSA (NEB) and 0.4 μM of
345 each primer. PCR reactions entailed: 94 °C for 3 min; and 30 cycles at 94 °C for 45 sec, 50 °C for 60 sec and
346 72 °C for 90 sec; with a final extension at 72 °C for 10 min. Triplicate reactions per cell were pooled prior to
347 Paired-end (PE) library sequencing (2 × 300 bp) and the resulting 18S V4 rRNA gene amplicons were
348 trimmed at Phred quality (Q) of 25 using a 10 bp running window using Sickle 1.33
349 (<https://github.com/najoshi/sickle>). Paired-end reads were merged using USEARCH v.9.0.2132 when reads
350 had a ≥40 bp overlap with max 5% mismatch. Merged reads were filtered to remove reads with maximum
351 error rate >0.001 or <200 bp length. Sequences with exact match to both primers were retained, primer
352 sequences were trimmed using Cutadapt v.1.13⁵¹, and the remaining sequences were *de novo* clustered at 99%
353 sequence similarity by UCLUST forming operational taxonomic units (OTUs). Each of the cells further
354 sequenced had a single abundant OTU that was taxonomically identified using BLASTn in GenBank's nr
355 database.

356

357 *Sequencing*

358 Sequencing libraries were prepared from 100 ng DNA using the TruSeq Nano DNA sample preparation kit (cat#
359 20015964/5, Illumina Inc.) targeting an insert size of 350bp. For six samples, less than 100ng was used (between
360 87 ng-97 ng). The library preparations were performed by SNP&SEQ Technology Platform at Uppsala
361 University according to the manufacturers' instructions. All samples were then multiplexed on one lane of an
362 Illumina HiSeqX instrument with 150 cycles paired-end sequencing using the v2.5 sequencing chemistry,
363 producing between 10,000 and 30,000,000 read pairs.

364

365 *Genome Assembly and 18S rRNA gene analysis*

366 The 43 Illumina datasets were trimmed using *Trim Galore* v0.6.1
367 (https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/) with default parameter and assembled into
368 genomic contigs with SPAdes v3.13.0⁵² in single-cell mode (--sc --careful -k 21,33,55,99). Open reading frames
369 (ORFs) were identified and translated using Prodigal v2.6.3 in ‘anonymous’ mode⁵³ and rRNA genes were
370 predicted using barrnap v0.9 (<https://github.com/tseemann/barrnap>) for eukaryotes. All 18S rRNA gene
371 sequences were, together with available reference sequences from the protist ribosomal reference database
372 (PR2), aligned with MAFFT E-INS-i v7.429⁵⁴ and trimmed with trimal⁵⁵ (gap threshold 0.01%). After
373 performing a modeltest using ModelFinder⁵⁶ (best model: GTR+R6+F), a phylogenetic tree was reconstructed in
374 IQ-TREE v2.1.1⁵⁷ with 1000 ultrafast bootstrap replicates (see Fig S11 for a tree with extended taxon sampling).
375 Additionally, we estimated the average nucleotide identity (ANI) for all pairs of SAGs using fastANI v1.2⁵⁸ (Fig
376 S12). Based on the 18S rRNA gene tree and the ANI value, groups of closely related SAGs with almost identical
377 18S rRNA gene sequences (sequence similarity above 99%) were identified for co-assembly. Co-assemblies
378 were generated in the same way as described above for single assemblies, pooling sequencing libraries from
379 closely related single cells. ORFs and rRNA genes were similarly extracted from the co-assemblies. The
380 completeness of the SAGs and CO-SAGs was then assessed using BUSCO v4.1.3¹⁵ with 255 markers for
381 eukaryotes (Fig S13) as well as using the 320 marker phylogenomic dataset as described below. General genome
382 characteristics were computed with QUAST v5.0.2⁵⁹. Alignments were reconstructed for the 18S rRNA genes
383 from the co-assemblies and those SAGs not included in any CO-SAG together with PR2 references for cryptists
384 and katablepharids (the closest groups to Picozoa in 18S rRNA gene phylogenies) in the same way as described
385 above. The tree was reconstructed using GTR+R4+F after model selection and support was assessed with 100
386 non-parametric bootstraps. The six CO-SAGs and the 11 individual SAGs were used in all subsequent analyses.

387 For each of these 17 assemblies we estimated the amount of prokaryotic/viral contamination by comparing the
388 predicted proteins against the NCBI nr database using DIAMOND in blastp mode⁶⁰. If at least 60% of all
389 proteins from a contig produced significant hits only to sequences annotated as prokaryotic or viral, we
390 considered that contig to be a putative contamination. In general only a small fraction of each assembly was
391 found to be such a contamination (Fig S14).

392

393 *Phylogenomics*

394 Existing untrimmed alignments for 320 genes and 763 taxa from¹ were used to create HMM profiles in HMMER
395 v3.2.1⁶¹, which were then used to identify homologous sequences in the protein sequences predicted from the
396 Picozoa assemblies (or co-assemblies) as well as in 20 additional, recently sequenced eukaryotic genomes and
397 transcriptomes (Table S2). Each single gene dataset was filtered using PREQUAL v1.02⁶² to remove non-
398 homologous residues prior to alignment, aligned using MAFFT E-INS-i, and filtered with Divvier -partial v1.0⁶³.

399 Alignments were then used to reconstruct gene trees with IQ-TREE (-mset LG, LG4X; 1000 ultrafast bootstraps
400 with the BNNI optimization). All trees were manually scrutinized to identify contamination and paralogs. These
401 steps were repeated at least two times, until no further contaminations or paralogs could be detected. We
402 excluded three genes that showed ambiguous groupings of Picozoa or rhodelphids in different parts of the trees.
403 From this full dataset of 317 genes and 794 taxa, we created a concatenated supermatrix alignment using the
404 cleaned alignments described above. This supermatrix was used to reconstruct a tree in IQ-TREE with the model
405 LG+G+F and ultrafast bootstraps (1000 UFBoots) estimation with the BNNI improvement.

406 We then prepared a reduced dataset with a more focused taxon sampling of 67 taxa, covering all major
407 eukaryotic lineages but focussing on the groups for which an affiliation to Picozoa had been reported previously.
408 For this dataset, closely related species were merged into OTUs in some cases in order to decrease the amount of
409 missing data per taxon (Table S3). The 317 single gene datasets were re-aligned using MAFFT E-INS-i, filtered
410 using both Divvier -partial and BMGE (-g 0.2 -b 10 -m BLOSUM75) and concatenated into two supermatrices.
411 Model selection of mixture models was performed using ModelFinder⁵⁶ for both datasets, and in both cases
412 LG+C60+G+F was selected as the best-fitting model. Trees for both datasets were reconstructed using the
413 Posterior Mean Site Frequency (PMSF)⁶⁴ approximation of this mixture model in IQ-TREE and support was
414 assessed with 100 non-parametric bootstraps (see Fig S15 for the Divvier derived tree)

415 In addition, we reconstructed a phylogenetic tree using the supermatrix alignment based on BMGE trimming in
416 PhyloBayes MPI v1.8⁶⁵ using the CAT+GTR+G model. We ran three independent chains for 3600 cycles, with
417 the initial 1500 cycles being removed as burnin from each chain. We then generated a consensus tree using the
418 bpcmp program of PhyloBayes. Partial convergence was achieved between chains 1 and 2 with a maxdiff value
419 of 0.26 (Fig S16). The third chain differed only in the position of haptists and *Ancoracysta twisti*, but not in the
420 relationships within Archaeplastida and the position of Picozoa (Fig S17).

421 In order to test the robustness of our results we additionally performed a fast-site removal analysis⁶⁶, iteratively
422 removing the 5000 fastest evolving sites (up to a total of 55000 removed sites). For each of these 11 alignments,
423 we reconstructed an ML tree using the model LG+C60+G+F in IQ-TREE with ultrafast bootstraps (1000
424 UFBoots) and evaluated the support for the branching of Picozoa with rhodelphids and red algae as well as for
425 other groupings (Fig S4). We also performed trimming of the 25% and 50% most heterogeneous sites based on
426 the χ^2 metric⁶⁷ and performed tree reconstruction using the same model as above (Fig S5, Fig S7). We also
427 prepared a supermatrix alignment (BMGE trimmed) from 224 genes with at least two Picozoa sequences in the
428 final dataset and performed similar tree reconstruction (model LG+C60+G+F in IQ-TREE with 1000 ultrafast
429 bootstraps, Fig S8).

430 Furthermore, we performed a supertree-based phylogenetic reconstruction using ASTRAL-III v5.7.3⁶⁸. We
431 reconstructed gene trees for each of the 317 alignments of the 67-taxon dataset using IQ-TREE (-m TEST -mset

432 LG -mrate G,R4 -madd LG4X,LG4X+F,LG4M,LG4M+F, using 1000 ultrafast bootstraps) and performed multi-
433 locus bootstrapping based on the bootstrap replicates (option -b in ASTRAL-III) (Fig S6).

434 Finally, we performed an approximately unbiased (AU) test in IQ-TREE of 15 topologies (see Table S4),
435 including previously recovered positions of Picozoa (as sister to red algae, cryptists, telonemids, Archaeplastida
436 etc.).

437

438 *Mitochondrial contig identification and annotation*

439

440 Using the published picozoan mitochondrial genome (Picozoa sp. MS584-11: MG202007.1 from¹⁹), BLAST
441 searches were performed on a dedicated sequenceServer⁶⁹ to identify mitochondrial contigs in the 43 picozoan
442 SAGs. Putative mitochondrial contigs were annotated using the MFannot server
443 (<https://megasun.bch.umontreal.ca/cgi-bin/mfannot/mfannotInterface.pl>). All contigs with predicted
444 mitochondrial genes or whose top hits in the NCBI nr database was the published picozoan mitochondrial
445 genome (MG202007.1) were considered to be *bona fide* mitochondrial contigs and retained (Supplementary
446 materials). Manual annotation was conducted as needed.

447

448 *Plastid Genes & EGT*

449 GetOrganelle v1.7.1⁷⁰ was used to identify organellar genomes. We searched the assemblies for putative plastid
450 contigs with the subcommand 'get_organelle_from_assembly.py -F embplant_pt,other_pt', while we attempted
451 to assemble such a genome directly using the command 'get_organelle_from_reads.py -R 30 -k 21,45,65,85,105
452 -F embplant_pt,other_pt'. We additionally searched the predicted proteins against available plastid protein
453 sequences from ncbi using DIAMOND v2.0.6⁶⁰ in blastp mode (--more-sensitive). Contigs that were identified
454 as putatively coming from a plastid genome were then checked manually by doing BLAST searches against NT,
455 and contigs that showed similarity only to bacterial genomes or to the picozoa mitochondrial assembly
456 MG202007.1 were rejected.

457 To search for known plastid pathways, we prepared Hidden Markov model (HMM) profiles for 32 gene
458 alignments that were shown to be retained in lineages with non-photosynthetic plastids and included a wide
459 diversity of plastid-bearing eukaryotes following a similar approach as in²². Using these profiles, we identified
460 homologues in the Picozoa SAGs, and aligned them together with the initial sequences used to create the profiles
461 using MAFFT E-INS-i. We trimmed the alignments using trimAl v1.4.rev15 '-gt 0.05' and reconstructed
462 phylogenetic trees using IQ-TREE (-m LG4X; 1000 ultrafast bootstraps with the BNNI optimization) from these
463 alignments. We then manually inspected the trees to assess whether picozoan sequences grouped with known
464 plastid-bearing lineages. We additionally used the sequences from these core plastid genes to search the raw
465 sequencing reads for any signs of homologues that could have been missed in the assemblies. We used the tool

466 PhyloMagnet⁷¹ to recruit reads and perform gene-centric assembly of these genes⁷². The assembled genes were
467 then compared to the NR database using DIAMOND in blastp mode (--more-sensitive --top 10).

468 To identify putative EGT, we prepared orthologous clusters for 419 species (128 bacteria and 291 eukaryotes)
469 with a focus on plastid-bearing eukaryotes and cyanobacteria, but also including other eukaryotes and bacteria,
470 using OrthoFinder v2.4.0⁷³. For Picozoa and a selection of 32 photosynthetic or heterotrophic lineages (Table
471 S8), we inferred trees for 2626 clusters that contained the species under consideration, at least one cyanobacterial
472 sequence, and at least one archaeplastid sequence of red algae, green algae or plants. Alignments for these
473 clusters were generated with MAFFT E-INS-i, filtered using trimAl '-gt 0.01' and phylogenetic trees were
474 reconstructed using IQ-TREE (-m LG4X; 1000 ultrafast bootstraps with the BNNI optimization). We then
475 identified trees where the target species grouped with other plastid-bearing lineages (allowing up to 10% non-
476 plastid sequences) and sister at least two cyanobacterial sequences. For Picozoa, we added the condition that
477 sequences from at least two SAG/COSAG assemblies must be monophyletic. For species with no known plastid
478 ancestry such as *Rattus* or *Phytophthora*, putative EGTs can be interpreted as false positives due to
479 contamination, poor tree resolution or other mechanisms, since we expect no EGTs from cyanobacteria to be
480 present at all in these species. This rough estimate of the expected false positive rate for this approach can give
481 us a baseline of false positives that can be expected for picozoa as well.

482 To put the number of putative EGTs into relation to the overall amount of gene transfers, we applied a very
483 similar approach to the one described above for detecting putative HGT events. We prepared additional trees (in
484 the same way as described for the detection of EGTs) for clusters that contained the taxon of interest and non-
485 cyanobacterial bacteria and identified clades of the taxon under consideration (including a larger taxonomic
486 group, e.g. Streptophyta for *Arabidopsis* or Metazoa for *Rattus*) that branched sister to a bacterial clade.

487

488 *Distribution of Picozoa in Tara Oceans*

489 We screened available OTUs that were obtained from V9 18S rRNA gene eukaryotic amplicon data generated
490 by *Tara Oceans*¹⁷ for sequences related to Picozoa. Using the V9 region of the 18S rRNA gene sequences from
491 the 17 Picozoa assemblies as well as from the picozoan PR2 references used to reconstruct the 18S rRNA gene
492 tree described above, we applied VSEARCH v2.15.1⁷⁴ (--usearch_global -iddef 1 --id 0.90) to find all OTUs
493 with at least 90 % similar V9 regions to any of these reference picozoan sequences. Using the relative abundance
494 information available for each *Tara Oceans* sampling location, we then computed the sum for all identified
495 Picozoa OTUs per station and plotted the relative abundance on a world map.

496

497 **Acknowledgments**

498 This work was supported by a grant from Science for Life Laboratory available to FB and a scholarship from
499 Carl Tryggers Stiftelse to VZ (PI: FB). TJGE thanks the European Research Council (ERC consolidator grant
500 817834); the Dutch Research Council (NWO-VICI grant VI.C.192.016); Moore–Simons Project on the Origin of
501 the Eukaryotic Cell (Simons Foundation 735925LPI, <https://doi.org/10.46714/735925LPI>); and the Marie
502 Skłodowska-Curie ITN project SINGEK (H2020-MSCA-ITN-2015-675752) which provided funding for MES.
503 PJK and VM were funded by an Investigator Grant from the Gordon and Betty Moore Foundation
504 (<https://doi.org/10.37807/GBMF9201>). The Pacific Ocean work was supported by GBMF3788 to AZW.
505 Sampling at the LMO station in the Baltic Sea was carried out by support from the Swedish Research Council
506 VR and the marine strategic research program EcoChange to JP. Sequencing was performed by the SNP&SEQ
507 Technology Platform in Uppsala, part of the National Genomics Infrastructure (NGI) Sweden and Science for
508 Life Laboratory. The SNP&SEQ Platform is also supported by the Swedish Research Council and the Knut and
509 Alice Wallenberg Foundation. Cell sorting and whole genome amplification was performed at the Microbial
510 Single Cell Genomics Facility (MSCG) at SciLifeLab. Computations were performed on resources provided by
511 the Swedish National Infrastructure for Computing (SNIC) at Uppsala Multidisciplinary Center for Advanced
512 Computational Science (UPPMAX) under Projects SNIC 2019/3-305, SNIC 2020/15-58, SNIC 2021/5-50,
513 Uppstore2018069. Finally, we thank Eunsoo Kim and Sally D. Warring for sharing peptide models from
514 *Palpitomonas bilix* and *Roombia truncata*.

515

516 **Author contributions**

517 FB and JGW conceived the study. For the Baltic samples, VVZ and JP performed sampling and cell sorting;
518 VVZ and JFHS performed genome amplification & sequencing preparation. For the Pacific samples, CP, SW
519 and AZW conceived, developed and implemented sort protocols; performed sampling, cell sorting and
520 sequencing; CP and AZW performed initial sequence analyses and phylogenetics. MES under the supervision of
521 TJGE and FB performed assembly of SAGs and Co-SAGs, phylogenomic analyses, and searched for plastid
522 evidence and gene transfers with the help of VM and PJK. RPS and JGW assembled mitochondrial genomes.
523 MES, FB, and JGW drafted the manuscript. TJGE, PKJ, AZW, JP, VVZ and JFHS contributed edits to the
524 manuscript. All authors read and approved the final version.

525

526 **Code & data availability**

527 All custom scripts used in this study are available at <https://github.com/maxemil/picozoa-scripts> under a MIT
528 license. All data used for the analyses as well as results files such as contigs and single gene trees are available at
529 figshare ([10.6084/m9.figshare.c.5388176](https://doi.org/10.6084/m9.figshare.c.5388176)). A sequenceServer BLAST server was set up for the SAG assemblies:

530 <http://evocellbio.com/SAGdb/burki/>. Raw sequencing reads were deposited in the Sequence Read Archive
531 (SRA) at NCBI under accession PRJNA747736 and will be available upon acceptance.

532

533 **References**

534 1. Strassert, J. F. H., Irisarri, I., Williams, T. A. & Burki, F. A molecular timescale for eukaryote evolution with
535 implications for the origin of red algal-derived plastids. *Nat Commun* 12, 1879 (2021).

536 2. Burki, F., Roger, A. J., Brown, M. W. & Simpson, A. G. B. The New Tree of Eukaryotes. *Trends Ecol Evol*
537 35, 43–55 (2019).

538 3. Marin, B., Nowack, E. C. M. & Melkonian, M. A Plastid in the Making: Evidence for a Second Primary
539 Endosymbiosis. *Protist* 156, 425–432 (2005).

540 4. Nowack, E. C. M. & Weber, A. P. M. Genomics-Informed Insights into Endosymbiotic Organelle Evolution
541 in Photosynthetic Eukaryotes. *Annu Rev Plant Biol* 69, 1–34 (2018).

542 5. Gawryluk, R. M. R. *et al.* Non-photosynthetic predators are sister to red algae. *Nature* 572, 240–243 (2019).

543 6. Li, L. *et al.* The genome of *Prasinoderma coloniale* unveils the existence of a third phylum within green
544 plants. *Nat Ecol Evol* 4, 1220–1231 (2020).

545 7. Burki, F. *et al.* Untangling the early diversification of eukaryotes: a phylogenomic study of the evolutionary
546 origins of Centrohelida, Haptophyta and Cryptista. *Proc Royal Soc B Biological Sci* 283, 20152802 (2016).

547 8. Strassert, J. F. H., Jamy, M., Mylnikov, A. P., Tikhonenkov, D. V. & Burki, F. New Phylogenomic Analysis
548 of the Enigmatic Phylum Telonemia Further Resolves the Eukaryote Tree of Life. *Mol Biol Evol* 36, 757–765
549 (2019).

550 9. Lax, G. *et al.* Hemimastigophora is a novel supra-kingdom-level lineage of eukaryotes. *Nature* 564, 410–414
551 (2018).

552 10. Irisarri, I., Strassert, J. F. H. & Burki, F. Phylogenomic Insights into the Origin of Primary Plastids.
553 *Systematic Biol* syab036- (2021) doi:10.1093/sysbio/syab036.

554 11. Yoon, H. S. *et al.* Single-Cell Genomics Reveals Organismal Interactions in Uncultivated Marine Protists.
555 *Science* 332, 714–717 (2011).

556 12. Not, F. *et al.* Picobiliphytes: A Marine Picoplanktonic Algal Group with Unknown Affinities to Other
557 Eukaryotes. *Science* 315, 253–255 (2007).

558 13. Cuvelier, M. L. *et al.* Widespread distribution of a unique marine protistan lineage. *Environ Microbiol* 10,
559 1621–1634 (2008).

560 14. Seenivasan, R., Sausen, N., Medlin, L. K. & Melkonian, M. *Picomonas judraskeda* Gen. Et Sp. Nov.: The
561 First Identified Member of the Picozoa Phylum Nov., a Widespread Group of Picoeukaryotes, Formerly Known
562 as ‘Picobiliphytes.’ *Plos One* 8, e59565 (2013).

- 563 15. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing
564 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31, 3210–3212
565 (2015).
- 566 16. Moreira, D. & López-García, P. The rise and fall of Picobiliphytes: How assumed autotrophs turned out to be
567 heterotrophs. *Bioessays* 36, 468–474 (2014).
- 568 17. Vargas, C. de *et al.* Eukaryotic plankton diversity in the sunlit ocean. *Science* 348, 1261605 (2015).
- 569 18. Kim, E. & Graham, L. E. EEF2 Analysis Challenges the Monophyly of Archaeplastida and Chromalveolata.
570 *Plos One* 3, e2621 (2008).
- 571 19. Janouškovec, J. *et al.* A New Lineage of Eukaryotes Illuminates Early Mitochondrial Genome Reduction.
572 *Curr Biol* 27, 3717–3724.e5 (2017).
- 573 20. Wideman, J. G. *et al.* Unexpected mitochondrial genome diversity revealed by targeted single-cell genomics
574 of heterotrophic flagellated protists. *Nat Microbiol* 5, 154–165 (2020).
- 575 21. Dorrell, R. G. *et al.* Principles of plastid reductive evolution illuminated by nonphotosynthetic chrysophytes.
576 *Proc National Acad Sci* 201819976 (2019) doi:10.1073/pnas.1819976116.
- 577 22. Mathur, V. *et al.* Multiple Independent Origins of Apicomplexan-Like Parasites. *Curr Biol* 29, 2936–2941.e5
578 (2019).
- 579 23. Reyes-Prieto, A., Weber, A. P. M. & Bhattacharya, D. The Origin and Establishment of the Plastid in Algae
580 and Plants. *Annu Rev Genet* 41, 147–168 (2007).
- 581 24. Gould, S. B., Waller, R. F. & McFadden, G. I. Plastid Evolution. *Annu Rev Plant Biol* 59, 491–517 (2008).
- 582 25. Shih, P. M. *et al.* Improving the coverage of the cyanobacterial phylum using diversity-driven genome
583 sequencing. *Proc National Acad Sci* 110, 1053–1058 (2013).
- 584 26. Ponce-Toledo, R. I. *et al.* An Early-Branching Freshwater Cyanobacterium at the Origin of Plastids. *Curr*
585 *Biol* 27, 386–391 (2017).
- 586 27. Yabuki, A. *et al.* *Palpitomonas bilix* represents a basal cryptist lineage: insight into the character evolution in
587 Cryptista. *Sci Rep-uk* 4, 4641 (2014).
- 588 28. Hehenberger, E., Gast, R. J. & Keeling, P. J. A kleptoplastidic dinoflagellate and the tipping point between
589 transient and fully integrated plastid endosymbiosis. *Proc National Acad Sci* 116, 17934–17942 (2019).
- 590 29. Sarai, C. *et al.* Dinoflagellates with relic endosymbiont nuclei as models for elucidating organellogenesis.
591 *Proc National Acad Sci* 117, 5364–5375 (2020).
- 592 30. Yamada, N., Sakai, H., Onuma, R., Kroth, P. G. & Horiguchi, T. Five Non-motile Dinotom Dinoflagellates
593 of the Genus *Dinotrix*. *Front Plant Sci* 11, 591050 (2020).
- 594 31. Stiller, J. W., Reel, D. C. & Johnson, J. C. A single origin of plastids revisited: convergent evolution in
595 organellar genome content. *J Phycol* 39, 95–105 (2003).

- 596 32. Larkum, A. W. D., Lockhart, P. J. & Howe, C. J. Shopping for plastids. *Trends Plant Sci* 12, 189–195
597 (2007).
- 598 33. Howe, C. J., Barbrook, A. C., Nisbet, R. E. R., Lockhart, P. J. & Larkum, A. W. D. The origin of plastids.
599 *Philosophical Transactions Royal Soc B Biological Sci* 363, 2675–2685 (2008).
- 600 34. Stiller, J. W. Toward an empirical framework for interpreting plastid evolution. *J Phycol* 50, 462–471
601 (2014).
- 602 35. Bhattacharya, D., Archibald, J. M., Weber, A. P. M. & Reyes-Prieto, A. How do endosymbionts become
603 organelles? Understanding early events in plastid evolution. *Bioessays* 29, 1239–1246 (2007).
- 604 36. Kim, E. & Maruyama, S. A contemplation on the secondary origin of green algal and plant plastids. *Acta Soc
605 Bot Pol* 83, 331–336 (2014).
- 606 37. Zhu, G., Marchewka, M. J. & Keithly, J. S. *Cryptosporidium parvum* appears to lack a plastid genome.
607 *Microbiology+* 146, 315–321 (2000).
- 608 38. Janouškovec, J. *et al.* Apicomplexan-like parasites are polyphyletic and widely but selectively dependent on
609 cryptic plastid organelles. *Elife* 8, e49662 (2019).
- 610 39. Gornik, S. G. *et al.* Endosymbiosis undone by stepwise elimination of the plastid in a parasitic dinoflagellate.
611 *Proc National Acad Sci* 112, 5767–5772 (2015).
- 612 40. Timmis, J. N., Ayliffe, M. A., Huang, C. Y. & Martin, W. Endosymbiotic gene transfer: organelle genomes
613 forge eukaryotic chromosomes. *Nat Rev Genet* 5, 123–135 (2004).
- 614 41. Archibald, J. M. Genomic perspectives on the birth and spread of plastids. *Proc National Acad Sci* 112,
615 10147–10153 (2015).
- 616 42. Burki, F. *et al.* Re-evaluating the Green versus Red Signal in Eukaryotes with Secondary Plastid of Red
617 Algal Origin. *Genome Biol Evol* 4, 626–635 (2012).
- 618 43. Deschamps, P. & Moreira, D. Reevaluating the Green Contribution to Diatom Genomes. *Genome Biol Evol*
619 4, 683–688 (2012).
- 620 44. Qiu, H., Yoon, H. S. & Bhattacharya, D. Algal endosymbionts as vectors of horizontal gene transfer in
621 photosynthetic eukaryotes. *Front Plant Sci* 4, 1–8 (2013).
- 622 45. Morozov, A. A. & Galachyants, Y. P. Diatom genes originating from red and green algae: Implications for
623 the secondary endosymbiosis models. *Mar Genom* 45, 72–78 (2019).
- 624 46. Sibbald, S. J. & Archibald, J. M. Genomic insights into plastid evolution. *Genome Biol Evol* 12, evaa096-
625 (2020).
- 626 47. Singer, A. *et al.* Massive Protein Import into the Early-Evolutionary-Stage Photosynthetic Organelle of the
627 *Amoeba Paulinella chromatophora*. *Curr Biol* 27, 2763–2773.e5 (2017).
- 628 48. Burki, F. *et al.* Endosymbiotic gene transfer in tertiary plastid-containing dinoflagellates. *Eukaryot Cell* 13,
629 246–255 (2014).

- 630 49. Hehenberger, E., Burki, F., Kolisko, M. & Keeling, P. J. Functional Relationship between a Dinoflagellate
631 Host and Its Diatom Endosymbiont. *Mol Biol Evol* 33, 2376–2390 (2016).
- 632 50. Needham, D. M. *et al.* A distinct lineage of giant viruses brings a rhodopsin photosystem to unicellular
633 marine predators. *Proc National Acad Sci* 116, 20574–20583 (2019).
- 634 51. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *Embnet J* 17, 10–
635 12 (2011).
- 636 52. Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell
637 Sequencing. *J Comput Biol* 19, 455–477 (2012).
- 638 53. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification. *Bmc*
639 *Bioinformatics* 11, 119 (2010).
- 640 54. Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7: Improvements in
641 Performance and Usability. *Mol Biol Evol* 30, 772–780 (2013).
- 642 55. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated alignment trimming
643 in large-scale phylogenetic analyses. *Bioinformatics* 25, 1972–1973 (2009).
- 644 56. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., Haeseler, A. von & Jermini, L. S. ModelFinder: fast
645 model selection for accurate phylogenetic estimates. *Nat Methods* 14, 587–589 (2017).
- 646 57. Minh, B. Q. *et al.* IQ-TREE 2: New models and efficient methods for phylogenetic inference in the genomic
647 era. *Mol Biol Evol* 37, 1530–1534 (2020).
- 648 58. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput ANI
649 analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun* 9, 5114 (2018).
- 650 59. Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for genome assemblies.
651 *Bioinformatics* 29, 1072–1075 (2013).
- 652 60. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat Methods*
653 12, 59–60 (2015).
- 654 61. Eddy, S. R. Accelerated Profile HMM Searches. *Plos Comput Biol* 7, e1002195 (2011).
- 655 62. Whelan, S., Irisarri, I. & Burki, F. PREQUAL: detecting non-homologous characters in sets of unaligned
656 homologous sequences. *Bioinformatics* 34, 3929–3930 (2018).
- 657 63. Ali, R. H., Bogusz, M. & Whelan, S. Identifying Clusters of High Confidence Homologies in Multiple
658 Sequence Alignments. *Mol Biol Evol* 36, 2340–2351 (2019).
- 659 64. Wang, H.-C., Minh, B. Q., Susko, E. & Roger, A. J. Modeling Site Heterogeneity with Posterior Mean Site
660 Frequency Profiles Accelerates Accurate Phylogenomic Estimation. *Systematic Biol* 67, 216–235 (2018).
- 661 65. Lartillot, N., Rodrigue, N., Stubbs, D. & Richer, J. PhyloBayes MPI: Phylogenetic Reconstruction with
662 Infinite Mixtures of Profiles in a Parallel Environment. *Systematic Biol* 62, 611–615 (2013).

- 663 66. Susko, E., Field, C., Blouin, C. & Roger, A. J. Estimation of Rates-Across-Sites Distributions in
664 Phylogenetic Substitution Models. *Systematic Biol* 52, 594–603 (2003).
- 665 67. Viklund, J., Ettema, T. J. G. & Andersson, S. G. E. Independent Genome Reduction and Phylogenetic
666 Reclassification of the Oceanic SAR11 Clade. *Mol Biol Evol* 29, 599–615 (2012).
- 667 68. Zhang, C., Rabiee, M., Sayyari, E. & Mirarab, S. ASTRAL-III: polynomial time species tree reconstruction
668 from partially resolved gene trees. *Bmc Bioinformatics* 19, 153 (2018).
- 669 69. Priyam, A. *et al.* Sequenceserver: a modern graphical user interface for custom BLAST databases. *Mol Biol*
670 *Evol* 36, 2922–2924 (2019).
- 671 70. Jin, J.-J. *et al.* GetOrganelle: a fast and versatile toolkit for accurate de novo assembly of organelle genomes.
672 *Genome Biol* 21, 241 (2020).
- 673 71. Schön, M. E., Eme, L. & Ettema, T. J. G. PhyloMagnet: Fast and accurate screening of short-read meta-
674 omics data using gene-centric phylogenetics. *Bioinformatics* (2019) doi:10.1093/bioinformatics/btz799.
- 675 72. Huson, D. H. *et al.* Fast and simple protein-alignment-guided assembly of orthologous gene families from
676 microbiome sequencing reads. *Microbiome* 5, 11 (2017).
- 677 73. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome*
678 *Biol* 20, 238 (2019).
- 679 74. Rognes, T., Flouri, T., Nichols, B., Quince, C. & Mahé, F. VSEARCH: a versatile open source tool for
680 metagenomics. *Peerj* 4, e2584 (2016).
- 681
- 682

683

684 **Figure Legends**

685 **Figure 1. a)** Assembly length in Mbp for 17 SAGs and CO-SAGs used for further analysis. **b)** Estimated
686 completeness of the 10 most complete SAGs and CO-SAGs as assessed using presence/absence of the BUSCO
687 dataset of 255 eukaryotic markers and a dataset of 317 Phylogenomic marker genes. These 10 assemblies were
688 used for the phylogenomic inference. **c)** Maximum likelihood tree of the 18S rRNA gene, reconstructed using
689 the model GTR+R4+F while support was estimated with 100 non-parametric bootstrap replicates in IQ-TREE.
690 Picozoa CO-SAGs and SAGs are written in bold, the sequences of *Picomonas judraskeda* and the SAGs from
691 Yoon et al. (2011)¹¹ in bold italic. The group labels BP1-3 are from Cuvelier et al. (2008)¹³ and deep branching
692 lineages from Moreira and López-García (2014)¹⁶.

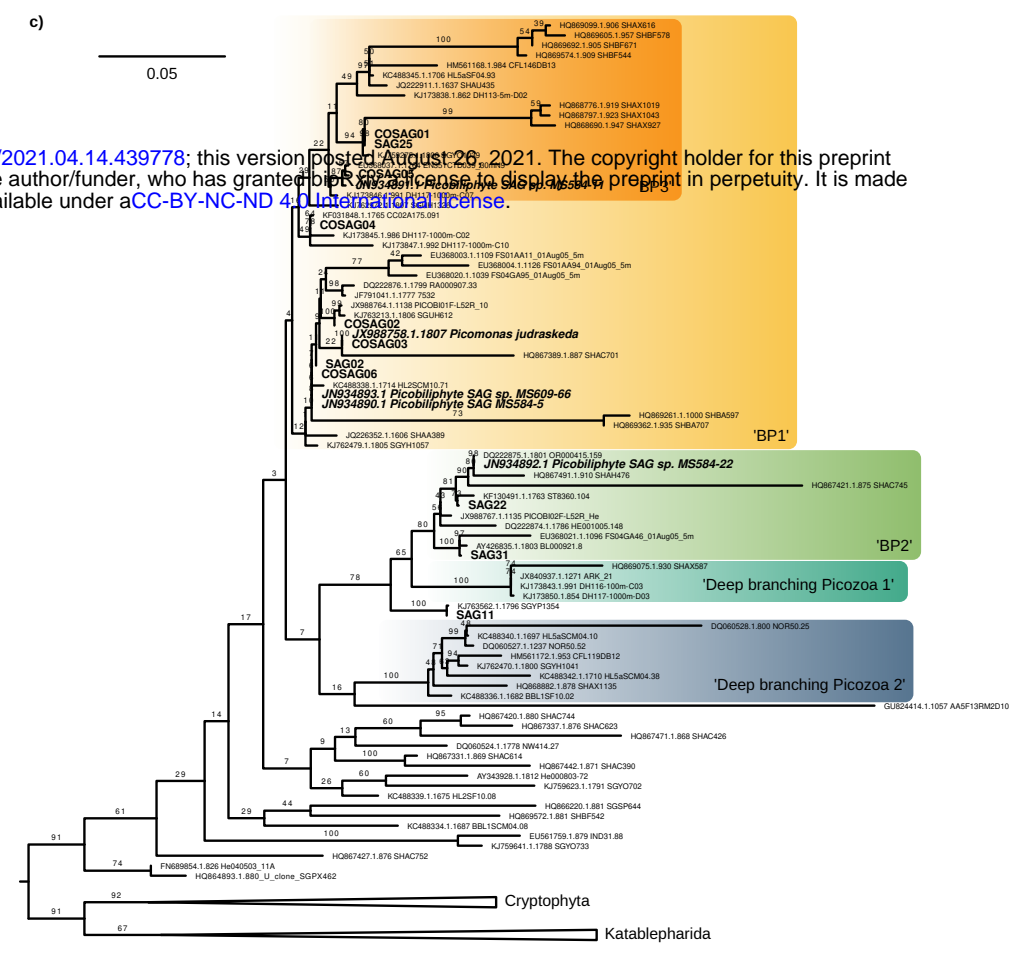
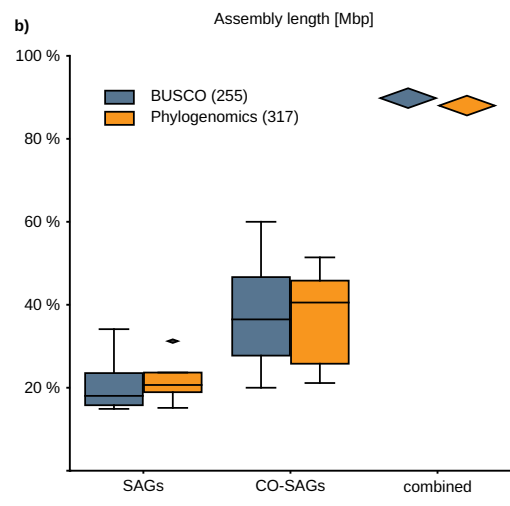
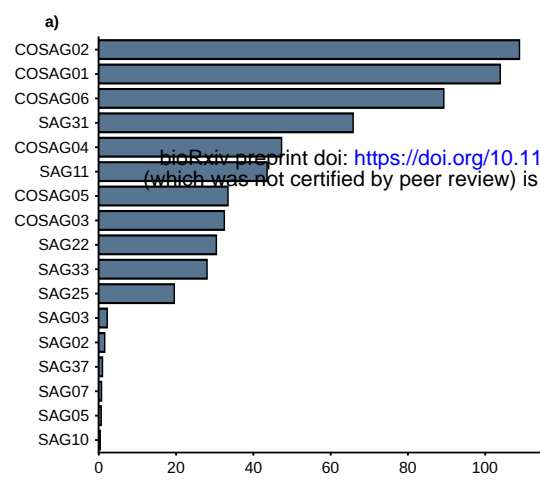
693

694 **Figure 2.** Maximum likelihood tree of eukaryotic species showing the position of Picozoa. The tree is based on
695 the concatenated alignment of 317 marker genes and was reconstructed using the site-heterogeneous model
696 LG+C60+F+G-PMSF. Support values correspond to 100 non-parametric bootstrap replicates/ posterior
697 probability values estimated using PhyloBayes CAT-GTR-G. Black circles denote full support (=100/1.0).

698

699 **Figure 3. a)** Number of inferred endosymbiotic gene transfers (EGT) across a selection of 33 species that
700 represent groups with a photosynthetic plastid (green), a non-photosynthetic plastid (blue), confirmed plastid
701 loss (yellow) and no known plastid ancestry (black). These species serve as a comparison to Picozoa (orange). **b)**
702 The number of EGTs from a) is related to the number of inferred HGT across the same 33 selected species. A
703 number below 1 indicates more HGT than EGT, while numbers above 1 indicate more EGT than HGT. No ratio
704 could be calculated for *Arabidopsis* because there were no detectable HGT events.

705



bioRxiv preprint doi: <https://doi.org/10.1101/2021.04.14.439778>; this version posted August 26, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

