

1 **A new lineage of non-photosynthetic green algae with extreme organellar genomes**

2

3 Tomáš Pánek^{1,2}, Dovilė Barcytė¹, Sebastian C. Treitli³, Kristína Záhonová¹, Martin Sokol¹, Tereza

4 Ševčíková¹, Eliška Zadrobílková², Karin Jaške¹, Naoji Yubuki², Ivan Čepička², Marek Eliáš^{1,*}

5

6 ¹Department of Biology and Ecology, Faculty of Science, University of Ostrava, 701 00 Ostrava,

7 Czech Republic

8 ²Department of Zoology, Faculty of Science, Charles University, 128 43 Prague, Czech Republic

9 ³Department of Parasitology, Faculty of Science, Charles University, BIOCEV, 252 42 Vestec,

10 Czech Republic

11

12 Correspondence: marek.elias@osu.cz

13

14

15

16

17

18

19

20

21

22

23

24 **Abstract**

25 **Background:** The plastid genomes of the green algal order Chlamydomonadales tend to expand
26 their non-coding regions, but this phenomenon is poorly understood. Here we shed new light on
27 organellar genome evolution in Chlamydomonadales by studying a previously unknown non-
28 photosynthetic lineage. We established cultures of two new *Polytoma*-like flagellates, defined their
29 basic characteristics and phylogenetic position, and obtained complete organellar genome
30 sequences and a transcriptome assembly for one of them.

31 **Results:** We discovered a novel deeply diverged chlamydomonadalean lineage that has no close
32 photosynthetic relatives and represents an independent case of photosynthesis loss. To
33 accommodate these organisms, we establish a new genus, *Leontynka*, with two species *L. pallida*
34 and *L. elongata* distinguished by both morphological and molecular characteristics. Notable
35 features of the colourless plastid of *L. pallida* deduced from the plastid genome (plastome)
36 sequence and transcriptome assembly include the retention of ATP synthase, thylakoid-associated
37 proteins, carotenoid biosynthesis pathway, and plastoquinone-based electron transport chain, the
38 latter two modules having an obvious functional link to the eyespot present in *Leontynka*. Most
39 strikingly, the *L. pallida* plastome with its ~362 kbp is by far the largest among non-photosynthetic
40 eukaryotes investigated to date. Instead of a high gene content, its size reflects extreme
41 proliferation of sequence repeats. These are present also in coding sequences, with one repeat type
42 found in exons of 11 out of 34 protein-coding genes and up to 36 copies per gene, affecting thus
43 the encoded proteins. The mitochondrial genome of *L. pallida* is likewise exceptionally large, with
44 its >104 kbp surpassed only by the mitogenome of *Haematococcus lacustris* among all members
45 of Chlamydomonadales studied so far. It is also bloated with repeats, yet completely different from
46 those in the *L. pallida* plastome, which contrasts with the situation in *H. lacustris* where both

47 organellar genomes have accumulated related repeats. Furthermore, the *L. pallida* mitogenome
48 exhibits an extremely high GC content in both coding and non-coding regions and, strikingly, a
49 high number of predicted G-quadruplexes.

50 **Conclusions:** With the unprecedented combination of plastid and mitochondrial genome
51 characteristics, *Leontynka* pushes the frontiers of organellar genome diversity and becomes an
52 interesting model for studying organellar genome evolution.

53 **Keywords:** Chlamydomonadales; GC content; G-quadruplex; green algae; mitochondrial genome;
54 non-photosynthetic algae; plastid genome; repeat expansion

55

56 **Background**

57 Secondary loss of photosynthesis has occurred numerous times across the diversity of plastid-
58 bearing eukaryotes, including land plants (Hadariová *et al.*, 2018; Sibbald & Archibald, 2020).
59 Among algae, photosynthesis loss has been most common among groups characterised by
60 secondary or higher-order plastids, with chrysophytes and myzozoans (including apicomplexans
61 as the best-studied non-photosynthetic “algae”) being the most prominent examples. In green
62 algae, loss of photosynthesis is restricted to several lineages within two classes, Trebouxiophyceae
63 and Chlorophyceae (Figuroa-Martinez *et al.*, 2015). Colourless trebouxiophytes are formally
64 classified in two genera, *Helicosporidium* and the polyphyletic *Prototheca*, collectively
65 representing three independent photosynthesis loss events (Suzuki *et al.*, 2018). While these
66 organisms live as facultative or obligate parasites of metazoans (including humans), non-
67 photosynthetic members of Chlorophyceae are all free-living osmotrophic flagellates. Two genera
68 of such colourless flagellates have been more extensively studied and are represented in DNA
69 sequence databases: the biflagellate *Polytoma* and the tetraflagellate *Polytomella*. They both fall

70 within the order Chlamydomonadales (Volvocales *sensu lato*), but are not closely related to each
71 other. Furthermore, *Polytoma* as presently circumscribed is polyphyletic, since *P. oviforme* does
72 not group with its congeners, including the type species *P. uvella* (Figuroa-Martinez *et al.*, 2015).
73 Hence, photosynthesis was lost at least three times in Chlamydomonadales, but the real number is
74 probably higher, since several other genera of colourless flagellates morphologically falling within
75 this group were historically described (Ettl, 1983), but remain to be studied by modern methods.
76 Indeed, a taxonomically unidentified non-photosynthetic chlamydomonadalean (strain NrCl902),
77 not related to any of the three known lineages, was reported recently (Kayama *et al.*, 2020);
78 whether it corresponds to any of the previously formally described taxa is yet to be investigated.

79 The non-photosynthetic chlamydomonadaleans are not only diverse phylogenetically, but
80 they also exhibit diversity in the features of their residual plastids. Most notably, *Polytomella*
81 represents one of the few known cases of a complete loss of the plastid genome (plastome) in a
82 plastid-bearing eukaryote (Smith & Lee, 2014). In contrast, *Polytoma uvella* harbours the largest
83 plastome amongst all non-photosynthetic eukaryotes studied to date (≥ 230 kbp). This is not due to
84 preserving a large number of genes, but because of the massive accumulation of long arrays of
85 short repeats in intergenic regions (Figuroa-Martinez *et al.*, 2017). The unusual architecture of
86 the *P. uvella* plastome seems to reflect a more general trend of plastome evolution in
87 Chlamydomonadales, i.e. a tendency to increase in size by the expansion of repetitive sequences.
88 An extreme manifestation of this trend was recently unveiled by sequencing the 1.35-Mbp
89 plastome of the photosynthetic species *Haematococcus lacustris*, a record-holder amongst all fully
90 sequenced plastomes to date (Bauman *et al.*, 2018; Smith, 2018). Interestingly, *H. lacustris* also
91 harbours the by far largest known mitochondrial genome (mitogenome) amongst all
92 Chlamydomonadales, which has expanded to 126.4 kbp by the accumulation of repeats highly

93 similar to those found in the plastome, suggesting an inter-organellar transfer of the repeats (Zhang
94 *et al.*, 2019). The mechanistic underpinnings of the repeat accumulation in chlamydomonadalean
95 organellar genomes are still not clear.

96 When studying protists living in hypoxic sediments, we obtained cultures of two colourless
97 flagellates that turned out to represent a novel, deeply separated lineage of Chlamydomonadales.
98 We here describe them formally as two species in a new genus. Using a combination of different
99 DNA sequencing technologies, we determined sequences of organellar genomes of one of the
100 isolates, which turned out to exhibit extreme features concerning the size and/or composition. Our
101 analysis of these genomes provides important new insights into the evolution of organelle genomes
102 in general.

103

104 **Results**

105 **A new lineage of non-photosynthetic Chlamydomonadales with two species**

106 Based on their 18S rRNA gene sequences, the two new isolates – AMAZONIE and MBURUCU
107 – constitute a clade (with full bootstrap support) that is nested within Chlamydomonadales, but
108 separate from all the principal chlamydomonadalean clades as demarcated by Nakada *et al.* (2008)
109 (Fig. 1a). Notably, this new lineage is clearly unrelated to all previously studied non-
110 photosynthetic chlamydomonadaleans, including *Polytomella* (branching off within the clade
111 *Reinhardtinia*), both lineages representing the polyphyletic genus *Polytoma* (*P. uvella* plus several
112 other species in the clade *Caudivolvoxa* and *P. oviforme* in the clade *Xenovolvoxa*), and the strain
113 NrCl902 (also in *Caudivolvoxa*; Additional file 1: Fig. S1). Our two strains were mutually
114 separated in the 18S rRNA gene tree as deeply as other chlamydomonadalean pairs classified as
115 separate species or even genera, and their 18S rRNA gene sequences differed in 13 positions (out

116 of 1703 available for comparison). In addition, the ITS1-5.8S-ITS2 rDNA regions of the two
117 strains exhibited only 88% identity and the differences included several compensatory base
118 changes (CBCs) in the helix II of the characteristic secondary structure of the ITS2 region
119 (Additional file 1: Fig. S2). This and morphology-based evidence presented below led us to
120 conclude that the two strains represent two different species of a new genus of
121 chlamydomonadalean algae, which we propose be called *Leontynka pallida* (strain AMAZONIE)
122 and *Leontynka elongata* (strain MBURUCU). Formal descriptions of the new taxa are provided in
123 Additional file 2: Note S1.

124 The phylogenetic position of *L. pallida* was also studied by using protein sequences
125 encoded by its plastome (see below). Phylogenomic analysis of a concatenated dataset of 24
126 conserved proteins encoded by plastomes of diverse members of Chlorophyceae, including a
127 comprehensive sample of available data from Chlamydomonadales, revealed *L. pallida* as a
128 separate lineage potentially sister to a fully supported broader clade comprising representatives of
129 the clades *Caudivolvoxa* and *Xenovolvoxa* (*sensu* Nakada *et al.*, 2008; Fig. 1b). This position of *L.*
130 *pallida* received moderate support in the maximum likelihood analysis (nonparametric bootstrap
131 value of 78%), but inconclusive support from the PhyloBayes analysis (posterior probability of
132 0.68). Importantly, *L. pallida* was unrelated to *Polytoma uvella* (nested within *Caudivolvoxa* with
133 full support). *Polytoma oviforme* and the genus *Polytomella* were missing from the analysis due to
134 lack of plastome data or the complete absence of the plastome, respectively.

135 Both *Leontynka* species lacked a green plastid (chloroplast). Instead, their cells were
136 occupied by a colourless leucoplast containing starch grains, typically filling most of its volume
137 (Fig. 2, Additional file 1: Figs S4 and S5). Two anterior, isokont flagella approximately as long as
138 the cell body emerged from a keel-shaped papilla (Additional file 1: Fig. S3c, d). Cells of both

139 species also contained two apical contractile vacuoles (Fig. 2c, f, Additional file 1: Fig. S3a, c,
140 Fig. S4a, d, h), a central or slightly posterior nucleus (Additional file 1: Fig. S3c, Fig. S4h),
141 inclusions of yellowish lipid droplets (Fig. 2h, Additional file 1: Fig. S4h), and one or occasionally
142 two eyespots (Fig. 2a, b, and f–j, Additional file 1: Fig. S3a, b, e, g, h, Fig. S4a, c–f, h).
143 Reproduction occurred asexually through zoospore formation, typically with up to four zoospores
144 formed per the mother cell (Additional file 1: Fig. S3g, h). The two species differed in the cell
145 shape and position of the eyespot, as described in more detail in Additional file 2: Note S2.

146 The plastid of both species was bounded by a double membrane and composed of
147 numerous separate compartments connected by narrow “bridges” (Fig. 2d, e, i, j, Additional file
148 1: Fig. S5a–c, h). Each compartment contained either a single large or two smaller starch grains,
149 leaving essentially no room for the stroma or thylakoids. Rarely, starch-free compartments
150 containing membranous inclusions were present (Fig. 2e). The eyespot globules were inside the
151 plastid and were associated with structures that we interpret as thylakoids (Fig. 2j). Mitochondria
152 were highly abundant and contained numerous cristae (Additional file 1: Fig. S5e, f, i). It was
153 impossible to unambiguously determine the crista morphology in *L. pallida* (Additional file 1: Fig.
154 S5i), but in *L. elongata*, the cristae were of the discoidal morphotype (Additional file 1: Fig. S5e).
155 Further details on the ultrastructure of *Leontynka* spp. are presented in Additional file 2: Note S2.

156

157 **The extremely bloated plastome of *Leontynka pallida***

158 A complete plastome sequence was assembled for *L. pallida* using a combination of Oxford
159 Nanopore and Illumina reads. It corresponds to a circular-mapping molecule comprising 362307
160 bp (Fig. 3a). Thirty-four protein-coding genes (including two intronic ORFs), 26 tRNA genes (a
161 standard set presumably allowing for translation of all sense codons), and genes for the three

162 standard rRNAs were identified and annotated in the genome. Three genes are interrupted by
163 introns: *atpA* with one group I intron that contains an ORF encoding a LAGLIDADG homing
164 endonuclease, *tufA* with one group II intron that contains an ORF encoding a reverse
165 transcriptase/maturase protein, and *rnl* with one group I and one group II intron, neither containing
166 an ORF. No putative pseudogenes or apparent gene remnants were identified in the *L. pallida*
167 plastome.

168 No genes that encode proteins directly associated with photosynthetic electron transport
169 components and CO₂ fixation were identified in the *L. pallida* plastome. The genes retained encode
170 proteins involved in transcription (RNA polymerase subunits), translation (*tufA* and ribosomal
171 subunit genes), protein turnover (*clpP*, *ftsH*), and a protein of an unclear function (*ycfI*). Nearly
172 all these genes have been preserved also in the plastome of *P. uvella* (Figuroa-Martinez *et al.*,
173 2017), except for *rps2*. The two non-photosynthetic *Leontynka* species share the absence of genes
174 for two ribosomal proteins: *rpl32*, which is, however, also missing from a subset of photosynthetic
175 representatives of Chlamydomonadales, and *rpl23* conserved in plastomes of all photosynthetic
176 chlorophytes investigated to date (Turmel & Lemieux, 2018). Whereas an Rpl32 protein with a
177 predicted plastid-targeting presequence is encoded by the *L. pallida* nuclear genome (Additional
178 file 3: Table S1), the loss of *rpl23* does not seem to be compensated in a similar way. Interestingly,
179 *rpl23* has been independently lost also from the plastome of *Helicosporidium* sp. (Figuroa-
180 Martinez *et al.*, 2017), which suggests that this ribosomal subunit may become dispensable upon
181 the loss of photosynthesis. In contrast to *P. uvella*, the *L. pallida* plastome has kept the same set
182 of genes encoding ATP synthase subunits as typical for photosynthetic green algae, i.e., *atpA*,
183 *atpB*, *atpE*, *atpF*, *atpH*, and *atpI*. As evidenced by the transcriptome data, the three missing ATP

184 synthase subunits (AtpC, AtpD, and AtpG) are encoded by the *L. pallida* nuclear genome and bear
185 predicted plastid-targeting signals (Additional file 3: Table S1).

186 What makes the *L. pallida* plastome truly peculiar are the intergenic regions. Their average
187 length is 4.7 kbp, which is 1.5 and five times more than the average length of the intergenic regions
188 in the plastomes of *P. uvella* and *C. reinhardtii*, respectively (Table 1). Furthermore, while the GC
189 content of the *P. uvella* intergenic regions is vastly different from that of coding regions (19%
190 versus 40%), the GC content of these two plastome partitions are highly similar in *L. pallida* (as
191 well as in *C. reinhardtii*; Table 1). A self-similarity plot generated for the *L. pallida* plastome
192 revealed a massive repetitiveness of the DNA sequence, with only short islands of unique
193 sequences scattered in the sea of repeats (Fig. 3b). The repeats are highly organised and occur in
194 various arrangements: tandem repeats, interspersed repeats, inverted repeats (palindromes), and
195 other higher-order composite repeated units. As an example, let us take the most abundant repeat,
196 the imperfect palindrome (IP) CAAACCAGT|NN|ACTGGTTAG. It is present in more than 1300
197 copies, with the dinucleotide AA as a predominant form of the internal spacer. This repeat is mostly
198 localised in clusters (>1200 cases) where its copies are interleaved by a repeat with the conserved
199 sequence TAACTAACTTC, so together they constitute a composite extremely abundant tandem
200 repeat. In a single region, the palindromic repeat combines with a different interspersed repeat
201 (TAACTACTT), together forming a small cluster of composite tandem repeats in 14 copies.
202 Besides, the same palindromic repeat is also part of another, 146 bp-long repeat present in 27
203 copies across the plastome (for details see Additional file 2: Note S3).

204 Apart from intergenic regions, sequence repeats are found also in the four introns present
205 in the *L. pallida* plastome. The most prominent is a cluster of 20 copies of the motif
206 TGGTTAGTAACTAACTTCCAAACCAGTAAAC in the intron inside the *atpA* gene that is

207 abundant also in intergenic regions (more than 1,000 copies typically located in huge clusters).
208 Strikingly, when analysing the distribution of the most abundant IPs, we noticed that the motif
209 AAGCCAGC|NNN|GCTGACTT and its variants are present also in coding regions, namely in
210 exons of 11 out of 34 protein-coding genes of *L. pallida* plastome (Fig. 4a-c). They can be present
211 in up to 36 copies per gene (“variant 8” in exons of *rpoC2*; Fig. 4b). In most cases the IP motif is
212 part of a longer repeat unit including extra nucleotides at both ends (“variant 4” to “variant 8” in
213 Fig. 4c). The most complex repeat unit variant is the following one (the IP core in round brackets;
214 square brackets indicate alternative nucleotides occurring at the same position): AAAGAT-
215 (AAGTCAGC|AGA|GCTGAC[AT]T)-
216 CCAGACCACTAAAGTGGTCAGTAACTAAAAGTTAT. It is restricted to coding sequences
217 (i.e., is absent from intergenic regions and introns) and occurs in eight copies inside three genes
218 (*rpoC2*, *rpoC1*, *ftsH*), resulting in an insertion of a stretch of 20 amino acid residues in the encoded
219 proteins. Other repeat variants (listed in Fig. 4c) have proliferated in exons as well as intergenic
220 regions and introns. However, none of the aforementioned nucleotide motifs were found in
221 plastomes of other chlamydomonadalean algae, indicating they have originated and diverged only
222 in the *Leontynka* lineage.

223 Manual inspection of protein sequence alignments including chlamydomonadalean
224 orthologs of the *L. pallida* proteins revealed that the intraexonic repeat insertions are located
225 mostly in poorly conserved regions (see Fig. 4d for an example). Preferential proliferation in
226 variable parts of coding regions is consistent with a high abundance of these motifs in proteins that
227 exhibit a general tendency for including rapidly evolving and poorly conserved regions, namely
228 FtsH (Additional file 1: Fig. S6), Ycf1, RpoC1, RpoC2, and RpoBb. Interestingly, the phase and
229 orientation of the intraexonic repeats with respect to the reading frame and the direction of

230 transcription is not random and is potentially biased such that not only termination codons, but
231 also codons generally rare in *L. pallida* plastid coding sequences are avoided from the actual frame
232 in which the insertion is read during translation (for details see Fig. 4c, e, Additional file 2: Note
233 S4). This bias does not merely reflect a possible bias in the orientation of the repeats relative to
234 the DNA strand of the genome, as the repeats are distributed roughly equally in both strands when
235 counted at the whole-genome level (Fig. 4c).

236

237 **A high number of potential quadruplex-forming sequences in the GC-rich mitogenome of**
238 ***Leontynka pallida***

239 The mitogenome sequence was assembled from Nanopore and Illumina reads as a linear molecule
240 of 110515 bp with long (~5770 bp) nearly perfect (97.7% identity) direct terminal repeats differing
241 primarily by the presence/absence of two short repetitive regions (13 and 70 bp) (Fig. 3c,
242 Additional file 1: Fig. S7). This possibly indicates that the *L. pallida* mitogenome is in fact circular,
243 with the slight differences in the terminal direct repeats of the assembled linear contigs reflecting
244 sequence variability of a particular genomic region between the different genome copies in *L.*
245 *pallida* or possibly sequencing or assembly artefacts. If circular, the mitogenome would then have
246 a length of ~104812 bp. The suspected circularity of the mitogenome is also compatible with the
247 absence of the *rtl* gene, which is present in all linear mitogenomes of Chlamydomonadales
248 characterised to date and encodes a reverse transcriptase-like protein implicated in the replication
249 of the mitogenome termini (Smith & Craig, 2021). Apart from *rtl*, the gene content of the *L. pallida*
250 mitogenome is essentially the same as in other chlamydomonadalean mitogenomes sequenced
251 before and includes seven protein-coding genes (with *cox1* interrupted by an ORF-free group II
252 intron), only three tRNA genes, and regions corresponding to the 16S and 23S rRNA genes. As in

253 other chlamydomonadales studied in this regard (Boer & Gray, 1988; Denovan-Wright *et al.*,
254 1994; Fan *et al.*, 2003), the mitochondrial 16S and 23S rRNA genes in *L. pallida* are fragmented,
255 consisting of multiple separately transcribed pieces. Four fragments, together constituting a
256 presumably complete 16S rRNA, were annotated by considering the sequence and secondary
257 structure conservation of the molecule. The number of the 16S rRNA fragments is thus the same
258 as in *Chlamydomonas reinhardtii*, but the breakpoints are not completely identical. Due to a lower
259 conservation of the 23S rRNA molecule, we could identify only a few of the presumed gene
260 fragments in the *L. pallida* mitogenome.

261 The large size and the low density of coding sequences of the *L. pallida* mitogenome
262 (~84.7% of its complete sequence is represented by intergenic regions) are atypical for
263 Chlamydomonadales, including the other non-photosynthetic species: the mitogenome of *P. uvella*
264 is 17.4 kbp long (Del Vasto *et al.*, 2015), and in *Polytomella* spp. the mitogenome size ranges from
265 ~13 to 24.4 kbp (Smith *et al.*, 2010; Smith *et al.*, 2013). In fact, the *L. pallida* mitogenome can be
266 compared only to the recently characterised mitogenome of *Haematococcus lacustris*, which with
267 the same gene content is even larger (126.4 kbp) yet with a similar representation of intergenic
268 regions (83.2%). A self-similarity plot generated for the *L. pallida* mitogenome revealed a highly
269 repetitive nature of the genome sequence (Fig. 3d), similar to the plastome. However, the repeats
270 are distributed less evenly than in the plastome, being present particularly in the terminal regions
271 of the assembled linear sequence and in several internal hotspots.

272 With the GC content 62.6% (as counted for the circularised version of the genome), the
273 mitogenome of *L. pallida* has the third highest documented mitochondrial GC content out of
274 11,077 examined mitogenomes available in GenBank, being surpassed only by the lycophyte
275 *Selaginella moellendorffii* (68.2%; Hecht *et al.*, 2011) and the green alga *Picocystis salinarum*

276 (67.7%). These values contrast sharply with the median GC content value for the whole set of the
277 mitogenomes examined, i.e. 38%. We also encountered an exceptionally high GC content (63.4%)
278 and a strong bias towards using GC-rich codons in all protein-coding genes in the *L. pallida*
279 mitogenome (see Additional file 3: Tables S2 and S3). Only two organisms are presently known
280 to have an even higher GC content of mitochondrial protein-coding genes: the sponge
281 *Leucosolenia complicata* (71.2%; Lavrov *et al.*, 2016) and *P. salinarum* (67.9%). Some *L. pallida*
282 mitogenome-encoded proteins, namely Nad2 and Nad5, also exhibit a higher relative content of
283 amino acids with GC-rich codons (G, A, R, and P) compared to most of their orthologs in other
284 species (Additional file 3: Table S4). Thus, not only the expanded GC-rich intergenic regions, but
285 also coding regions of the *L. pallida* mitogenome contribute to its extremely high GC content.

286 The repeats in the plastome and mitogenome of *H. lacustris* are nearly identical (Zhang *et*
287 *al.*, 2019), so it was interesting to compare the two *L. pallida* organellar genomes to find out
288 whether they behave similarly. However, as follows from the respective similarity plot (Fig. 3e)
289 and comparison of most abundant inverted repeats and palindromes (Fig. 4a), the repeats in the
290 two genomes do not resemble each other. The proliferation of different repeats in the two
291 organellar genomes of *L. pallida* at least partially accounts for their strikingly different GC content
292 (62.6% vs 37%). Interestingly, the most abundant IP in the *L. pallida* mitogenome contains the
293 GGGG motif (Fig. 4a), which prompted us to bioinformatically investigate the possible occurrence
294 of G-quadruplexes, unusual secondary structures in nucleic acids formed by guanine-rich regions
295 (Burge *et al.*, 2006). Indeed, the *L. pallida* mitogenome was suggested to include up to 14.7
296 potential quadruplex-forming sequences (PQS) per 1,000 bp. A similar value was inferred for the
297 *S. moellendorffii* mitogenome (15.6 PQS per 1,000 bp), whereas the other mitochondrial and
298 plastid genomes that we analysed for comparison (for technical reasons focusing on GC-rich

299 genomes only) exhibited a much lower values (0.0-6.9 PQS per 1,000 bp; see Additional file 3:
300 Table S2).

301

302 **Discussion**

303 Both 18S rRNA and plastid gene sequence data concur on the conclusion that the two strains
304 investigated in this study, AMAZONIE and MBURUCU, represent a phylogenetically novel
305 lineage within Chlamydomonadales that is unrelated to any of the previously known non-
306 photosynthetic lineages in this order, i.e. *Polytomella*, *Polytoma sensu stricto* (including the type
307 species *P. uvella*), *Polytoma oviforme*, and the recently reported strain NrCI902. However,
308 morphological features of AMAZONIE and MBURUCU, including the cell shape and the
309 presence of two flagella, papilla, eyespot, and starch granules, make our organisms highly
310 reminiscent of the genus *Polytoma* (Ettl, 1983). This is consistent with the previous insight that
311 the *Polytoma* morphotype does not define a coherent phylogenetic unit (Figueroa-Martinez *et al.*,
312 2015). All other historically described genera of colourless flagellates assigned formerly to
313 Chlamydomonadales are sufficiently different from our strains as to consider them a potential
314 taxonomic home for AMAZONIE and MBURUCU (see Additional file 2: Note S5), justifying the
315 erection of the new genus *Leontynka* to accommodate the two strains. Furthermore, these strains
316 clearly differ from each other in morphology (cell shape and size, position of the eyespot) and are
317 genetically differentiated, as apparent from the comparison of the 18S rRNA gene and ITS2 region
318 sequences. Indeed, given the presence of several CBCs in the helix II of the conserved ITS2
319 secondary structure, the two strains are predicted to be sexually incompatible and hence
320 representing separate “biological species” (Coleman, 2000; Wolf *et al.*, 2013). We considered a
321 possibility that AMAZONIE and MBURUCU may represent some of the previously described

322 *Polytoma* species, but as detailed in Additional file 2: Note S5, none seems to be close enough in
323 morphology as reported in the original descriptions. Given the fact that the majority of *Polytoma*
324 species have been isolated and described from central Europe whereas our strains both come from
325 tropical regions of South America, it is not so surprising that we encountered organisms new to
326 science.

327 *Leontynka* spp. exhibit a number of ultrastructural similarities to the previously studied
328 *Polytoma* species (Lang, 1963; Siu *et al.*, 1976; Gaffal & Schneider, 1980). For example, although
329 photosynthetic chlamydomonadalean flagellates usually contain only a few mitochondria
330 squeezed between the nucleus and the plastid, the cells of non-photosynthetic taxa, including
331 *Leontynka*, are mitochondria-rich. It is possible that the proliferation of mitochondria compensates
332 for the loss of the energetic function of the plastid in the non-photosynthetic species. Previous
333 ultrastructural studies of *Polytoma obtusum* (Siu *et al.*, 1976) and *Polytomella* sp. (Dudkina *et al.*,
334 2010) showed that their mitochondria possess lamellar or irregular tubulo-vesicular cristae,
335 respectively. The cristae of *L. pallida* resemble the latter morphotype, whereas *L. elongata* most
336 probably possesses discoidal cristae (Additional file 1: Fig. S5e, f). Discoidal cristae are a very
337 rare morphotype within the supergroup Archaeplastida, although they apparently evolved several
338 times independently during the eukaryote evolution (Pánek *et al.*, 2020) and were previously
339 noticed in several other non-photosynthetic chlorophytes (*Polytoma uvella*, *Polytomella agilis*, and
340 *Prototheca zopfii*; Webster *et al.*, 1967).

341 A particularly notable feature of *Leontynka* spp. is the presence of two eyespots. These
342 were more frequent in *L. elongata* (about half of the cells had two eyespots), whereas in the *L.*
343 *pallida* cultures, such cells were rather rare. Variation in the number of eyespots (from none to
344 multiple) in *Chlamydomonas reinhardtii* was shown to be a result of genetic mutations (Lamb *et*

345 *al.*, 1999), but the factors behind the eyespot number variation observed in *Leontynka* spp. are
346 unknown. The reddish colour of the *Leontynka* eyespots suggests the presence of carotenoids
347 (similar to the eyespot of *C. reinhardtii*; Böhm & Kreimer, 2020). In addition, searches of the *L.*
348 *pallida* transcriptome assembly revealed the presence of a homolog of the *C. reinhardtii* eyespot-
349 associated photosensor channelrhodopsin 1 (ChR1) that is the requires a carotenoid derivative,
350 retinal, as a chromophore (Petroustos, 2017; Additional file 3: Table S1). The preservation of the
351 plastid-localized carotenoid biosynthetic pathway in non-photosynthetic eyespot-bearing
352 chlamydomonadales, namely certain *Polytomella* species and the strain NrCI902, has been noted
353 before (Asmail & Smith, 2016; Kayama *et al.*, 2020), and the same holds true for *L. pallida* based
354 on our analysis of its transcriptome assembly (Additional file 3: Table S1). Notably, like the
355 *Polytomella* species and the strain NrCI902 (Kayama *et al.*, 2020), *L. pallida* has also retained
356 enzymes for the synthesis of plastoquinone, which serves as an electron acceptor in two reactions
357 of carotenoid biosynthesis, and the plastid terminal oxidase (PTOX), which recycles plastoquinone
358 (from its reduced form plastoquinole) by passing the electrons further to molecular oxygen
359 (Additional file 3: Table S1). *Leontynka* thus represents an independent case supporting the notion
360 that retention of the eyespot constrains the reductive evolution of a non-photosynthetic plastid.

361 *Leontynka* is significant not only as a novel non-photosynthetic group *per se*, but also as
362 an independent lineage within Chlamydomonadales lacking any close photosynthetic relatives.
363 Specifically, based on the phylogenetic analysis of plastome-encoded proteins, *Leontynka*
364 branches off between two large assemblages, each comprised of several major
365 chlamydomonadalean clades defined by Nakada *et al.* (2008). One of these assemblages
366 (potentially sister to *Leontynka*) is comprised of the *Caudivolvax* and *Xenolvovax* clades, the
367 other includes *Reinhardtinia*, *Oogamochlamydia*, and the genus *Desmotetra* (Fig. 1a). The

368 radiation of the *Reinhardtinia* clade itself was dated to ~300 MYA (Herron *et al.*, 2009), so the
369 last common ancestor of *Leontynka* and any of its presently known closest photosynthetic relative
370 must have existed even earlier. In other words, it is possible that *Leontynka* has been living without
371 photosynthesis for hundreds of millions of years. The loss of photosynthesis in the four other
372 known colourless chlamydomonadalean lineages certainly does not trace that far in the past.
373 Specifically, the origin of *Polytomella* must postdate the radiation of *Reinhardtinia*, owing to the
374 position of the genus with this clade, whereas *Polytoma sensu stricto* (*P. uvella* and relatives) has
375 close photosynthetic relatives (*Chlamydomonas leiostraca*, *C. applanata* etc.) within the clade
376 *Polytominia* in *Caudivolvox* (Fig. 1, Additional file 1: Fig. S1). *Polytoma oviforme* is specifically
377 related to the photosynthetic *Chlamydomonas chlamydogama*, together constituting a clade in
378 *Xenovolvox* that has not been formally recognised before and which we here designate
379 “*Oviforminia*” (Additional file 1: Fig. S1). Finally, the recently reported strain NrCI902 is closely
380 related to the photosynthetic *Chlamydomonas pseudoplanoconvexa* (Fig. 1A; Additional file 1:
381 Fig. S1). The independent phylogenetic position of *L. pallida* based on plastome-encoded proteins
382 is unlikely an artefact stemming from increased substitution rate of *L. pallida* plastid genes
383 manifested by the markedly longer branch of *L. pallida* in the tree compared to most other species
384 included in the analysis. Indeed, the branches of *P. uvella* and the strain NrCI902 are even longer
385 (Fig. 1B), yet both organisms are placed at positions consistent with the 18S rRNA gene tree (Fig.
386 1A; Additional file 1: Fig. S1). Nevertheless, whether *Leontynka* represents a truly ancient non-
387 photosynthetic lineage or whether it diverged from a photosynthetic ancestor rather recently needs
388 to be tested by further sampling of the chlamydomonadalean diversity, as we cannot rule out the
389 possibility that photosynthetic organisms closely related to the genus *Leontynka* are eventually
390 discovered.

391 The presented considerations about the different ages of the separately evolved non-
392 photosynthetic chlamydomonadalean lineages are somewhat at odds with features of their
393 plastomes. Despite the presumably more recent loss of photosynthesis compared to *Leontynka*,
394 both *P. uvella* and the strain NrCl902 exhibit a more reduced set of plastid genes (Table 1), whereas
395 in *Polytomella*, plastome reduction triggered by photosynthesis loss has reached its possible
396 maximum, i.e., a complete disappearance of the genome. It is likely that factors other than
397 evolutionary time are contributing to the different degrees of plastome reduction in different
398 evolutionary lineages, although little is known in this regard. Compared to *P. uvella* and the strain
399 NrCl902, *L. pallida* has preserved one gene for a plastidial ribosomal protein (*rps2*) and,
400 intriguingly, all standard plastidial genes for ATP synthase subunits, complemented by three more
401 subunits encoded by the nuclear genome to allow for the assembly of a complete and presumably
402 functional complex. It was proposed that the retention of ATP synthase in certain non-
403 photosynthetic plastids is functionally linked to the retention of the twin-arginine protein
404 translocase (Tat; Kamikawa *et al.*, 2015). The function of the translocase depends on a
405 transmembrane proton gradient, which in photosynthetic plastids is primarily generated by the
406 photosynthetic electron transport chain, whereas in non-photosynthetic ones its build-up would
407 depend solely on the function of ATP synthase working in the opposite direction, i.e. pumping
408 protons against the gradient at the expense of ATP. Interestingly, we found homologs of all three
409 Tat subunits (TatA, TatB, TatC) in the nuclear transcriptome of *L. pallida* (Additional file 3: Table
410 S1), providing further support to the hypothesis by Kamikawa *et al.* (2015). However, it must be
411 noted that certain members of the non-photosynthetic trebouxiophyte genus *Prototheca* possess
412 the plastidial ATP synthase in the absence of the Tat translocase (Suzuki *et al.*, 2018), suggesting
413 that ATP synthase may be retained by a non-photosynthetic plastid for roles other than just

414 supporting the function of the Tat system. Directly relevant for the retention of ATP synthase in
415 *L. pallida* might be its role in the functioning of the eyespot hypothesized in *C. reinhardtii*
416 (Schmidt *et al.*, 2007).

417 The Tat translocase and the ATP synthase are both normally localised to the thylakoid
418 membrane. While thylakoids may seem dispensable in a non-photosynthetic plastid, it seems there
419 are putative thylakoids present in *Leontynka*, associated with the eyespot (Fig. 2j). Indeed, in the
420 well-studied cases of *C. reinhardtii* and some other chlamydomonadalean algae, the layers of
421 pigment granules are organized on the surface of thylakoids closely apposed to the plastid envelope
422 (Kreimer, 1994; Böhm & Kreimer, 2020). Interestingly, our searches of the *L. pallida*
423 transcriptome assembly revealed the presence of homologs of additional proteins functionally
424 associated with thylakoids. These include components of several additional thylakoid-associated
425 protein targeting or translocation systems (Schünemann *et al.*, 2007; Skalitzky *et al.*, 2011; Ziehe
426 *et al.*, 2017), namely the plastidial SRP pathway (cpSRP54 and cpFtsY), ALB3 protein insertase,
427 and thylakoid-specific Sec translocase (Additional file 3: Table S1). Furthermore, we also found
428 in *L. pallida* homologs of proteins implicated in thylakoid biogenesis, such as VIPP1, FZL, THF1,
429 or SCO2 (Mechela *et al.*, 2019; Additional file 3: Table S1). Interestingly, some of the
430 corresponding transcripts have very low read coverage or are even represented by incomplete
431 sequences, suggesting a low level of gene expression and presumably low abundance of the
432 respective proteins. These observations support the notion that the thylakoid system is preserved
433 in *Leontynka* plastids, however inconspicuous and likely reduced. Nevertheless, we cannot rule
434 out that at least some of these proteins or complexes may have relocalised to the inner bounding
435 membrane of the *Leontynka* plastid, or even to a cellular compartment other than the plastid (as
436 suggested for some of these proteins by the results of *in silico* targeting prediction (Additional file

437 3: Table S1). The exact localisation of these complexes, the actual substrates of the plastidial SRP
438 pathway, ALB3 insertase, and Tat and Sec translocases, and indeed, the physiological functions
439 of the *L. pallida* leucoplast as a whole remain subjects for future research.

440 The most intriguing feature of *L. pallida* is the extreme expansion of its organellar
441 genomes. Generally, organellar genomes show a remarkable variation in the gene content,
442 architecture, and nucleotide composition, with most of them being AT-rich. The *L. pallida*
443 plastome is no exception in this respect, since its GC content is only ~37%. As noted by Smith
444 (2018), 98 % of plastomes are under 200 kbp and harbour modest amounts (<50 %) of non-coding
445 DNA. The *L. pallida* plastome, reaching 362.3 kbp, may not seem that impressive in comparison
446 with the giant plastomes recently reported from some photosynthetic species, including a distantly
447 related chlamydomonadalean *Haematococcus lacustris* (1.35 Mbp; Bauman *et al.*, 2018) or certain
448 red algae (up to 1.13 Mbp; Muñoz-Gómez *et al.*, 2017). However, it by far dwarfs plastomes of
449 all non-photosynthetic eukaryotes studied to date. The previous record holder, the plastome of
450 *Polytoma uvella* with ~230 kbp (Figuroa-Martinez *et al.*, 2017), accounts for only two thirds of
451 the size of the *L. pallida* plastome. The difference is not only because of a higher number of genes
452 in the latter, but primarily because of a more extreme expansion of intergenic regions in *L. pallida*
453 (4.7 kbp on average) than in *P. uvella* (3.0 kbp on average; Table 1). The plastome of the strain
454 NrCl902 with its size of 176.4 kbp, while exhibiting the same gene content as the *P. uvella*
455 plastome, is much less extreme (Kayama *et al.*, 2020), although still with the intergenic regions
456 substantially expanded as compared to the plastomes of non-photosynthetic trebouxiophytes
457 (Table 1).

458 Thus, despite its uniqueness, the organisation of the *L. pallida* plastome fits into the general
459 pattern observed in chlamydomonadalean algae, where plastomes in different lineages tend to

460 increase in size by accumulating repetitive sequences (Gaouda *et al.*, 2018; Smith, 2018). It was
461 suggested that the repeats are prone to double-strand breaks, which are then repaired by an error-
462 prone mechanism favouring repeat expansion (Smith, 2020a). However, the plastome of *L. pallida*
463 is bloated not only due to extreme proliferation of repetitive DNA in intergenic regions, but also
464 due to the expansion of some of them into the intronic regions and, much more surprisingly, even
465 into exons (Fig. 4). The biased orientation and phase of the insertions with respect to the coding
466 sequence and the reading frame avoid introduction of termination codons as well as rare codons
467 or codons for rare amino acids (C, W) into the coding sequences (Fig. 4 c, d, Additional file 2:
468 Note S4), which suggests that purifying selection eliminates those insertions that would disrupt or
469 reduce the efficiency of translation of the respective mRNAs. Still, exons provide an important
470 niche for the repeats: for example, for the “variant 8” repeat, the exonic copies constitute ~12.2%
471 of the whole repeat population (compared to protein-coding sequences constituting ~17.2% of the
472 total plastome length)! Such a massive proliferation of repeats to coding regions is unprecedented
473 to our knowledge, although a much less extensive invasion of a different repeat into coding
474 sequences was recently noticed in the plastome of another chlamydomonadalean alga,
475 *Chlorosarcinopsis eremi* (Smith 2020a). Here the repeats are found in small numbers in the genes
476 *ftsH*, *rpoC2*, and *ycf1*, paralleling the situation in *L. pallida* and consistent with the notion that
477 genes encoding proteins rich in poorly conserved regions are most likely to tolerate the invasion
478 of the repeats.

479 Recent sequencing of the mitogenome of *H. lacustris*, which is inflated by the
480 accumulation of repeats highly similar to those found in the plastome of the same species (Zhang
481 *et al.*, 2019), provided the first evidence that error-prone repair of double-strand breaks leading to
482 repeat proliferation may operate also in chlamydomonadalean mitochondria. Smith (2020b)

483 recently reported the presence of highly similar repeats in the mitogenome of another
484 chlamydomonadalean alga, *Stephanosphaera pluvialis*, and proposed horizontal gene transfer
485 between the *H. lacustris* and *S. pluvialis* lineages as a possible explanation for the sharing of
486 similar mitochondrial repeats by the two organisms. Our characterisation of the *L. pallida*
487 mitogenome, which is also repeat-rich and larger than any chlamydomonadalean algal
488 mitogenome sequenced to date except that from *H. lacustris*, revealed that mitogenome inflation
489 may be more common in Chlamydomonadales. However, in contrast to *H. lacustris*, the GC
490 content as well as the repeats in the two organellar genomes of *L. pallida* differ significantly
491 (Additional file 3: Table S2), so the evolutionary path leading to the parallel inflation of both
492 genomes in this lineage may have been completely different from the one manifested in *H.*
493 *lacustris*. Strikingly, the specific nature of the mitochondrial repeats in *L. pallida* entails the high
494 abundance of PQS in the mitogenome. G-quadruplexes are increasingly recognised as regulatory
495 structures (Hänsel-Hertsch *et al.*, 2016), and they can form also in the mitogenomes, although their
496 role in mtDNA still needs to be elucidated (Falabella *et al.*, 2019). However, the PQS abundance
497 in the *L. pallida* mitogenome is truly extreme and comparable only with the situation in the
498 mitogenome of the lycophyte *S. moellendorffii* (Additional file 3: Table S2). Both species are thus
499 interesting candidates for studying the role of G-quadruplexes in mitochondrial DNA.

500

501 **Conclusions and future directions**

502 Our study indicates that continued sampling of microbial eukaryotes is critical for further progress
503 in our knowledge of the phylogenetic diversity of life and for better understanding of the general
504 principles governing the evolution of organellar genomes. The specific factors contributing to the
505 propensity of chlamydomonadalean organellar genomes to accumulate repetitive sequences,

506 reaching one of its extremes in *L. pallida*, remain unknown and may not be easy to define.
507 However, future research on *Leontynka*, including characterisation of organellar genomes of *L.*
508 *elongata*, may bring additional insights into the molecular mechanisms and evolutionary forces
509 shaping the organellar genomes in this group. It will also be important to perform a detailed
510 comparative analysis of the molecular machinery responsible for genome replication and
511 maintenance in Chlamydomonadales and other green algae. The transcriptome assembly reported
512 here for *L. pallida* will be instrumental not only in this enterprise, but will also serve as a resource
513 for exploring the full range of physiological roles of the plastid in the *Leontynka* lineage and may
514 help to further clarify the phylogenetic position of *Leontynka* within Chlamydomonadales. We
515 posit that *Leontynka* may become an important model system for analysing the evolutionary and
516 functional aspects of photosynthesis loss in eukaryotes with primary plastids.

517

518 **Methods**

519 **Isolation, cultivation, and basic characterisation of new protist strains**

520 Two strains, AMAZONIE and MBURUCU, were obtained from freshwater hypoxic sediment
521 samples collected in Peru and Argentina, respectively. The strains were cultivated and
522 morphologically characterised by light and transmission electron microscopy, using routine
523 methods. Basic molecular characterisation was achieved by determining partial sequences of the
524 rDNA operon. Further details are provided in Additional file 2: Methods S1-S3.

525

526 **Organellar genome and nuclear transcriptome sequencing**

527 Bacterial contamination in the AMAZONIE culture was minimised by filtration, and DNA and
528 RNA were extracted using standard protocols detailed in Additional file 2: Methods S4. Nanopore

529 sequencing was performed using 4 µg of genomic DNA. The DNA was sheared at 20 kbp using
530 Covaris g-TUBE (Covaris) according to the manufacturer's protocol. After shearing, two libraries
531 were prepared using Ligation Sequencing Kit from Oxford Nanopore Technologies (SQK-
532 LSK108). The prepared library was loaded onto a R9.4.1 Spot-On Flow cell (FLO-MIN106).
533 Sequencing was performed on a MinION Mk1B machine for 48 hours using the MinKNOW 2.0
534 software. Basecalling was performed using Guppy 3.0.3 with the Flip-flop algorithm. Illumina
535 sequencing of the genomic DNA was performed using 1 µg of genomic DNA with the Illumina
536 HiSeq 2000 (2x150bp) paired-end technology with libraries prepared using TruSeq DNA PCR-
537 Free (Illumina, San Diego, CA) at Macrogen Inc. (Seoul, South Korea). The transcriptome was
538 sequenced using the HiSeq 2000 (2x100bp) paired-end technology with libraries prepared using
539 the TruSeq RNA sample prep kit v2 (Illumina, San Diego, CA) at Macrogen Inc. (Seoul, South
540 Korea).

541

542 **Organelle genome and nuclear transcriptome assembly**

543 Raw Illumina sequencing reads were trimmed with Trimmomatic v0.32 (Bolger *et al.*, 2014).
544 Initial assembly of the Oxford Nanopore data was performed using Canu v1.7 with the
545 corMaxEvidenceErate set to 0.15 (Koren *et al.*, 2017). After assembly, the plastome-derived
546 contigs were identified using BLAST (Altschul *et al.*, 1997) with the *Chlamydomonas reinhardtii*
547 plastome as a query. Nine putative plastid genome sequences were selected and polished using the
548 raw nanopore reads with Nanopolish (Loman *et al.*, 2015) followed by polishing with Illumina
549 reads with Pilon v1.22 (Walker *et al.*, 2014). After polishing of the contigs, the Illumina reads
550 were re-mapped onto them, and the mapped reads were extracted and used as an input in Unicycler
551 v0.4.8 (Wick *et al.*, 2017) together with the nanopore reads. Unicycler generated a single circular

552 contig of 362,307 bp. For the mitogenome, a single linear contig was identified in the Canu
553 assembly with BLAST with standard mitochondrial genes as queries; the contig sequence was
554 polished using the same method as described above, but it remained linear after a subsequent
555 Unicycler run. However, direct inspection of the contig revealed highly similar regions (about
556 5,600 bp in length) at both termini. The terminal regions were further polished by mapping of
557 Illumina genomic reads using BWA (Li, 2013) and SAMtools (Li *et al.*, 2009) followed by manual
558 inspection in Tablet (Milne *et al.*, 2016), which increased the sequence similarity of the termini to
559 97.7% (along the region of 5,771 bp).

560 Illumina genomic reads were also assembled separately with the SPAdes Genome
561 assembler v3.10.1 (Bankevich *et al.*, 2012) and used for cleaning the transcriptomic data as
562 follows. Contaminant bacterial contigs > 400,000 bp that were identified with BLAST in the
563 SPAdes (16 contigs) and Canu assemblies (11 contigs), together with published genome
564 assemblies of close relatives of bacteria identified in the AMAZONIE culture (*Curvibacter*
565 *lanceolatus* ATCC 14669, *Bacteroides luti* strain DSM 26991, and *Paludibacter jiangxiensis*
566 strain NM7), were used for RNA-seq read mapping (Hisat2 2.1.0; Kim *et al.*, 2015) to identify and
567 remove bacterial transcriptomic reads that survived the filtration of the culture and polyA
568 selection. This procedure removed ~4 % of the reads. Cleaned reads were used for transcriptome
569 assembly with the rnaSPAdes v3.13.0 using k-mer size of 55bp (Bushmanova *et al.*, 2019).

570

571 **Annotation of organellar genomes and other sequence analyses**

572 Initial annotation of both the plastid and mitochondrial genomes of the strain AMAZONIE were
573 obtained using MFannot (http://megasun.bch.umontreal.ca/cgi-bin/dev_mfa/mfannotInterface.pl).
574 The program output was carefully checked manually, primarily by relying on BLAST searches, to

575 find possible missed genes, to validate or correct the assessment of the initiation codons, to fix the
576 delimitation of introns, and to ensure that all genes were properly named. ORFs lacking discernible
577 homologs (as assessed by HHpred; Zimmermann *et al.*, 2018) encoding proteins shorter than 150
578 amino acid residues and ORFs consisting mostly of sequence repeats were omitted from the
579 annotation. Distribution of repeats within the organellar genomes and comparison of repeats
580 between organellar genomes of *L. pallida* and other selected chlamydomonadales were analysed
581 using the dottup programme from the EMBOSS package ([http://www.bioinformatics.nl/cgi-](http://www.bioinformatics.nl/cgi-bin/emboss/dottup)
582 [bin/emboss/dottup](http://www.bioinformatics.nl/cgi-bin/emboss/dottup)). Detailed analyses of imperfect palindromes and G-quadruplexes were
583 performed using the Palindrome analyzer (Brázda *et al.*, 2016) and the G4hunter web-based server
584 (Brázda *et al.*, 2019). The Palindrome analyzer was used to search for motifs 8-100 bp in length
585 with spacers 0-10 bp, and a maximum of one mismatch in the palindrome. The G4hunter web-
586 based server was used under the default settings, i.e., window=25 and threshold=1.2.

587 To understand the position of amino acid stretches encoded by the characteristic repeats
588 that have invaded the coding sequence of the *ftsH* gene, the tertiary structure of the encoded protein
589 was predicted by homology modelling using the Phyre2 program
590 (<http://www.sbg.bio.ic.ac.uk/phyre2/html/page.cgi?id=index>; Kelley *et al.*, 2015). The secondary
591 structure of the ITS2 region was modelled manually according to the consensus secondary ITS2
592 structure of two green algae (Caisová *et al.*, 2013), visualised by VARNA software (Darty *et al.*,
593 2009), and manually edited in a graphical editor. Homologs of nucleus-encoded plastidial proteins
594 of specific interest were searched in the *L. pallida* transcriptome assembly by using TBLASTN
595 and the respective proteins sequences from *Arabidopsis thaliana* or *C. reinhardtii* (selected based
596 on the information from the literature or keyword database searches). Significant hits (E-value
597 $\leq 1e-5$) were evaluated by BLASTX searches against the NCBI non-redundant protein sequence

598 database to filter out bacterial contaminants and sequences corresponding to non-orthologous
599 members of broader protein families. Subcellular localization (for complete sequences only) was
600 assessed by using TargetP-2.0 (<https://services.healthtech.dtu.dk/service.php?TargetP-2.0>;
601 Almagro Armenteros *et al.*, 2019) and PredAlgo ([http://lobosphaera.ibpc.fr/cgi-
602 bin/predalgotdb2.perl?page=main](http://lobosphaera.ibpc.fr/cgi-bin/predalgotdb2.perl?page=main); Tadrif *et al.*, 2012).

603

604 **Phylogenetic analyses**

605 Multiple sequence alignment of 18S rRNA gene relying on a total of 201 chlorophyte OTUs was
606 computed using MAFFT v7 (Kato *et al.*, 2019) and trimmed manually. The 18S rRNA sequence
607 from *Polytoma oviforme* available in GenBank (U22936.1) was proposed to be chimeric (Nakada
608 *et al.*, 2008), but given the relevance of this organism for our analysis, we included it, masking the
609 regions putatively derived from a different source by strings of N. Maximum likelihood tree
610 inference was performed using IQ-TREE multicore v1.6.12 (Nguyen *et al.*, 2015) under
611 TIM2+F+I+G4 model with 100 non-parametric bootstrap replicates. For multigene analysis,
612 alignments of conserved plastome-encoded proteins used previously (Fučíková *et al.*, 2019) were
613 updated by adding the respective homologs from *L. pallida* and thirteen additional relevant
614 chlorophycean taxa not represented in the initial dataset. On the other hand, sequences representing
615 the OCC clade of Chlorophyceae (evidently only distantly related to *L. pallida* based on the 18S
616 rRNA gene phylogeny and morphological features) to keep the size of the dataset easier to analyse
617 with a complex substitution model. For the final matrix, a subset of 24 proteins (all having their *L.*
618 *pallida* representative) were used. Multiple alignments of the homologous amino acid sequences
619 were built using MAFFT v7.407 with the L-INS-i algorithm (Kato & Standley, 2013) and
620 manually trimmed to exclude unreliably aligned regions. The final concatenated matrix comprised

621 5,020 amino acid residues. The tree was built using PhyloBayes v4.1 (Lartillot *et al.*, 2013) under
622 the CAT+GTR model of sequence evolution, with two independent chains that converged at
623 15,298 generations with the largest discrepancy in posterior probabilities (PPs) (maxdiff) of
624 0.0535238 (at burn-in of 20%). The maximum likelihood (ML) tree was inferred with IQ-TREE
625 multicore v1.6.12 using the LG+C60+F+G4 substitution model. Statistical support was assessed
626 with 100 IQ-TREE non-parametric bootstraps with correction and PhyloBayes posterior
627 probabilities.

628

629 **Supplementary information**

630 **Supplementary information** accompanies this paper at XXXXXXXX.

631

632 **Additional file 1: Supplementary Figs S1-S9. Fig. S1.** Maximum likelihood phylogenetic tree
633 (RAxML, GTRGAMMA+I substitution model) of 18S rRNA gene sequences from
634 Chlorophyceae. **Fig. S2.** Predicted secondary structure of the ITS2 region of *Leontynka pallida*,
635 with differences in the corresponding region of *Leontynka elongata* mapped onto it. **Fig. S3.**
636 *Leontynka pallida* under the light microscope. **Fig. S4.** *Leontynka elongata* under the light
637 microscope. **Fig. S5.** Ultrastructure of *Leontynka elongata* (a–f) and *Leontynka pallida* (g–i). **Fig.**
638 **S6.** Occurrence of the “variant 8” repeat in the FtsH protein of *Leontynka pallida* mapped on its
639 predicted structure. **Fig. S7.** Alignment of the highly similar terminal regions of the originally
640 assembled linear mitogenome contig. **Fig. S8.** Occurrence of the “variant 8” repeat (translated in
641 reading frame +0 as KDKPANLTS and -0 as KEVSFAGLSL) in variable region of protein
642 sequence of the ribosomal protein Rps8 from *Leontynka pallida* (full protein alignment together
643 with representatives of other chlamydomonadalean algae).

644

645 **Additional file 2: Supplementary Notes S1-S5 and supplementary Methods S1-S4. Note S1.**

646 Taxonomic descriptions. **Note S2.** Further details on the morphology and ultrastructure of

647 *Leontynka* spp. **Note S3.** Further details on various kinds of repeats in the plastome of *L. pallida*.

648 **Note S4.** Further details on the repeat insertions in *L. pallida* plastid coding sequences. **Note S5.**

649 Differential diagnosis of *Leontynka* spp. with regard to previously described colourless

650 chlamydomonadalean taxa. **Methods S1.** Isolation and cultivation of strains. **Methods S2.** Light

651 and transmission electron microscopy. **Methods S3.** Amplification and sequencing of 18S and ITS

652 rDNA regions. **Methods S4.** DNA and RNA isolation.

653

654 **Additional file 3: Supplementary Tables S1-S7. Table S1.** Nuclear transcripts from *Leontynka*

655 *pallida* specifically discussed in the paper. **Table S2.** Comparison of GC content, number of

656 imperfect palindromes, and potential quadruplex-forming sequences in selected organellar

657 genomes. **Table S3.** Strong codon usage bias in the mitochondrial genome of *Leontynka pallida*.

658 **Table S4.** Relative frequency of amino acids with GC-rich codons (G, A, R, P) in proteins encoded

659 by different mitogenomes. **Table S5.** Relative frequency of codons in plastid genes of *Leontynka*

660 *pallida*. **Table S6.** Relative frequency of amino acids in proteins encoded by the plastome of

661 *Leontynka pallida*. **Table S7.** The most abundant imperfect palindrome in the *Leontynka pallida*

662 plastome that is missing in exons.

663

664 **Declarations**

665 **Ethics approval and consent to participate**

666 Not applicable.

667

668 **Consent for publication**

669 Not applicable.

670

671 **Competing interests**

672 The authors declare no competing interests.

673

674 **Availability of data and materials**

675 Sequences determined in this study are available from GenBank with the following accession
676 numbers: ##### and ##### – partial nuclear rDNA sequences (18S rRNA-ITS1-5.8S rRNA-ITS2)
677 from *L. pallida* and *L. elongata*, respectively; ##### and ##### – plastid and mitochondrial genome
678 sequence from *L. pallida*; ##### – transcriptome assembly from *L. pallida*. The cultures of
679 *Leontynka* spp. investigated in this study are available upon request.

680

681 **Funding**

682 This work was supported by the Czech Science Foundation project 17-21409S (to M.E.) and the
683 project “CePaViP”, supported by the European Regional Development Fund, within the
684 Operational programme for Research, Development and Education
685 (CZ.02.1.01/0.0/0.0/16_019/0000759). TP was also supported by the Charles University (UNCE
686 204069).

687

688 **Authors' contribution**

689 TP, DB, IČ and ME conceived the original research plans; TP, SCT, KZ, EZ, and KJ obtained
690 nucleic acids for sequencing; SCT and TP obtained Oxford Nanopore data and generated
691 organellar genome assemblies; EZ and IČ isolated the strains; TP and DB carried out the
692 morphological characterisation of the strains; DB and NY obtained the TEM data; MS assembled
693 the transcriptome; TP, DB, and TŠ carried out phylogenetic analyses; TP, IČ and ME analysed
694 and annotated the organellar genome sequences; IČ and ME supervised the work of junior
695 researchers and obtained funding; TP, DB and ME drafted the manuscript; all authors contributed
696 to the final version of the text; ME agreed to serve as the author responsible for contact and
697 ensuring communication.

698

699 **Acknowledgements**

700 We are thankful to Karolina Fučíková for providing alignments of chlorophycean plastid proteins
701 and Petr Janšta for collecting the sample MBURUCU.

702

703 **References**

704 **Almagro Armenteros JJ, Salvatore M, Emanuelsson O, Winther O, von Heijne G, Elofsson**

705 **A, Nielsen H. 2019.** Detecting sequence signals in targeting peptides using deep learning.

706 *Life Science Alliance* **2**: e201900429.

707 **Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997.**

708 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.

709 *Nucleic Acids Research* **25**: 3389–3402.

710 **Asmail SR, Smith DR. 2016.** Retention, erosion, and loss of the carotenoid biosynthetic pathway

711 in the nonphotosynthetic green algal genus *Polytomella*. *New Phytologist* **209**: 899–903.

- 712 **Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,**
713 **Nikolenko SI, Pham S, Prjibelski AD *et al.* 2012.** SPAdes: a new genome assembly
714 algorithm and its applications to single-cell sequencing. *Journal of Computational*
715 *Biology* **19**: 455–477.
- 716 **Bauman N, Akella S, Hann E, Morey R, Schwartz AS, Brown R, Richardson TH. 2018.** Next-
717 generation sequencing of *Haematococcus lacustris* reveals an extremely large 1.35-
718 megabase chloroplast genome. *Genome Announcements* **6**: e00181–e001818.
- 719 **Böhm M, Kreimer G. 2020.** Orient in the World with a Single Eye: The Green Algal Eyespot and
720 Phototaxis. In: Cánovas FM, Lüttge U, Risueño MC, Pretzsch H (eds) Progress in Botany
721 Vol. 82. Progress in Botany, vol 82. Springer, Cham.
- 722 **Bolger AM, Lohse M, Usadel B. 2014.** Trimmomatic: a flexible trimmer for Illumina sequence
723 data. *Bioinformatics* **30**: 2114–2120.
- 724 **Boer PH, Gray MW. 1988.** Scrambled ribosomal RNA gene pieces in *Chlamydomonas*
725 *reinhardtii* mitochondrial DNA. *Cell* **55**: 399–411.
- 726 **Brázda V, Kolomazník J, Lýsek J, Bartas M, Fojta M, Št'astný J, Mergny J. 2019.** G4Hunter
727 web application: a web server for G-quadruplex prediction. *Bioinformatics* **35**: 3493–3495.
- 728 **Brázda V, Kolomazník J, Lýsek J, Hároníková L, Coufal J, Št'astný J. 2016.** Palindrome
729 analyser - A new web-based server for predicting and evaluating inverted repeats in
730 nucleotide sequences. *Biochemical and Biophysical Research Communications* **478**: 1739–
731 1745.
- 732 **Burge S, Parkinson GN, Hazel P, Todd AK, Neidle S. 2006.** Quadruplex DNA: sequence,
733 topology and structure *Nucleic Acids Research* **34**: 5402–5415.

- 734 **Bushmanova E, Antipov D, Lapidus A, Prjibelski AD. 2019.** rnaSPAdes: a de novo
735 transcriptome assembler and its application to RNA-Seq data. *GigaScience* **8**: giz100.
- 736 **Caisová L, Marin B, Melkonian M. 2013.** A consensus secondary structure of ITS2 in the
737 Chlorophyta identified by phylogenetic reconstruction. *Protist* **164**: 482–496.
- 738 **Coleman AW. 2000.** The significance of a coincidence between evolutionary landmarks found in
739 mating affinity and a DNA sequence. *Protist* **151**: 1–9.
- 740 **Darty K, Denise A, Ponty Y. 2009.** VARNA: Interactive drawing and editing of the RNA
741 secondary structure. *Bioinformatics* **25**: 1974–1975.
- 742 **Del Vasto M, Figueroa-Martinez F, Featherston J, González MA, Reyes-Prieto A, Durand
743 PM, Smith DR. 2015.** Massive and widespread organelle genomic expansion in the green
744 algal genus *Dunaliella*. *Genome Biology and Evolution* **7**: 656–663.
- 745 **Denovan-Wright EM, Lee RW. 1994.** Comparative structure and genomic organization of the
746 discontinuous mitochondrial ribosomal RNA genes of *Chlamydomonas eugametos* and
747 *Chlamydomonas reinhardtii*. *Journal of Molecular Biology* **241**: 298–311.
- 748 **Dudkina NV, Oostergetel GT, Lewejohann D, Braun HP, Boekema EJ. 2010.** Row-like
749 organization of ATP synthase in intact mitochondria determined by cryo-electron
750 tomography. *Biochimica et Biophysica Acta (BBA)-Bioenergetics* **1797**: 272–277.
- 751 **Ettl H. 1983.** Chlorophyta I: Phytomonadina. In: Ettl H, Gerloff J, Heynig H, Mollenhauer D
752 (eds.), Süßwasserflora von Mitteleuropa. Band 9. Stuttgart: Gustav Fischer Verlag. 807 p.
- 753 **Falabella M, Fernandez RJ, Johnson FB, Kaufman BA. 2019.** Potential roles for G-
754 quadruplexes in mitochondria. *Current Medicinal Chemistry* **26**: 2918–2932.
- 755 **Fan J, Schnare MN, Lee RW. 2003.** Characterization of fragmented mitochondrial ribosomal
756 RNAs of the colorless green alga *Polytomella parva*. *Nucleic Acids Research* **31**: 769–778.

- 757 **Figuroa-Martinez F, Nedelcu AM, Smith DR, Reyes-Prieto A. 2015.** When the lights go out:
758 the evolutionary fate of free-living colorless green algae. *New Phytologist* **206**: 972–982.
- 759 **Figuroa-Martinez F, Nedelcu AM, Smith DR, Reyes-Prieto A. 2017.** The plastid genome of
760 *Polytoma uvella* is the largest known among colorless algae and plants and reflects
761 contrasting evolutionary paths to nonphotosynthetic lifestyles. *Plant Physiology* **173**: 932–
762 943.
- 763 **Fučíková K, Lewis PO, Neupane S, Karol KG, Lewis LA. 2019.** Order, please! Uncertainty in
764 the ordinal-level classification of Chlorophyceae. *PeerJ* **7**: e6899.
- 765 **Gaffal KP, Schneider GJ. 1980.** Numerical, morphological and topographical heterogeneity of
766 the chondriome during the vegetative life cycle of *Polytoma papillatum*. *Journal of Cell*
767 *Science* **46**: 299–312.
- 768 **Gaouda H, Hamaji T, Yamamoto K, Kawai-Toyooka H, Suzuki M, Noguchi H, Minakuchi**
769 **Y, Toyoda A, Fujiyama A, Nozaki H et al. 2018.** Exploring the limits and causes of plastid
770 genome expansion in volvocine green algae. *Genome Biology and Evolution* **10**: 2248–2254.
- 771 **Greiner S, Lehwark P, Bock R. 2019.** OrganellarGenomeDRAW (OGDRAW) version 1.3.1:
772 expanded toolkit for the graphical visualization of organellar genomes. *Nucleic Acids*
773 *Research* **47**: W59–W64
- 774 **Gruber HE, Rosario B. 1974.** Variation in eyespot ultrastructure in *Chlamydomonas reinhardi*
775 (ac-31). *Journal of Cell Science* **15**: 481–494.
- 776 **Hadariová L, Vesteg M, Hampl V, Krajčovič J. 2018.** Reductive evolution of chloroplasts in
777 non-photosynthetic plants, algae and protists. *Current Genetics* **64**: 365–387.

- 778 **Hänsel-Hertsch R, Beraldi D, Lensing SV, Marsico G, Zyner K, Parry A, Di Antonio M, Pike**
779 **J, Kimura H, Narita M *et al.* 2016.** G-quadruplex structures mark human regulatory
780 chromatin. *Nature Genetics* **48**: 1267–1272.
- 781 **Hecht J, Grewe F, Knoop V. 2011.** Extreme RNA editing in coding islands and abundant
782 microsatellites in repeat sequences of *Selaginella moellendorffii* mitochondria: The root of
783 frequent plant mtDNA recombination in early tracheophytes. *Genome Biology and Evolution*
784 **3**: 344–358.
- 785 **Herron MD, Hackett JD, Aylward FO, Michod RE. 2009.** Triassic origin and early radiation of
786 multicellular volvocine algae. *Proceedings of the National Academy of Sciences, USA* **106**:
787 3254–3258.
- 788 **Kamikawa R, Tanifuji G, Ishikawa SA, Ishii K, Matsuno Y, Onodera NT, Ishida K,**
789 **Hashimoto T, Miyashita H, Mayama S *et al.* 2015.** Proposal of a twin arginine translocator
790 system-mediated constraint against loss of ATP synthase genes from nonphotosynthetic
791 plastid genomes. [Corrected]. *Molecular Biology and Evolution* **32**: 2598–2604.
- 792 **Katoh K, Rozewicki J, Yamada KD. 2019.** MAFFT online service: multiple sequence alignment,
793 interactive sequence choice and visualization. *Briefings in Bioinformatics* **20**: 1160–1166.
- 794 **Katoh K, Standley DM. 2013.** MAFFT multiple sequence alignment software version 7:
795 improvements in performance and usability. *Molecular Biology and Evolution* **30**: 772–780.
- 796 **Kayama M, Chen JF, Nakada T, Nishimura Y, Shikanai T, Azuma T, Miyashita H, Takaichi**
797 **S, Kashiyama Y, Kamikawa R. 2020.** A non-photosynthetic green alga illuminates the
798 reductive evolution of plastid electron transport systems. *BMC Biology* **18**: 126.
- 799 **Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJ. 2015.** The Phyre2 web portal for
800 protein modeling, prediction and analysis. *Nature Protocols* **10**: 845–858.

- 801 **Kim D, Langmead B, Salzberg S. 2015.** HISAT: a fast spliced aligner with low memory
802 requirements. *Nature Methods* **12**: 357–360.
- 803 **Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. 2017.** Canu: scalable
804 and accurate long-read assembly via adaptive k-mer weighting and repeat
805 separation. *Genome Research* **27**: 722–736.
- 806 **Kreimer G. 1994.** Cell biology of phototaxis in flagellate algae. *International Review of Cytology*
807 **148**: 229–310.
- 808 **Lamb MR, Dutcher SK, Worley CK, Dieckmann CL. 1999.** Eyespot-assembly mutants in
809 *Chlamydomonas reinhardtii*. *Genetics* **153**: 721–729.
- 810 **Lang NJ. 1963.** Electron microscopic demonstration of plastids in *Polytoma*. *Journal of*
811 *Protozoology* **10**: 333–339.
- 812 **Lavrov DV, Adamski, M, Chevaldonne, P, Adamska, M. 2016.** Extensive mitochondrial
813 mRNA editing and unusual mitochondrial genome organization in calcarean
814 sponges. *Current Biology* **26**: 86-92.
- 815 **Lartillot N, Rodrigue N, Stubbs D, Richer J. 2013.** PhyloBayes MPI: phylogenetic
816 reconstruction with infinite mixtures of profiles in a parallel environment. *Systematic*
817 *Biology* **62**: 611–615.
- 818 **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin**
819 **R, 1000 Genome Project Data Processing Subgroup. 2009.** The sequence alignment/map
820 format and SAMtools. *Bioinformatics* **25**: 2078–2079.
- 821 **Li H. 2013.** Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
822 arXiv preprint arXiv:1303.3997.

- 823 **Loman NJ, Quick J, Simpson JT. 2015.** A complete bacterial genome assembled de novo using
824 only nanopore sequencing data. *Nature Methods* **12**: 733–735.
- 825 **Mechela A, Schwenkert S, Soll J. 2019.** A brief history of thylakoid biogenesis. *Open Biology* **9**:
826 180237.
- 827 **Milne I, Bayer M, Stephen G, Cardle L, Marshall D. 2016.** Tablet: visualizing next-generation
828 sequence assemblies and mappings. *Methods in Molecular Biology* **1374**: 253–268.
- 829 **Muñoz-Gómez SA, Mejía-Franco FG, Durnin K, Colp M, Grisdale CJ, Archibald JM,
830 Slamovits CH. 2017.** The new red algal subphylum Proteorhodophytina comprises the
831 largest and most divergent plastid genomes known. *Current Biology* **27**: 1677–1684.e4.
- 832 **Nakada T, Misawa K, Nozaki H. 2008.** Molecular systematics of Volvocales (Chlorophyceae,
833 Chlorophyta) based on exhaustive 18S rRNA phylogenetic analyses. *Molecular
834 Phylogenetics and Evolution* **48**: 281–291.
- 835 **Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. 2015.** IQ-TREE: a fast and effective
836 stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology
837 and Evolution* **32**: 268–274.
- 838 **Petroutsos D. 2017.** Chlamydomonas Photoreceptors: Cellular Functions and Impact on
839 Physiology. In: Hippler M. (eds) Chlamydomonas: Biotechnology and Biomedicine.
840 Microbiology Monographs, vol 31. Springer, Cham.
- 841 **Schmidt M, Luff M, Mollwo A, Kaminski M, Mittag M, Kreimer G. 2007.** Evidence for a
842 specialized localization of the chloroplast ATP-synthase subunits α , β , and γ in the eyespot
843 apparatus of *Chlamydomonas reinhardtii* (Chlorophyceae). *Journal of Phycology* **43**: 284–
844 294.

- 845 **Schünemann D. 2007.** Mechanisms of protein import into thylakoids of chloroplasts. *Biological*
846 *Chemistry* **388**: 907–915.
- 847 **Sibbald SJ, Archibald JM. 2020.** Genomic insights into plastid evolution. *Genome Biology and*
848 *Evolution* **12**: 978–990.
- 849 **Siu C, Swift H, Chiang K. 1976.** Characterization of cytoplasmic and nuclear genomes in the
850 colorless alga *Polytoma*. *Journal of Cell Biology* **69**: 352–370.
- 851 **Skalitzky CA, Martin JR, Harwood JH, Beirne JJ, Adamczyk BJ, Heck GR, Cline K,**
852 **Fernandez DE. 2011.** Plastids contain a second sec translocase system with essential
853 functions. *Plant Physiology* **155**: 354–369.
- 854 **Smith, DR. 2012.** Updating our view of organelle genome nucleotide landscape. *Frontiers in*
855 *Genetics* **3**: 175.
- 856 **Smith DR, Hua J, Lee RW. 2010.** Evolution of linear mitochondrial DNA in three known
857 lineages of *Polytomella*. *Current Genetics* **56**: 427–438.
- 858 **Smith DR, Hua J, Archibald JM, Lee RW. 2013.** Palindromic genes in the linear mitochondrial
859 genome of the nonphotosynthetic green alga *Polytomella magna*. *Genome Biology and*
860 *Evolution* **5**: 1661–1667.
- 861 **Smith DR, Lee RW. 2014.** A plastid without a genome: evidence from the nonphotosynthetic
862 green algal genus *Polytomella*. *Plant Physiology* **164**: 1812–1819.
- 863 **Smith DR. 2018.** *Haematococcus lacustris*: the makings of a giant-sized chloroplast genome. *AoB*
864 *Plants* **10**: ply058.
- 865 **Smith DR. 2020a.** Can green algal plastid genome size be explained by DNA repair mechanisms?
866 *Genome Biology and Evolution* **12**: 3797–3802.

- 867 **Smith DR. 2020b.** Common repeat elements in the mitochondrial and plastid genomes of green
868 algae. *Frontiers in Genetics* **11**: 465.
- 869 **Smith DR, Craig, R. J. 2021.** Does mitochondrial DNA replication in *Chlamydomonas* require a
870 reverse transcriptase? *New Phytologist* **229**: 1192–1195.
- 871 **Suzuki S, Endoh R, Manabe RI, Ohkuma M, Hirakawa Y. 2018.** Multiple losses of
872 photosynthesis and convergent reductive genome evolution in the colourless green algae
873 *Prototheca*. *Scientific Reports* **8**: 940.
- 874 **Tardif M, Atteia A, Specht M, Cogne G, Rolland N, Brugière S, Hippler M, Ferro M, Bruley
875 C, Peltier G, Vallon O, Cournac L. 2012.** PredAlgo: a new subcellular localization
876 prediction tool dedicated to green algae. *Molecular Biology and Evolution* **29**: 3625–3639.
- 877 **Turmel M, Lemieux C. 2018.** Chapter six - Evolution of the plastid genome in green algae.
878 *Advances in Botanical Research* **85**: 157–193.
- 879 **Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,
880 Wortman J, Young SK et al. 2014.** Pilon: an integrated tool for comprehensive microbial
881 variant detection and genome assembly improvement. *PLoS One* **9**: e112963.
- 882 **Webster DA, Hackett DP, Park RB. 1968.** The respiratory chain of colorless algae: III. Electron
883 microscopy. *Journal of Ultrastructural Research* **21**: 514–523.
- 884 **Wick RR, Judd LM, Gorrie CL, Holt KE. 2017.** Unicycler: resolving bacterial genome
885 assemblies from short and long sequencing reads. *PLoS Computational Biology* **13**:
886 e1005595.
- 887 **Wolf M, Chen S, Song J, Ankenbrand M, Müller T. 2013.** Compensatory base changes in ITS2
888 secondary structures correlate with the biological species concept despite intragenomic
889 variability in ITS2 sequences--a proof of concept. *PLoS One* **8**: e66726.

890 **Zhang X, Bauman N, Brown R, Richardson TH, Akella S, Hann E, Morey R, Smith DR.**

891 **2019.** The mitochondrial and chloroplast genomes of the green alga *Haematococcus* are

892 made up of nearly identical repetitive sequences. *Current Biology* **29**: R736–R737.

893 **Ziehe D, Dünschede B, Schünemann D. 2017.** From bacteria to chloroplasts: evolution of the

894 chloroplast SRP system. *Biological Chemistry* **398**: 653–661.

895 **Zimmermann L, Stephens A, Nam SZ, Rau D, Kübler J, Lozajic M, Gabler F, Söding J,**

896 **Lupas AN, Alva V. 2018.** A completely reimplemented MPI Bioinformatics Toolkit with a

897 new HHpred server at its core. *Journal of Molecular Biology* **430**: 2237–2243.

898

899

900

901

902

903

904

905

906

907

908

909

910

911 **Table 1.** Basic characteristics of plastomes of *Leontynka pallida*, selected other non-
 912 photosynthetic chlorophytes and the photosynthetic *Chlamydomonas reinhardtii* and
 913 *Haematococcus lacustris* for comparison.

914

species	plastomes		genes		coding DNA**		intergenic regions		
	accession number	size (bp)	protein-coding*	tRNAs	GC (%)	total length (%)	GC (%)	average length (bp)	total length (%)
<i>Leontynka pallida</i>	#####	362,307	32	26	35	18	37	4,726	80
<i>Polytoma uwella</i>	KX828177.1	230,207	25	27	40	20	19	2,998	68
Volvocales sp. NrCI902	LC516060.1	176,432	25	30	36.1	32.6	42.8	1,944	60.9
<i>Helicosporidium</i> sp.	DQ398104.1	37,454	26	25	27.6	95.4	13.4	33.8	4.4
<i>Chlamydomonas reinhardtii</i>	NC_005353.1	203,828	69	29	34.6	43	34	923	48.5
<i>Haematococcus lacustris</i> (†)	NC_037007.1	1,352,306	63	28	38.6	6.8	51	12,326	91.15

915

916 *does not include ORFs inside introns

917 **does not include introns and ORFs inside introns, but includes tRNAs and rRNAs

918 (†) since the annotation of the *H. lacustris* plastome available in the respective GenBank record is highly incomplete, the values presented are based
 919 on a reannotation obtained by using MFannot (with a genetic code translating UGA as Trp, based on the insights by Smith, 2018).

920

921

922

923

924

925

926

927

928 **Figure legends**

929 **Fig. 1** Phylogenetic position of the new genus of non-photosynthetic green algae, *Leontynka*. Non-
930 photosynthetic taxa in Chlamydomonadales are highlighted in violet. (a) Maximum likelihood
931 phylogenetic tree (IQ-TREE, TIM2+F+I+G4 substitution model) of 18S rRNA gene sequences
932 from Chlamydomonadales and related chlorophytes. Non-parametric bootstrap support values
933 calculated from 100 replicates are shown when ≥ 50 . Previously demarcated main clades (Nakada
934 *et al.*, 2008) are collapsed and the outgroup taxa are not shown for simplicity. The full version of
935 the tree is provided as Additional file 1: Fig. S1. (b) Phylogenetic analysis of a concatenated dataset
936 of 24 conserved plastome-encoded proteins (5,020 amino acid positions) from
937 Chlamydomonadales, including *Leontynka pallida*, and the sister order Sphaeropleales (*sensu lato*;
938 see Fučíková *et al.*, 2019). The tree topology was inferred using PhyloBayes (CAT+GTR
939 substitution model), branch support values correspond to posterior probability (from PhyloBayes)
940 / maximum likelihood bootstrap analysis (IQ-TREE, LG+C60+F+G4 substitution model, 100 non-
941 parametric bootstrap replicates). Black dots represent full support obtained with both methods,
942 asterisks denote bootstrap support values < 50 .

943
944 **Fig. 2** Light and transmission electron microscopy of *Leontynka pallida* (a–e) and *Leontynka*
945 *elongata* (f–j). Note the difference in the cell shape between the two species and the presence of a
946 single (a, f, h) or two (b, d, g, i, j) eyespots. Lipid droplets were also detected within the cells (d,
947 h, j). Abbreviations: bb – basal body; cv – contractile vacuole; e – eyespot; L – lipid droplet; m –
948 mitochondrion; N – nucleus; n – nucleolus; s – starch. Arrows point to thylakoids; asterisks mark
949 plastid “bridges” connecting separate compartments; arrowheads show membranous inclusions.
950 Scale bars: a-c, f-g = 10 μm ; d, i = 0.2 μm ; e = 0.5 μm ; j = 1 μm .

951
952 **Fig. 3** Characteristics of the organellar genomes of *Leontynka pallida*. (a) Gene map of the *L.*
953 *pallida* plastid genome. Genes are shown as squares (coloured according to the functional
954 category; see the graphical legend in the left bottom corner) on the inner or outer side of the outer
955 circle depending on their orientation (transcription in the clockwise or counter-clockwise direction,
956 respectively; see the grey arrows). Genes marked with an asterisk contain introns. The inner circle
957 represents a GC content plot. (b) Sequence self-similarity plot of the *L. pallida* plastome (ptDNA).
958 (c) Map of the *L. pallida* mitochondrial genome. The display convention is the same as for the
959 plastid genome. (d) Sequence self-similarity plot of the *L. pallida* mitochondrial genome
960 (mtDNA). (e) *L. pallida* plastome-mitogenome similarity plot. All similarity plots were generated
961 by using the word size of 15 bp and black dots represent the occurrence of the same word at the
962 places compared. The organellar genome maps (a, c) were visualised by using OGDRAW v1.3.1
963 (Greiner *et al.*, 2019).

964
965 **Fig. 4** Distribution of repeats in organellar genomes of *Leontynka pallida*. (a) Most abundant
966 imperfect palindromes and their characteristics. The “Spacer” corresponds to the presumed loop
967 separating the palindromic regions presumably pairing to form a stem structure. The “Mismatch”
968 column indicates the number of positions that deviate from a perfect palindrome. The occurrence
969 of the repeats is given for the plastome (ptDNA) and mitogenome (mtDNA), with the number of
970 cases indicated for the whole organellar genome and separately for exons in protein-coding genes.
971 In two cases of mitogenome repeats, two variants – a shorter and a longer – are considered, with
972 the latter indicated in parentheses. (b) Distribution of the imperfect palindrome
973 AAGCCAGC|NNN|GCTGACTT and its most common variants within exons of the plastome. The

974 numbers show the abundance of the given repeat in direct / reverse complement orientation
975 (relative to the coding sequence). In the case of “variant 1”, the repeat has the same sequence in
976 both directions, so only one number per gene is presented. Note that the variants considered are
977 not mutually exclusive alternatives but correspond to nested categories with a different degree or
978 relaxation of the sequence pattern. (c) Characterisation of repeats from (b) and their abundance in
979 various regions of the plastome and the mitogenome of *L. pallida* as well as other plastomes of
980 Chlamydomonadales deposited in NCBI databases. The numbers show the abundance of the given
981 repeat in direct / reverse complement orientation (relative to the coding sequence in the case of
982 exons, or relative to the DNA strand corresponding to the reference organellar sequence in case of
983 the values for the whole organellar genome). (d) Occurrence of the “variant 8” repeat (translated
984 in the reading frame +0 as KDKPANLTS) in a variable region of the ribosomal protein Rps8
985 (detail; the full alignment is available as Additional file 1: Fig. S8). (e) Occurrence of the “variant
986 4” repeat in protein-coding sequences and its translation in all six reading frames. The category of
987 rare codons (“rare 2%”) is defined as the sum of the least used codons together representing less
988 than 2% of all codons in the plastome (100% = 19,899 codons); the categories of the 4%, 10%,
989 and 20% rarest codons and more than 50% of most frequent codons are defined similarly (listed
990 in Additional file 3: Table S5). The numbers indicated for the “codon usage” correspond to the
991 minimal number of the codons of the respective category present in the respective reading frame,
992 with the “max X” numbers indicating the maximal number of such codons depending on the actual
993 nucleotide sequence of the degenerated “variant 4” repeat. Note that in some cases the theoretical
994 maximal number is not observed in the actual *L. pallida* plastid gene sequences (see the asterisks).
995 The column “Rare AA” indicates the occurrence of amino acids belonging to the category of amino
996 acids generally rarely used in plastome-encoded proteins in *L. pallida* (see Additional file 3: Table

997 S6). The occurrence of the repeat variants indicated for coding sequences (CDS) correspond to
998 their occurrence as counted at the nucleotide level, whereas the occurrence in proteins is counted
999 at the amino acid sequence level (and may be higher due to different nucleotide sequences
1000 encoding the same amino acid sequence). The analysis of intraexonic repeat insertions is
1001 commented in more detail in Additional file 2: Note S4.







