

1 High quality genome assembly and annotation (v1) of the
2 eukaryotic terrestrial microalga *Coccomyxa viridis* SAG 216-4

3

4 Anton Kraege¹, Edgar A. Chavarro-Carrero¹, Nadège Guiguelmoni², Eva Schnell¹, Joseph
5 Kirangwa², Stefanie Heilmann-Heimbach^{3,4}, Kerstin Becker^{5,6}, Karl Köhrer⁶, WGGC Team,
6 DeRGA Community, Philipp Schiffer², Bart P. H. J. Thomma^{1,7} and Hanna Rovenich¹

7

8 ¹ Institute for Plant Sciences, University of Cologne, Germany; ² Institute of Zoology, University
9 of Cologne, Germany; ³ Institute of Human Genetics, University Hospital of Bonn & University
10 of Bonn, Bonn, Germany; ⁴ NGS Core Facility, Medical Faculty of the University of Bonn, Bonn,
11 Germany; ⁵ Cologne Center for Genomics (CCG), Medical Faculty, University of Cologne,
12 50931 Cologne, Germany; ⁶ Biological and Medical Research Centre (BMFZ), Genomics &
13 Transcriptomics Laboratory, Heinrich-Heine-University Düsseldorf, Germany; ⁷ Cluster of
14 Excellence on Plant Sciences (CEPLAS), Germany

15

16 corresponding author: hanna.roevenich@uni-koeln.de

17

18

19 Abstract

20 Unicellular green algae of the genus *Coccomyxa* are recognized for their worldwide
21 distribution and ecological versatility. Most species described to date live in close association
22 with various host species, such as in lichen associations. However, little is known about the
23 molecular mechanisms that drive such symbiotic lifestyles. We generated a high-quality
24 genome assembly for the lichen photobiont *Coccomyxa viridis* SAG 216-4 (formerly *C.*
25 *mucigena*). Using long-read PacBio HiFi and Oxford Nanopore Technologies in combination
26 with chromatin conformation capture (Hi-C) sequencing, we assembled the genome into 21
27 scaffolds with a total length of 50.9 Mb, an N50 of 2.7 Mb and a BUSCO score of 98.6%. While
28 19 scaffolds represent full-length nuclear chromosomes, two additional scaffolds represent
29 the mitochondrial and plastid genomes. Transcriptome-guided gene annotation resulted in the
30 identification of 13,557 protein-coding genes, of which 68% have annotated PFAM domains
31 and 962 are predicted to be secreted.

32

33 Keywords: *Coccomyxa viridis*, EBP, ERGA, long-read sequencing, genome assembly,
34 genome annotation, Trebouxiophyceae

35

36

37 Introduction

38 Green algae are photosynthesizing eukaryotic organisms that differ greatly in terms of
39 morphology and colonize a large variety of aquatic and terrestrial habitats. Phylogenetically,
40 green algae form a paraphyletic group that has recently been proposed to comprise three
41 lineages including the Prasinodermophyta in addition to the Chlorophyta and Streptophyta (Li
42 et al., 2020). This new phylum diverged before the split of the Chlorophyta and Streptophyta
43 that occurred between 1,000 and 700 million years ago (Morris et al., 2018). While the
44 streptophyte lineage encompasses charophyte green algae as well as land plants, the
45 chlorophyte lineage consists of 7 prasinophyte classes, which gave rise to 4 phycoplast-
46 containing core chlorophyte classes (Chlorodendrophyceae, Trebouxiophyceae,
47 Ulvophyceae, Chlorophyceae) with one independent sister class (Pedinophyceae) (Leliaert et
48 al., 2012; Marin, 2012).

49 The *Coccomyxa* genus is represented by coccoid unicellular green algae that belong
50 to the class of Trebouxiophyceae. Morphologically, *Coccomyxa* spp. are characterized by
51 irregular elliptical to globular cells that range from 6–14 x 3–6 µm in size, with a single parietal
52 chloroplast lacking pyrenoids and the absence of flagellate stages (Schmidle, 1901). Members
53 of this genus are found in freshwater, marine, and various terrestrial habitats where they occur
54 free-living or in symbioses with diverse hosts (Darienکو et al., 2015; Gustavs et al., 2017;
55 Malavasi et al., 2016). Several *Coccomyxa* species establish stable, mutualistic associations
56 with fungi that result in the formation of complex three-dimensional architectures, known as
57 lichens (Faluaburu et al., 2019; Gustavs et al., 2017; Jaag, 1933; Yahr et al., 2015; Zoller and
58 Lutzoni, 2003). Others associate with vascular plants or lichens as endo- or epiphytes,
59 respectively (Cao et al., 2018a; Cao et al., 2018b; Tagirdzhanova et al., 2023; Trémouillaux-
60 Guiller et al., 2002), and frequently occur on the bark of trees (Kulichová et al., 2014; Štifterová
61 and Neustupa, 2015) where they may interact with other microbes. One novel species was
62 recently found in association with carnivorous plants, even though the nature of this
63 relationship remains unclear (Sciuto et al., 2019). Besides, *Coccomyxa* also establishes
64 parasitic interactions with different mollusk species affecting their filtration ability and
65 reproduction (Gray et al., 1999; Sokolnikova et al., 2016; Sokolnikova et al., 2022; Vaschenko
66 et al., 2013).

67 Despite this ecological versatility, little is known about the molecular mechanisms that
68 determine the various symbiotic lifestyles in *Coccomyxa*. One short read-based genome is
69 available for *C. subellipsoidea* C-169 that was isolated on Antarctica where it occurred on
70 dried algal peat (Blanc et al., 2012), whereas another high-quality genome has recently been
71 made available for a non-symbiotic strain of *C. viridis* that was isolated from a lichen thallus
72 (Tagirdzhanova et al., 2023). For *Coccomyxa* sp. Obi, LA000219 and SUA001 chromosome-

73 , scaffold- and contig-level assemblies are available on NCBI, respectively, as well as two
74 metagenome-assembled genomes of *C. subellipsoidea*. To facilitate the study of *Coccomyxa*
75 symbiont-associated traits and their evolutionary origin, we here present the generation of a
76 high-quality chromosome-scale assembly of the phycobiont *C. mucigena* SAG 216-4 using
77 long-read PacBio HiFi and Oxford Nanopore Technology (ONT) combined with Hi-C and RNA
78 sequencing. Recent SSU and ITS rDNA sequencing-based re-evaluations of the *Coccomyxa*
79 phylogeny placed the SAG 216-4 isolate in the clade of *C. viridis* (Dariencko et al., 2015;
80 Malavasi et al., 2016). Hence, this isolate will be referred to as *C. viridis* here and data have
81 been deposited under the corresponding Taxonomy ID.

82

83 **Materials & Methods**

84 ***Sample information***

85 *Coccomyxa viridis* (formerly *Coccomyxa mucigena*) SAG 216-4 was ordered from the Culture
86 Collection of Algae at the Georg-August-University Göttingen (*Sammlung von Algenkulturen*
87 *der Universität Göttingen*, international acronym SAG), Germany. The stock culture was
88 reactivated in liquid modified Waris-H growth medium (McFadden and Melkonian, 1986) with
89 soil extract and 3x vitamins (0.15 nM vitamin B12, 4.1 nM biotin, 0.3 μM thiamine-HCl, 0.8 nM
90 niacinamide), and maintained through regular medium replacement. Cultures were grown at
91 ~ 15 μmol photons m⁻² s⁻¹ (fluorescent light tubes: L36W/640i energy saver cool white and
92 L58W/956 BioLux, Osram, Munich, Germany) in a 14/10 h light/dark cycle at 20°C.

93

94 ***DNA and RNA extraction***

95 Cells of a 7-week-old *C. viridis* culture were harvested over 0.8 μm cellulose nitrate filters
96 (Sartorius, Göttingen, Germany) using a vacuum pump. Material was collected with a spatula,
97 snap-frozen and ground in liquid nitrogen using mortar and pestle. The ground material was
98 used for genomic DNA extraction with the RSC Plant DNA Kit (Promega, Madison, WI, USA)
99 using the Maxwell® RSC device according to manufacturer's instructions. To prevent shearing
100 of long DNA fragments, centrifugation was carried out at 10,000 g during sample preparation.
101 Following DNA extraction, DNA fragments <10,000 bp were removed using the SRE XS kit
102 (Circulomics, Baltimore, MD, USA) according to manufacturer's instructions. DNA quantity and
103 quality were assessed using the Nanodrop 2000 spectrometer and Qubit 4 fluorometer with
104 the dsDNA BR assay kit (Invitrogen, Carlsbad, CA, USA), and integrity was confirmed by gel
105 electrophoresis. High-molecular weight DNA was stored at 4°C.

106 For total RNA extraction, algal cells were collected from a dense nine-day-old culture
107 and ground in liquid nitrogen using mortar and pestle. RNA was extracted with the Maxwell
108 RSC Plant RNA kit (Promega, Madison, WI, USA) using the Maxwell® RSC device according

109 to manufacturer's instructions. RNA quality and quantity was determined using the Nanodrop
110 2000 and stored at -80°C.

111

112 ***Pacific Biosciences High-Fidelity (PacBio HiFi) sequencing***

113 HiFi libraries were prepared with the Express 2.0 Template kit (Pacific Biosciences, Menlo
114 Park, CA, USA) and sequenced on a Sequel II/Sequel IIe instrument with 30h movie time. HiFi
115 reads were generated using SMRT Link (v10; (Pacific Biosciences, Menlo Park, CA, USA)
116 with default parameters.

117

118 ***Oxford Nanopore Technologies (ONT) sequencing***

119 Library preparation with the Rapid Sequencing Kit (SQK-626 RAD004) was performed with
120 ~400 ng HMW DNA according to manufacturer's instructions (Oxford Nanopore Technologies,
121 Oxford, UK). The sample was loaded onto an R9.4.1 flow cell in a minION Mk1B device
122 (Oxford Nanopore Technologies, Oxford, UK), which was run for 24 h. Subsequent base
123 calling was performed using Guppy (version 630 3.1.3; Oxford Nanopore Technologies,
124 Oxford, UK). Adapter sequences were removed using Porechop (version 0.2.4 with default
125 settings) (Wick, 2018), and the reads were self-corrected and trimmed using Canu (version
126 1.8) (Koren et al., 2017).

127

128 ***Chromosome conformation capture (Hi-C) and sequencing***

129 *C. viridis* cells were cross-linked in 3% formaldehyde for 1 hour at room temperature. The
130 reaction was quenched with glycine at a final concentration of 250 mM. Cells were collected
131 by centrifugation at 16,000 g for 10 min. Pellets were flash-frozen in liquid nitrogen and ground
132 using mortar and pestle. Hi-C libraries were prepared using the Arima-HiC+ kit (Arima
133 Genomics, Carlsbad, CA, USA) according to manufacturer's instructions, and subsequently
134 paired-end (2x150 bp) sequenced on a NovaSeq 6000 instrument (Illumina, San Diego, CA,
135 USA).

136

137 ***RNA sequencing***

138 Library preparation for full-length mRNASeq was performed using the NEB Ultra II Directional
139 RNA Library Prep with NEBNext Poly(A) mRNA Magnetic Isolation Module and 500 ng total
140 RNA as starting material, except for W-RNA Lplaty, where library prep was based on 100 ng
141 total RNA as starting material. Sequencing was performed on an Illumina NovaSeq 6000
142 device with 2x150 bp paired-end sequencing protocol and >50 M reads per sample.

143

144 ***Genome assembly***

145 PacBio HiFi reads were assembled using Raven (v1.8.1) (Vaser and Šikić, 2021) with default
146 settings. Hi-C reads were mapped onto this assembly with Juicer (v2.0) using the “assembly”
147 option to skip the post-processing steps and generate the merged_nodups.txt file (Durand et
148 al., 2016b). For the juicer pipeline, restriction site maps were generated using the *DpnII*
149 (*GATC*) and *Hinfl* (*GANTC*) restriction site profile and the assembly was indexed with BWA
150 index (v0.7.17-r1188) (Li and Durbin, 2009), and used to polish the assembly using 3d-dna
151 (v180922) (Dudchenko et al., 2017). Afterwards, Juicebox (v1.11.08) was used to manually
152 curate the genome assembly by splitting contigs and rearranging them according to the Hi-C
153 pattern (Durand et al., 2016a). Contigs were merged to scaffolds according to the Hi-C map
154 and Ns were introduced between contigs within scaffolds, gaps between contigs were
155 removed and contigs were merged. Subsequently, ONT reads were mapped to the assembly
156 using Minimap2 (v2.24-r1122) and Samtools (v1.10) and mapped reads were visualized in
157 Integrative Genome Viewer (v2.11.2) (Danecek et al., 2021; Li, 2021; Robinson et al., 2011).
158 Whenever gaps between contigs were spanned by at least five reads with a mapping quality
159 of 30, the contigs were fused in the assembly.

160 Potential telomeres were identified using tapestry (v1.0.0) with “AACCCCT” as telomere
161 sequence (Davey et al., 2020). To check for potential contaminations, Blobtools (v1.1.1) and
162 BLAST (v2.13.0+) were used to create a Blobplot including taxonomic annotation at genus
163 level (Camacho et al., 2009; Laetsch and Blaxter, 2017). To check completeness of the
164 assembly and retrieve ploidy information, kat comp from the Kmer Analysis Toolkit (v2.4.2)
165 was used, and results were visualized using the kat plot spectra-cn function with the -x 800
166 option to extend the x-axis (Mapleson et al., 2016). Genome synteny to the closest sequenced
167 relative *C. subellipsoidea* C-169 was determined using Mummer3 (Blanc et al., 2012; Kurtz et
168 al., 2004). In detail, the two assemblies were first aligned using Nucmer, followed by a filtering
169 step with Delta-filter using the many-to-many option (-m). Finally, the alignment was visualized
170 with Mummerplot.

171

172 **Annotation**

173 To annotate repetitive elements **in the nuclear genome**, a database of simple repeats was
174 created with RepeatModeler (v2.0.3) that was expanded with transposable elements (TE) from
175 the TransposonUltimate resonaTE (v1.0) pipeline (Flynn et al., 2020; Riehl et al., 2022). This
176 pipeline uses multiple tools for TE prediction and combines the prediction output. For the
177 prediction of TEs in *Coccomyxa viridis* helitronScanner, ltrHarvest, mitefind, mitetracker,
178 RepeatModeler, RepeatMasker, sinefind, tirvish, transposonPSI and NCBI CDD1000 were
179 used within TransposonUltimate resonaTE and TEs that were predicted by at least two tools
180 were added to the database. TEclass (v2.1.3) was used for classification (Abrusán et al.,
181 2009). To softmask the genome and obtain statistics on the total TE and repetitive element

182 content in the genome, RepeatMasker (v4.1.2-p1)(Smit et al., 2012) was used with excln
183 option to exclude Ns in the masking.

184 Gene annotation **in the nuclear genome** was performed making use of RNA
185 sequencing data. To this end, the genome was indexed, and reads were mapped with HiSat2
186 (v2.2.1) using default settings (Kim et al., 2019). Afterwards, **BRAKER1** (v2.1.6) was used for
187 transcriptome-guided gene prediction **based on the RNA sequencing data** with default settings
188 (Hoff et al., 2016). To generate protein and coding sequence files the Braker output was
189 transformed with Gffread (v0.12.7) (Pertea and Pertea, 2020). PFAM domain annotation was
190 performed with InterProScan (v5.61) (Paysan-Lafosse et al., 2023). To estimate the number
191 of secreted proteins, SignalP (v6.0) was run in the slow-sequential mode on the annotated
192 proteins (Teufel et al., 2022). Finally, BUSCO (v5.3.2) was run with the Chlorophyta database
193 (chlorophyta_odb10) to estimate the completeness of the gene annotation (Manni et al.,
194 2021). The circos plot visualization of the annotation was created with R (v4.2.0) and Circilize
195 (v0.4.14) (Gu et al., 2014). All software and tools used for the genome assembly and
196 annotation are summarized in **Table S1**.

197 **Organelle genomes were annotated separately. Scaffolds were identified as organelle**
198 **genomes based on their lower GC content and smaller size. The mitochondrial genome was**
199 **annotated using MFannot (Lang et al., 2023) as well as GeSeq (Tillich et al., 2017) and the**
200 **annotation was combined within the GeSeq platform. The plastid genome was annotated**
201 **using GeSeq alone. The annotations were visualized using the OGDRAW webserver (Greiner**
202 **et al., 2019).**

203

204

205 Results

206 The version 1 genome of *C. viridis* was assembled from 32.2 Gbp of PacBio HiFi reads with
207 a mean read length of 15 kb, 0.95 Gbp Nanopore reads with a mean read length of 8.8 kb and
208 15 million pairs of Hi-C seq data. The PacBio HiFi reads were first assembled using Raven
209 (Vaser and Šikić, 2021), yielding 27 contigs. These contigs were scaffolded and manually
210 curated using Hi-C data (Dudchenko et al., 2017; Durand et al., 2016a; Durand et al., 2016b;
211 Li and Durbin, 2009). To close the remaining gaps between contigs within scaffolds, ONT
212 reads were mapped onto the assembly (Danecek et al., 2021; Li, 2021) and gaps that were
213 spanned by at least 5 ONT reads with a mapping quality >30 were manually closed, finally
214 resulting in 21 scaffolds consisting of 26 contigs with a total length of 50.9 Mb and an N50 of
215 2.7 Mb (Figure 1, **Table 1**). Using Tapestry (Davey et al., 2020), telomeric regions
216 ([AACCCT]*n*) were identified at both ends of nine of the 21 scaffolds (≥5 repeats) (Figure 1a),
217 suggesting that these represent full-length chromosomes, which was confirmed by Hi-C

218 analysis (Figure 1b). Additionally, the Hi-C contact map indicated centromeres for some of the
219 chromosomes. However, the determination of exact centromere locations on all chromosomes
220 will require ChIP-seq analysis and CenH3 mapping. While Tapestry detected telomeric
221 sequences at only one end of eight other scaffolds and none for scaffold 18 and 19, the Hi-C
222 map points towards the presence of telomeric repeats at both ends of all scaffolds 1-19 (Figure
223 1b), suggesting that the v1 assembly contains 19 full-length chromosomes that compose the
224 nuclear genome. Scaffolds 20 and 21 were considerably shorter with ~162 kb and ~70 kb and
225 displayed a markedly lower GC content at 41-42% (Figure 1a), suggesting that these scaffolds
226 represent the chloroplast and mitochondrial genomes, respectively. BLAST analyses
227 confirmed the presence of plastid and mitochondrial genes on the respective scaffolds, and
228 the overall scaffold lengths corresponded with the sizes of the plastid and mitochondrial
229 genomes of *Coccomyxa subellipsoidea* C-169 with 175 kb and 65 kb, respectively (Blanc et
230 al., 2012). Full annotation of scaffolds 20 and 21 showed that they indeed represent
231 chloroplast and mitochondrial genomes, respectively (Figure 2).

232 To rule out the presence of contaminants, the assembly and PacBio HiFi raw reads
233 were used to produce a Blobplot (Camacho et al., 2009; Laetsch and Blaxter, 2017), which
234 indicates that 98.76% of the reads match only the *Coccomyxa* genus (Figure 3) and,
235 consequently, that the original sample was free of contaminating organisms. Finally, a KAT
236 analysis showed a single peak of k-mer multiplicity based on HiFi reads that were represented
237 once in the assembly (Figure 4) (Mapleson et al., 2016), indicative of a high-quality, haploid
238 genome.

239 To annotate the nuclear genome, we first assessed the presence of repetitive
240 elements. In total, we found 8.9% of the genome to be repetitive (Table 2), comparable to the
241 7.2% of repetitive sequences found in the genome of *C. subellipsoidea* C-169 (Blanc et al.,
242 2012). These 8.9% repetitive elements were annotated as either simple repeats (2.3%) or
243 transposable elements (6.6%). Of the transposable elements, 36% were annotated as
244 retrotransposons and 64% as DNA transposons. The distribution of the repetitive elements
245 was even across the genome with only a few repeat-rich regions (Figure 5). Next, we aimed
246 to produce a high-quality genome annotation using RNA sequencing data. In total 13,557
247 genes were annotated with an average length of 3.1 kb (Table 2). The amount of alternative
248 splicing in the genome is predicted to be very low, given the average of one transcript per
249 gene model. To confirm the actual amount of alternative splicing, however, further analyses
250 will be required. Of the 13,557 genes, 68% have annotated PFAM domains and 962 are
251 predicted to carry a signal peptide for secretion. A total of 1,489 (98.6 %) complete gene
252 models among 1,519 conserved Benchmarking Universal Single-Copy Orthologs (BUSCO)
253 (Manni et al., 2021) in the chlorophyta_odb10 database were identified (Table 2), suggesting
254 a highly complete genome annotation.

255 Until recently, the taxonomic classification and definition of *Coccomyxa* species was
256 based on environmentally variable morphological and cytological characteristics. This
257 classification was reviewed based on the phylogenetic analyses of nuclear SSU and ITS rDNA
258 sequences, which resulted in the definition of 27 currently recognized *Coccomyxa* species
259 (Darienko et al., 2015; Malavasi et al., 2016). Dot plot analysis of the high-quality genome
260 assembly of *C. viridis* SAG216-4 with the assembly of the most closely related sequenced
261 relative *C. subellipsoidea* C-169 revealed a lack of synteny since the few identified orthologous
262 sequences were < 1 kb and, therefore, do not represent full-length genes (Figure 6a, Table
263 2). This lack of synteny was no technical artifact since the *C. viridis* assembly could be fully
264 aligned to itself (Figure 6b), and BLAST analyses with five out of six non-identical ITS
265 sequences identified in the *C. viridis* SAG 216-4 assembly confirmed its species identity. A
266 comparison of the assembly of *C. subellipsoidea* C-169 to that of *Chlorella variabilis*
267 (Chlorophyte, *Trebouxiophyceae*) has previously identified few syntenic regions which
268 displayed poor gene collinearity (Blanc et al., 2012). Future studies will help to clarify whether
269 the absence of synteny between *C. viridis* and *C. subellipsoidea* is due to the quality of the
270 available assemblies or whether it has biological implications.

271

272 Data availability

273 Data for *C. viridis* SAG 216-4 with the ToLID ucCocViri1 is available via the European
274 Nucleotide Archive (ENA) under the study accession number PRJNA1054215. Fastqc reports
275 of raw data can be found in (Kraege et al., 2023).

276

277 Acknowledgements

278 This genome project is part of the DeRGA pilot study ([https://www.erga-biodiversity.eu/pilot-](https://www.erga-biodiversity.eu/pilot-project)
279 [project](https://www.erga-biodiversity.eu/pilot-project)) and we acknowledge coordination and support by Dr. Astrid Böhne and Prof. Dr. Ann-
280 Marie Waldvogel. Development of the pilot ERGA Data Portal ([https://portal.erga-](https://portal.erga-biodiversity.eu/)
281 [biodiversity.eu/](https://portal.erga-biodiversity.eu/)) was funded by the European Molecular Biology Laboratory. We thank the
282 West German Genome Center (WGGC) and the Cologne Center for Genomics (CCG) for
283 library generation and quality control as well as their support for PacBio Hifi, Hi-C and RNA
284 sequencing.

285

286 Conflict of interest

287 The authors declare no conflict of interest.

288

289 Funding information

290 BPHJT acknowledges funding by the Alexander von Humboldt Foundation in the framework
291 of an Alexander von Humboldt Professorship endowed by the German Federal Ministry of
292 Education and Research and is furthermore supported by the Deutsche
293 Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence
294 Strategy – EXC 2048/1 – Project ID: 390686111. This research was also funded by the
295 Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – SFB1535 -
296 Project ID 458090666. PHS was funded by a DFG ENP grant (grant number: 434028868),
297 which also funded JK's position. NG was first funded through a DFG grant to PHS
298 (458953049) and subsequently through the European Union's Horizon Europe Research and
299 Innovation program under the Marie Skłodowska-Curie grant agreement No. 101110569.

300

301 References

302 Abrusán, G., Grundmann, N., DeMester, L., and Makalowski, W. (2009). TEclass--a tool for
303 automated classification of unknown eukaryotic transposable elements. *Bioinformatics* 25, 1329-1330.

304 Blanc, G., Agarkova, I., Grimwood, J., Kuo, A., Brueggeman, A., Dunigan, D.D., Gurnon, J., Ladunga,
305 I., Lindquist, E., Lucas, S., *et al.* (2012). The genome of the polar eukaryotic microalga *Coccomyxa*
306 *subellipsoidea* reveals traits of cold adaptation. *Genome Biol* 13, R39.

307 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T.L.
308 (2009). BLAST+: architecture and applications. *BMC Bioinformatics* 10, 421.

309 Cao, S., Zhang, F., Zheng, H., Liu, C., Peng, F., and Zhou, Q. (2018a). *Coccomyxa antarctica* sp.
310 nov. from the Antarctic lichen *Usnea aurantiacoatra*. *PhytoKeys*, 107-115.

311 Cao, S., Zhang, F., Zheng, H., Peng, F., Liu, C., and Zhou, Q. (2018b). *Coccomyxa greatwallensis* sp.
312 nov. (Trebouxiophyceae, Chlorophyta), a lichen epiphytic alga from Fildes Peninsula, Antarctica.
313 *PhytoKeys*, 39-50.

314 Danecek, P., Bonfield, J.K., Liddle, J., Marshall, J., Ohan, V., Pollard, M.O., Whitwham, A., Keane, T.,
315 McCarthy, S.A., Davies, R.M., *et al.* (2021). Twelve years of SAMtools and BCFtools. *Gigascience* 10.

316 Darienko, T., Gustavs, L., Eggert, A., Wolf, W., and Proschold, T. (2015). Evaluating the species
317 boundaries of green microalgae (*Coccomyxa*, *Trebouxiophyceae*, *Chlorophyta*) using integrative
318 taxonomy and DNA barcoding with further implications for the species identification in environmental
319 samples. *PLoS One* 10, e0127838.

320 Davey, J.W., Davis, S.J., Mottram, C., and Ashton, P.D. (2020). Tapestry: validate and edit small
321 eukaryotic genome assemblies with long reads. *bioRxiv*.

322 Dudchenko, O., Batra, S.S., Omer, A.D., Nyquist, S.K., Hoeger, M., Durand, N.C., Shamim, M.S.,
323 Machol, I., Lander, E.S., Aiden, A.P., *et al.* (2017). *De novo* assembly of the *Aedes aegypti* genome
324 using Hi-C yields chromosome-length scaffolds. *Science* 356, 92-95.

325 Durand, N.C., Robinson, J.T., Shamim, M.S., Machol, I., Mesirov, J.P., Lander, E.S., and Aiden, E.L.
326 (2016a). Juicebox provides a visualization system for Hi-C contact maps with unlimited zoom. *Cell*
327 *Syst* 3, 99-101.

328 Durand, N.C., Shamim, M.S., Machol, I., Rao, S.S.P., Huntley, M.H., Lander, E.S., and Aiden, E.L.
329 (2016b). Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell*
330 *Systems* 3, 95-98.

- 331 Faluaburu, M.S., Nakai, R., Imura, S., and Naganuma, T. (2019). Phylotypic characterization of
332 mycobionts and photobionts of rock tripe lichen in East Antarctica. *Microorganisms* 7.
- 333 Flynn, J.M., Hubley, R., Goubert, C., Rosen, J., Clark, A.G., Feschotte, C., and Smit, A.F. (2020).
334 RepeatModeler2 for automated genomic discovery of transposable element families. *Proc Natl Acad*
335 *Sci U S A* 117, 9451-9457.
- 336 Gray, A.P., Lucas, I.A.N., Seed, R., and Richardson, C.A. (1999). *Mytilus edulis chilensis* infested with
337 *Coccomyxa parasitica* (Chlorococcales, Coccomyxaceae). *Journal of Molluscan Studies* 65, 289-294.
- 338 Greiner, S., Lehwark, P., and Bock, R. (2019). OrganellarGenomeDRAW (OGDRAW) version 1.3.1:
339 expanded toolkit for the graphical visualization of organellar genomes. *Nucleic Acids Res* 47, W59-
340 W64.
- 341 Gu, Z., Gu, L., Eils, R., Schlesner, M., and Brors, B. (2014). Circlize implements and enhances
342 circular visualization in R. *Bioinformatics* 30, 2811-2812.
- 343 Gustavs, L., Schiefelbein, U., Darienko, T., and Pröschold, T. (2017). Symbioses of the green algal
344 genera *Coccomyxa* and *Elliptochloris* (Trebouxiophyceae, Chlorophyta). In *Algal and Cyanobacteria*
345 *Symbioses*, M. Grube, J. Seckbach, and L. Muggia, eds. (Europe: World Scientific), pp. 169-208.
- 346 Hoff, K.J., Lange, S., Lomsadze, A., Borodovsky, M., and Stanke, M. (2016). BRAKER1:
347 Unsupervised RNA-Seq-based genome annotation with GeneMark-ET and AUGUSTUS.
348 *Bioinformatics* 32, 767-769.
- 349 Jaag, O. (1933). *Coccomyxa* Schmidle- Monographie einer Algengattung (Bern: Gebrüder Fretz).
- 350 Kim, D., Paggi, J.M., Park, C., Bennett, C., and Salzberg, S.L. (2019). Graph-based genome
351 alignment and genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol* 37, 907-915.
- 352 Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H., and Phillippy, A.M. (2017). Canu:
353 scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation.
354 *Genome Res* 27, 722-736.
- 355 Kraege, A., Thomma, B.P.H.J., and Rovenich, H. (2023). Fastqc reports of sequencing data from
356 *Coccomyxa viridis* SAG 216-4.
- 357 Kulichová, J., Škaloud, P., and Neustupa, J. (2014). Molecular diversity of green corticolous
358 microalgae from two sub-Mediterranean European localities. *European Journal of Phycology* 49, 345-
359 355.
- 360 Kurtz, S., Phillippy, A., Delcher, A.L., Smoot, M., Shumway, M., Antonescu, C., and Salzberg, S.L.
361 (2004). Versatile and open software for comparing large genomes. *Genome Biol* 5, R12.
- 362 Laetsch, D.R., and Blaxter, M.L. (2017). BlobTools: Interrogation of genome assemblies. *F1000*
363 *Research* 6, 1287.
- 364 Lang, B.F., Beck, N., Prince, S., Sarrasin, M., Rioux, P., and Burger, G. (2023). Mitochondrial genome
365 annotation with MFannot: a critical analysis of gene identification and gene model prediction. *Front*
366 *Plant Sci* 14, 1222186.
- 367 Leliaert, F., Smith, D.R., Moreau, H., Herron, M.D., Verbruggen, H., Delwiche, C.F., and De Clerck, O.
368 (2012). Phylogeny and molecular evolution of the green algae. *Crit Rev Plant Sci* 31, 1-46.
- 369 Li, H. (2021). New strategies to improve minimap2 alignment accuracy. *Bioinformatics* 37, 4572-4574.
- 370 Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
371 transform. *Bioinformatics* 25, 1754-1760.

- 372 Li, L., Wang, S., Wang, H., Sahu, S.K., Marin, B., Li, H., Xu, Y., Liang, H., Li, Z., Cheng, S., *et al.*
 373 (2020). The genome of *Prasinoderma coloniale* unveils the existence of a third phylum within green
 374 plants. *Nat Ecol Evol* 4, 1220-1231.
- 375 Malavasi, V., Skaloud, P., Rindi, F., Tempesta, S., Paoletti, M., and Pasqualetti, M. (2016). DNA-
 376 Based taxonomy in ecologically versatile microalgae: A re-evaluation of the species concept within the
 377 coccoid green algal genus *Coccomyxa* (Trebouxiophyceae, Chlorophyta). *PLoS One* 11, e0151137.
- 378 Manni, M., Berkeley, M.R., Seppey, M., Simao, F.A., and Zdobnov, E.M. (2021). BUSCO update:
 379 novel and streamlined workflows along with broader and deeper phylogenetic coverage for scoring of
 380 eukaryotic, prokaryotic, and viral genomes. *Mol Biol Evol* 38, 4647-4654.
- 381 Mapleson, D., Accinelli, G.G., Kettleborough, G., Wright, J., and Clavijo, B.J. (2016). KAT: A K-mer
 382 Analysis Toolkit to quality control NGS datasets and genome assemblies. *Bioinformatics* 33, 574-576.
- 383 Marin, B. (2012). Nested in the Chlorellales or independent class? Phylogeny and classification of the
 384 Pedinophyceae (Viridiplantae) revealed by molecular phylogenetic analyses of complete nuclear and
 385 plastid-encoded rRNA operons. *Protist* 163, 778-805.
- 386 McFadden, G.I., and Melkonian, M. (1986). Use of Hepes buffer for microalgal culture media and
 387 fixation for electron microscopy. *Phycologia* 25, 551-557.
- 388 Morris, J.L., Puttick, M.N., Clark, J.W., Edwards, D., Kenrick, P., Pressel, S., Wellman, C.H., Yang, Z.,
 389 Schneider, H., and Donoghue, P.C.J. (2018). The timescale of early land plant evolution. *Proc Natl*
 390 *Acad Sci U S A* 115, E2274-E2283.
- 391 Paysan-Lafosse, T., Blum, M., Chuguransky, S., Grego, T., Pinto, B.L., Salazar, G.A., Bileschi, M.L.,
 392 Bork, P., Bridge, A., Colwell, L., *et al.* (2023). InterPro in 2022. *Nucleic Acids Res* 51, D418-D427.
- 393 Perteua, G., and Perteua, M. (2020). GFF utilities: GffRead and GffCompare. *F1000 Research* 9, 304.
- 394 Riehl, K., Riccio, C., Miska, E.A., and Hemberg, M. (2022). TransposonUltimate: software for
 395 transposon classification, annotation and detection. *Nucleic Acids Res* 50, e64.
- 396 Robinson, J.T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E.S., Getz, G., and Mesirov,
 397 J.P. (2011). Integrative genomics viewer. *Nat Biotechnol* 29, 24-26.
- 398 Schmidle, W. (1901). Über drei Algengenera. *Berichte der deutschen botanischen Gesellschaft* 19,
 399 10-24.
- 400 Sciuto, K., Baldan, B., Maracto, S., and Moro, I. (2019). *Coccomyxa cimbria* sp. nov., a green
 401 microalga found in association with carnivorous plants of the genus *Drosera* L. *European Journal of*
 402 *Phycology* 54, 531-547.
- 403 Smit, A.F., Hubley, R., and Green, P. (2012). RepeatMasker (Retrieved from
 404 <https://repeatmasker.org>).
- 405 Sokolnikova, Y., Magarlamov, T., Stenkova, A., and Kumeiko, V. (2016). Permanent culture and
 406 parasitic impact of the microalga *Coccomyxa parasitica*, isolated from horse mussel *Modiolus*
 407 *kurilensis*. *J Invertebr Pathol* 140, 25-34.
- 408 Sokolnikova, Y., Tumas, A., Stenkova, A., Slatvinskaya, V., Magarlamov, T., and Smagina, E. (2022).
 409 Novel species of parasitic green microalgae *Coccomyxa veronica* sp. nov. infects *Anadara*
 410 *broughtonii* from Sea of Japan. *Symbiosis* 87, 293-305.
- 411 Štifterová, A., and Neustupa, J. (2015). Community structure of corticolous microalgae within a single
 412 forest stand: evaluating the effects of bark surface pH and tree species. *Fottea Olomouc* 15, 113-122.

- 413 Tagirdzhanova, G., Scharnagl, K., Yan, X., and Talbot, N.J. (2023). Genomic analysis of *Coccomyxa*
414 *viridis*, a common low-abundance alga associated with lichen symbioses. *Sci Rep* 13, 21285.
- 415 Teufel, F., Almagro Armenteros, J.J., Johansen, A.R., Gislason, M.H., Pihl, S.I., Tsirigos, K.D.,
416 Winther, O., Brunak, S., von Heijne, G., and Nielsen, H. (2022). SignalP 6.0 predicts all five types of
417 signal peptides using protein language models. *Nat Biotechnol* 40, 1023-1025.
- 418 Tillich, M., Lehwark, P., Pellizzer, T., Ulbricht-Jones, E.S., Fischer, A., Bock, R., and Greiner, S.
419 (2017). GeSeq - versatile and accurate annotation of organelle genomes. *Nucleic Acids Res* 45, W6-
420 W11.
- 421 Trémouillaux-Guiller, J., Rohr, T., Rohr, R., and Huss, V.A.R. (2002). Discovery of an endophytic alga
422 in *Ginkgo biloba*. *Am J Bot* 89, 727-733.
- 423 Vaschenko, M.A., Kovaleva, A.L., Syasina, I.G., and Kukhlevsky, A.D. (2013). Reproduction-related
424 effects of green alga *Coccomyxa* sp. infestation in the horse mussel *Modiolus modiolus*. *J Invertebr*
425 *Pathol* 113, 86-95.
- 426 Vaser, R., and Šikić, M. (2021). Time- and memory-efficient genome assembly with Raven. *Nature*
427 *Computational Science* 1, 332-336.
- 428 Wick, R. (2018). Porechop (Retrieved from <https://github.com/rrwick/Porechop>).
- 429 Yahr, R., Florence, A., Škaloud, P., and Voytsekhovich, A. (2015). Molecular and morphological
430 diversity in photobionts associated with *Micarea* s. str. (*Lecanorales*, Ascomycota). *The Lichenologist*
431 47, 403-414.
- 432 Zoller, S., and Lutzoni, F. (2003). Slow algae, fast fungi: exceptionally high nucleotide substitution rate
433 differences between lichenized fungi Omphalina and their symbiotic green algae *Coccomyxa*. *Mol*
434 *Phylogenet Evol* 29, 629-640.
435
436

437 Figure legends

438 **Figure 1. Genome assembly of *Coccomyxa viridis* SAG 216-4.** (a) An overview of the *C.*
439 *viridis* genome assembly depicts chromosome-scale scaffolds. Green bars indicate scaffold
440 sizes and red bars represent telomeres. Variations in color intensities correlate with read
441 coverage. Read coverage per scaffold is determined by mapping PacBio HiFi reads onto the
442 assembly. Scaffolds 20 and 21 were identified as chloroplast and mitochondrial genomes
443 based on size and low GC contents, and BLAST analyses. (b) Hi-C contact map showing
444 interaction frequencies between regions in the nuclear genome of *Coccomyxa viridis*.
445 Scaffolds are framed by blue lines while contigs within scaffolds are depicted in green.

446

447 **Figure 2 Scaffolds 20 and 21 represent the plastid and mitochondrial genomes of *C.***
448 ***viridis* SAG 216-4.** Gene maps of the chloroplast (a) and mitochondrial (b) genomes. The
449 inner circles indicate the GC content and mapped genes are shown on the outer circles. Genes
450 that are transcribed clockwise are placed inside the outer circles, and genes that are
451 transcribed counterclockwise at the outside of the outer circles.

452

453 **Figure 3. Taxonomic annotation indicates absence of contaminations in the genome**
454 **assembly.** (b) Taxon-annotated GC coverage scatter plot (Blobplot) of the contigs from the
455 genome assembly shows that all scaffolds are taxon-annotated as *Coccomyxa* and all
456 scaffolds that belong to the nuclear genome have similar GC contents (~54%). The GC
457 content of the mitochondrial and plastid genomes are considerably lower (~41%). (b) In total
458 98.76% of the reads can be mapped onto the assembly and are therefore classified as
459 *Coccomyxa* reads.

460

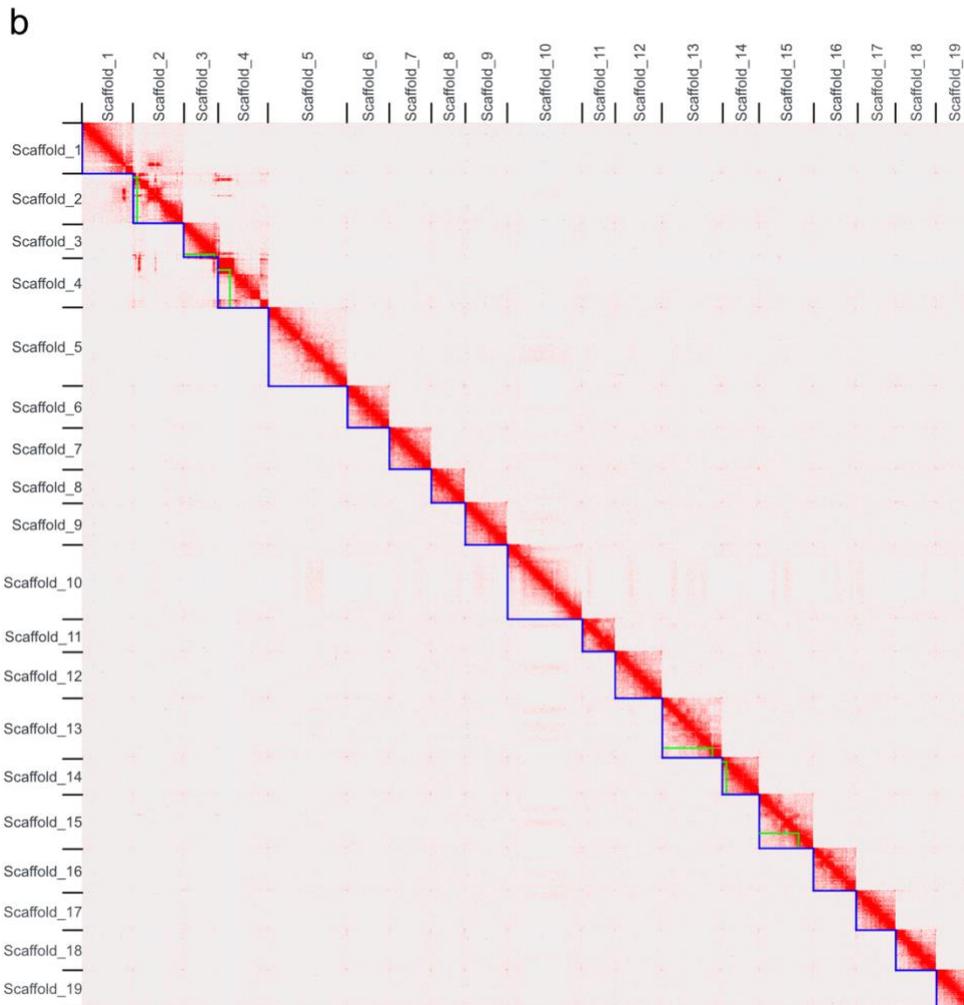
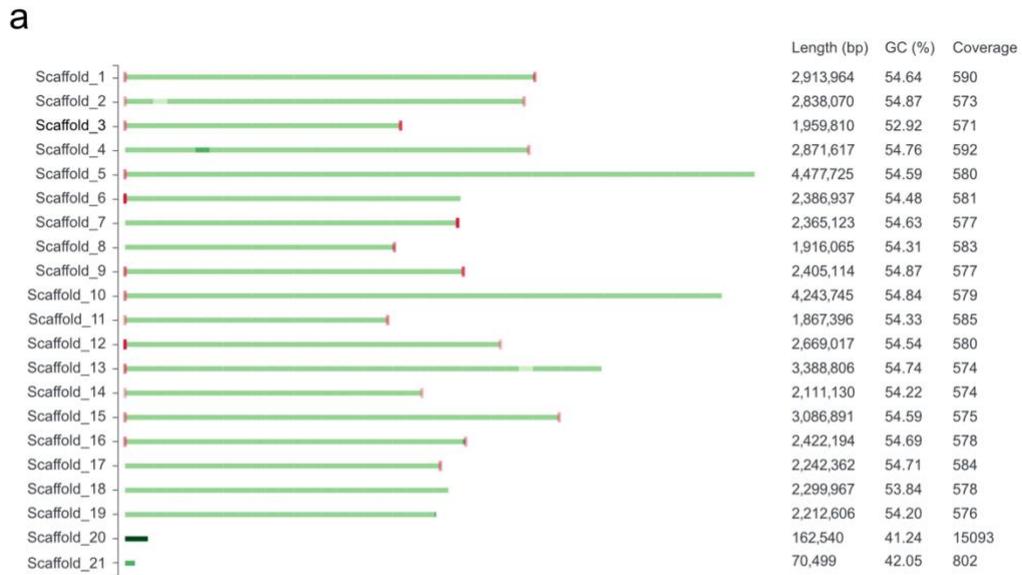
461 **Figure 4. The *Coccomyxa viridis* SAG 216-4 genome is haploid.** The KAT specra-cn plot
462 depicts the 27-mer multiplicity of the PacBio HiFi reads against the genome assembly. Black
463 areas under the peaks represent k-mers present in the reads but absent from the assembly,
464 colored peaks indicate k-mers that are present once to multiple times in the assembly. The
465 single red peak in the KAT specra-cn plot suggests that *Coccomyxa viridis* has a haploid
466 genome, while the black peak at low multiplicity shows that the assembly is highly complete
467 and that all reads are represented in the assembly.

468

469 **Figure 5. Circos plot summarizing the nuclear genome annotation of *Coccomyxa viridis***
470 **SAG 216-4.** From outside to inside the tracks display: GC content (over 1-kb windows), gene
471 density (blue) and repetitive element density (red).

472

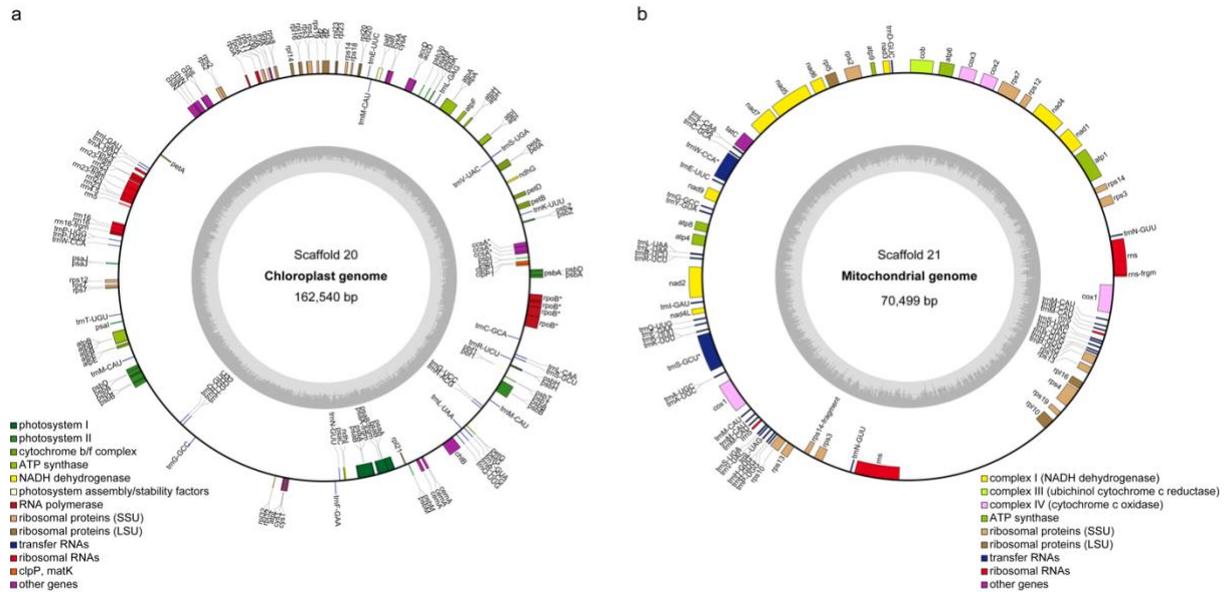
473 **Figure 6. No synteny detected between related *Coccomyxa* species.** (a) Dot plot of
474 orthologous sequences in the genome assemblies of *C. viridis* SAG 216-4 and *C.*
475 *subellipsoidea* C-169. Violet and blue dots represent orthologous sequences on same and
476 opposite strands, respectively. Dot sizes does not correlate with the length of the sequences
477 they represent, which were all < 1 kb. The width of each box corresponds to the length (bp) of
478 the respective scaffold. (b) Dot plot of the genome assembly of *C. viridis* SAG216-4 against
479 itself.
480
481



482

483 **Figure 1. Genome assembly of *Coccomyxa viridis* SAG 216-4.** (a) An overview of the *C.*
 484 *viridis* genome assembly depicts chromosome-scale scaffolds. Green bars indicate scaffold
 485 sizes and red bars represent telomeres. Variations in color intensities correlate with read

486 coverage. Read coverage per scaffold is determined by mapping PacBio HiFi reads onto the
487 assembly. Scaffolds 20 and 21 were identified as chloroplast and mitochondrial genomes
488 based on size and low GC contents, and BLAST analyses. (b) Hi-C contact map showing
489 interaction frequencies between regions in the nuclear genome of *Coccomyxa viridis*.
490 Scaffolds are framed by blue lines while contigs within scaffolds are depicted in green.
491
492



493

494

Figure 2 Scaffolds 20 and 21 represent the plastid and mitochondrial genomes of *C.*

495

***viridis* SAG 216-4. Gene maps of the chloroplast (a) and mitochondrial (b) genomes. The**

496

inner circles indicate the GC content and mapped genes are shown on the outer circles. Genes

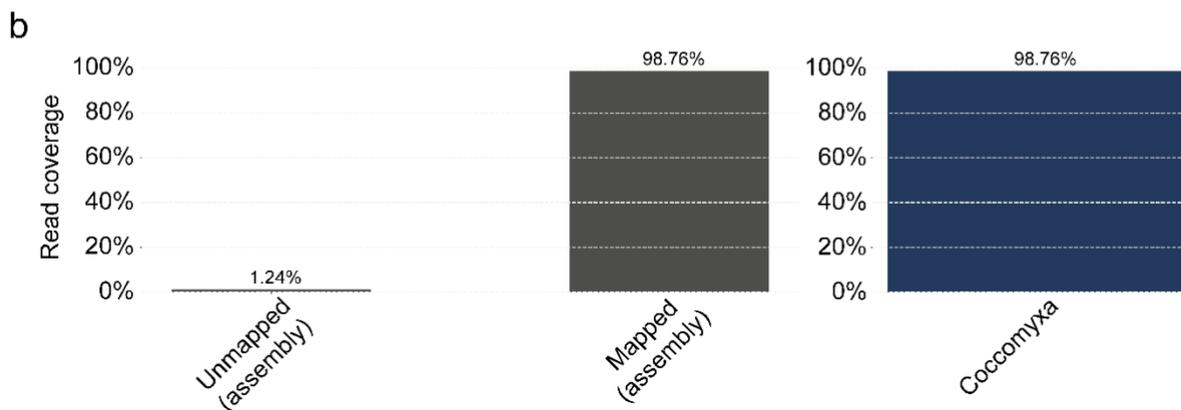
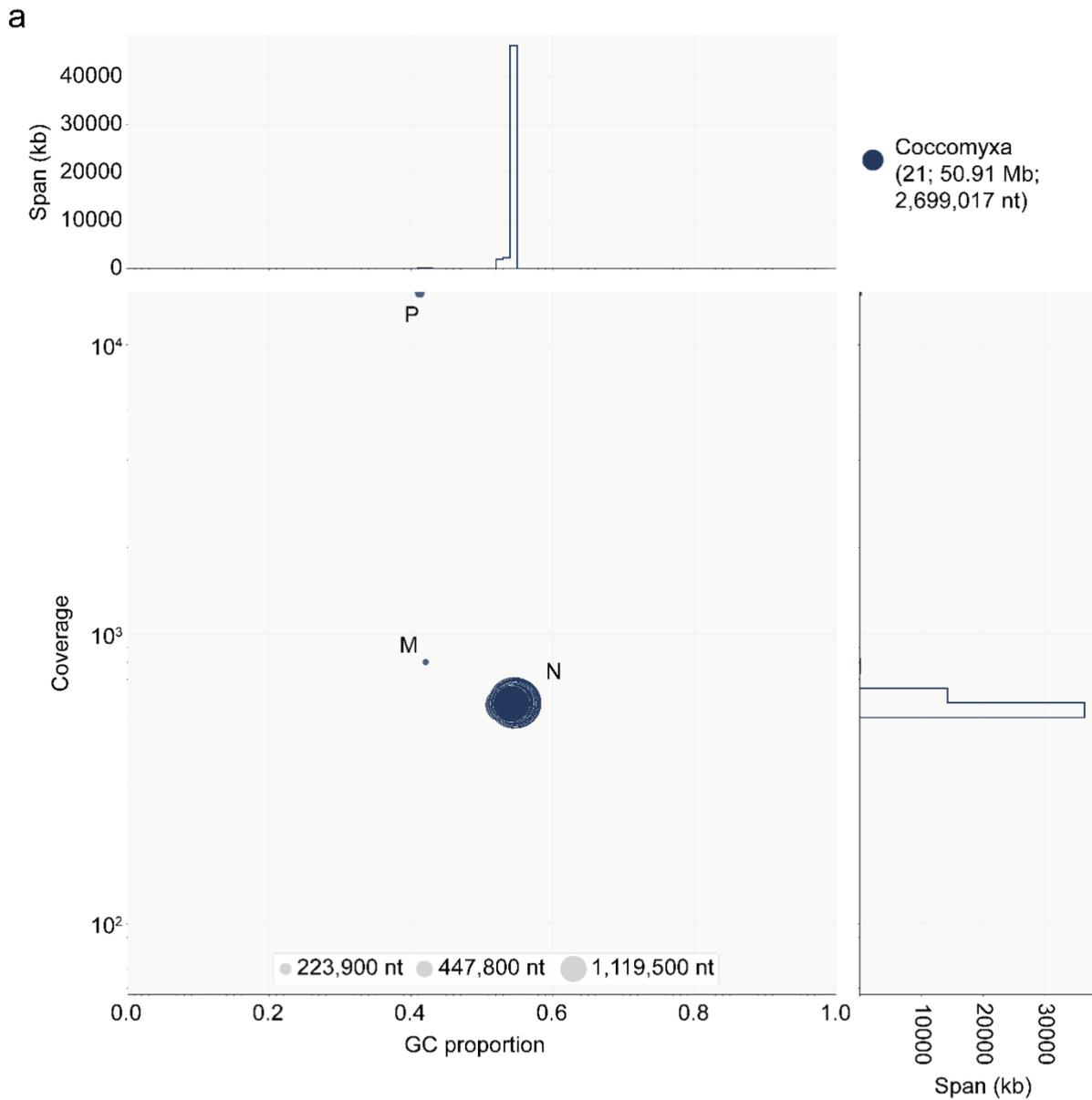
497

that are transcribed clockwise are placed inside the outer circles, and genes that are

498

transcribed counterclockwise at the outside of the outer circles.

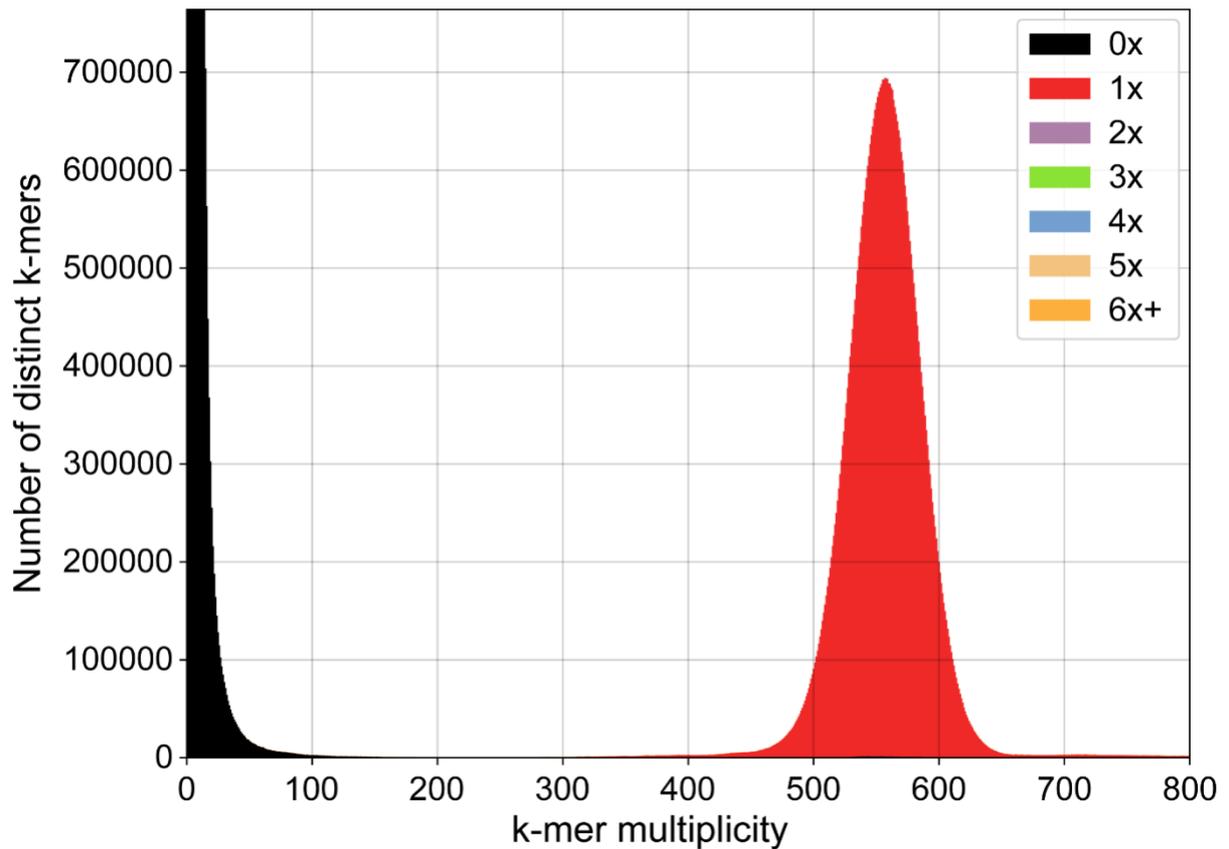
499



500

501 **Figure 3. Taxonomic annotation indicates absence of contaminations in the genome**
 502 **assembly.** (b) Taxon-annotated GC coverage scatter plot (Blobplot) of the contigs from the
 503 genome assembly shows that all scaffolds are taxon-annotated as *Coccomyxa* and all
 504 scaffolds that belong to the nuclear genome (N) have similar GC contents (~54%). The GC

505 content of the mitochondrial (M) and plastid (P) genomes are considerably lower (~41%). (b)
506 In total 98.76% of the reads can be mapped onto the assembly and are therefore classified
507 as *Coccomyxa* reads.
508
509

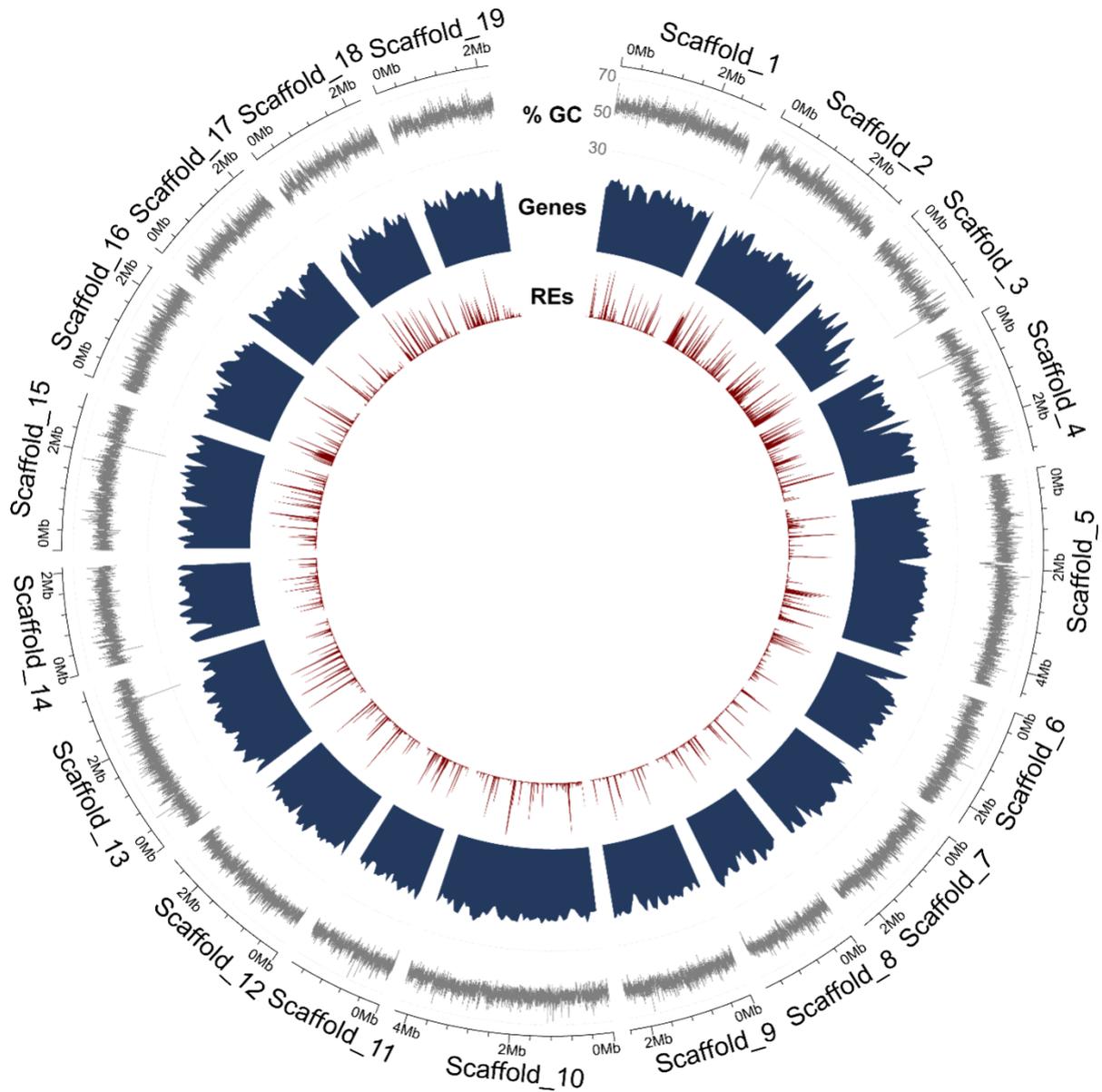


510

511 **Figure 4.** The *Coccoomyxa viridis* SAG 216-4 **nuclear** genome is haploid. The KAT spectra-
 512 cn plot depicts the 27-mer multiplicity of the PacBio HiFi reads against the **nuclear** genome
 513 assembly. Black areas under the peaks represent k-mers present in the reads but absent from
 514 the assembly, colored peaks indicate k-mers that are present once to multiple times in the
 515 assembly. The single red peak in the KAT spectra-cn plot suggests that *Coccoomyxa viridis* has
 516 a haploid genome, while the black peak at low multiplicity shows that the assembly is highly
 517 complete and that all reads are represented in the assembly.

518

519



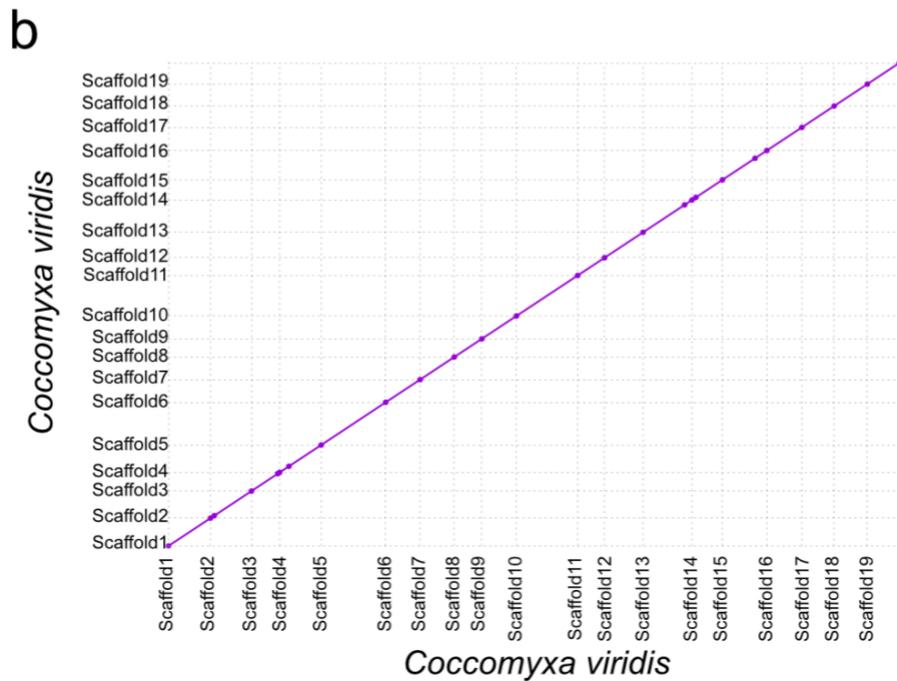
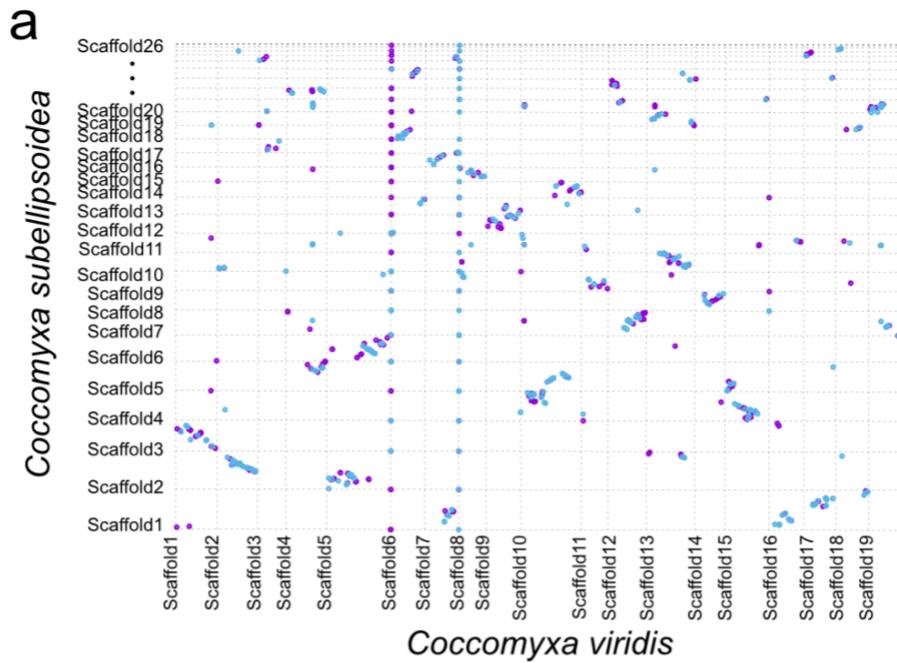
520

521 **Figure 5.** Circos plot summarizing the **nuclear** genome annotation of *Coccomyxa viridis*
 522 **SAG 216-4.** From outside to inside the tracks display: GC content (over 1-kb windows), gene
 523 density (blue) and repetitive element density (red).

524

525

526



527
528
529
530
531
532
533
534
535
536

Figure 6. No synteny detected between related *Coccomyxa* species. (a) Dot plot of orthologous sequences in the genome assemblies of *C. viridis* SAG 216-4 and *C. subellipsoidea* C-169. Violet and blue dots represent orthologous sequences on same and opposite strands, respectively. Dot sizes does not correlate with the length of the sequences they represent, which were all < 1 kb. The width of each box corresponds to the length (bp) of the respective scaffold. (b) Dot plot of the genome assembly of *C. viridis* SAG216-4 against itself.

537 **Table 1.** Genome features of *C. viridis* SAG 216-4 including the mitochondrial and plastid genomes.

Assembly ID	<i>C. viridis</i> SAG 216-4 genomes
Total length (bp)	50,911,578
No. of contigs	27
No. of scaffolds	21
Longest scaffold (bp)	4,477,725
N50 (bp)	2,669,017
L50	8
GC content (%)	54.5

538

539

540 **Table 2** Annotation features of the *C. viridis* SAG 216-4 nuclear genome.

Genome annotation	
Repeat content (%)	8.85
Retrotransposons	2.4
DNA transposons	4.2
Simple repeats	2.25
No. gene models	13,557
Average gene length (bp)	3146
No. exons	122,978
Average no. exons per gene model	9
Average exon length (bp)	158
No. transcripts	14,024
Average no. transcripts/gene model	1
No. gene models <200 bp length	0
No. proteins with ≥ 1 PFAM domain	9205
No. proteins with signal peptide	962
BUSCO (chlorophyta_odb10)	C: 98.6% [S: 82.5%, D: 16.1%], F: 0.1%, M: 1.3%, N: 1519

541

542

543 **Table S1.** Summary of bioinformatics tools used for genome assembly and annotation.

Assembly		Annotation	
Tool	Version	Tool	Version
Raven	v1.8.1	RepeatModeler	v2.0.3
Juicer	v2.0	TransposonUltimate	v1.0
BWA	v0.7.17-r1188	TEclass	v2.1.3
3d-dna	v180922	RepeatMasker	v4.1.2-p1
Juicebox	v1.11.08	HiSat2	v2.2.1
Minimap2	v2.24-r1122	Braker	v2.1.6
Samtools	v1.10	Gffread	v0.12.7
Integrative Genome Viewer	v2.11.2	SignalP	v6.0
Tapestry	v1.0.0	BUSCO	v5.3.2
Blobtools	v1.1.1	R	v4.2.0
BLAST	2.13.0+	Circilize	v0.4.14
Kmer Analysis Toolkit	V2.4.2	InterProScan	v5.61
Mummer	C3.23		

544

545

546