

**The plastid genome in Cladophorales green algae
is encoded by hairpin chromosomes**

Andrea Del Cortona^{1,2,3,4,a}, Frederik Leliaert^{1,5,a}, Kenny A. Bogaert¹, Monique Turmel⁶,
Christian Boedeker⁷, Jan Janouškovec⁸, Juan M. Lopez-Bautista⁹, Heroen Verbruggen¹⁰,
Klaas Vandepoele^{2,3,4} & Olivier De Clerck^{1*}

¹Department of Biology, Phycology Research Group, Ghent University, Krijgslaan 281, 9000 Ghent, Belgium

²Department of Plant Biotechnology and Bioinformatics, Ghent University, Technologiepark 927, 9052 Zwijnaarde, Belgium

³VIB Center for Plant Systems Biology, Technologiepark 927, 9052 Zwijnaarde, Belgium

⁴Bioinformatics Institute Ghent, Ghent University, Technologiepark 927, 9052 Zwijnaarde, Belgium

⁵Botanic Garden Meise, Nieuwelaan 38, 1860 Meise, Belgium

⁶Institut de biologie intégrative et des systèmes, Département de biochimie, de microbiologie et de bio-informatique, Pavillon Charles-Eugène-Marchand

1030, avenue de la Médecine, Université Laval, G1V 0A6 , Québec (QC), Canada

⁷School of Biological Sciences, New Kirk Building Kelburn Parade, Victoria University of Wellington, PO Box 600 Wellington, New Zealand

⁸Department of Genetics, Evolution and Environment, Gower Street, University College London, London WC1E 6BT, United Kingdom

⁹Department of Biological Sciences, 300 Hackberry Lane, The University of Alabama, Tuscaloosa, AL35484-0345, USA

¹⁰School of BioSciences, Professors Walk, University of Melbourne, Victoria 3010, Australia

^a These authors contributed equally to this work

* Correspondence and requests for materials should be addressed to the Lead Contact O.D.C. (olivier.declerck@ugent.be).

Summary

Virtually all plastid (chloroplast) genomes are circular double-stranded DNA molecules, typically between 100-200 kb in size and encoding circa 80-250 genes. Exceptions to this universal plastid genome architecture are very few and include the dinoflagellates where genes are located on DNA minicircles. Here we report on the highly deviant chloroplast genome of Cladophorales green algae, which is entirely fragmented into hairpin chromosomes. Short and long read high-throughput sequencing of DNA and RNA demonstrated that the chloroplast genes of *Boodlea composita* are encoded on 1-7 kb DNA contigs with an exceptionally high GC-content, each containing a long inverted repeat with one or two protein-coding genes and conserved non-coding regions putatively involved in replication and/or expression. We propose that these contigs correspond to linear single-stranded DNA molecules that fold onto themselves to form hairpin chromosomes. The *Boodlea* chloroplast genes are highly divergent from their corresponding orthologs, and display an alternative genetic code. The origin of this highly deviant chloroplast genome likely occurred before the emergence of the Cladophorales, and coincided with an elevated transfer of chloroplast genes to the nucleus. A chloroplast genome that is composed only of linear DNA molecules is unprecedented among eukaryotes and highlights unexpected variation in the plastid genome architecture.

Introduction

Photosynthetic eukaryotes possibly originated 1.9 billion years ago following an endosymbiotic event in which a heterotrophic ancestor of the Archaeplastida engulfed a cyanobacterium that became stably integrated and evolved into a membrane-bound organelle, the plastid [1, 2]. Following this primary endosymbiosis, an intricate history of plastid acquisition via eukaryote-eukaryote endosymbioses resulted in the spread of plastids to distantly related eukaryotic lineages [3].

Plastids have retained a reduced version of the genome inherited from their cyanobacterial ancestor. A core set of genes involved in the light reactions of photosynthesis, ATP generation, and functions related to transcription and translation is typically retained [4]. Many genes have been lost or transferred to the nuclear genome and, as a result, plastids are dependent on nuclear-encoded, plastid-targeted proteins for the maintenance of essential biochemical pathways and other functions such as genome replication, gene expression, and DNA repair [5]. Nearly all plastid genomes consist of a single circular-mapping chromosome, typically between 100-200 kb and encoding circa 80-250 genes [4, 6]. Diversity in size, gene content, density and organisation of plastid genomes among different eukaryotic lineages [7-9] is by and large limited, especially when compared to mitochondria.

While fragmented mitochondrial genomes evolved several times independently during the evolution of eukaryotes [8, 10], fragmented plastid genomes are only known in dinoflagellates [11] and a single green algal species [12]. In peridinin-containing dinoflagellates, the chloroplast genome is fragmented into DNA minicircles of 2-3 kb, most of which carry one gene only [11, 13]. Larger minicircles of up to 12 kb have also been described [14], as well as minicircles containing two genes [15], and 'empty' minicircles without genes [16]. The genes located on these minicircles mostly encode key components of the major photosynthetic complexes, including subunits of photosystems I and II, the cytochrome b6f complex, and ATP synthase, as well as rRNAs and a few tRNAs [11]. The only other alga with a fragmented chloroplast genome is the green alga *Koshicola spirodelophila*, but here the level of fragmentation is minor: the plastid genome is divided into three large circular chromosomes totalling 385 kb, with a gene content comparable to other green algae [12]. In addition, plastid minicircles that coexist with a conventional plastid genome have been observed in a few algae, including dinoflagellates with haptophyte-derived plastids [17] and the green alga *Acetabularia* [18, 19].

Although plastid genomes generally assemble as circular-mapping DNAs, they can take multiple complex conformations *in vivo*, including multigenomic, linear-branched structures with discrete termini [20, 21]. The alveolate *Chromera velia* is the only known alga with a linear-mapping plastid genome with telomeric arrangement [22], and is also atypical in that several core photosynthesis genes are fragmented. Linear plastid genomes, however, may be more widespread as several plastid genomes currently do not map as a circle [23].

Currently, and in stark contrast to other algae [9, 24-26], little is known about the gene content and structure of the chloroplast genome in the Cladophorales (Ulvophyceae), an ecologically important group of marine and freshwater green algae, which includes several hundreds of species. These macroscopic multicellular algae have giant, multinucleate cells containing numerous chloroplasts (Figure 1A-C). Most attempts to amplify common chloroplast genes have failed [27], with only one highly divergent *rbcL* sequence published thus far for *Chaetomorpha valida* [28]. An atypical plastid genome in the Cladophorales is suggested by the presence of abundant plasmid-like DNA that has been observed in the chloroplasts of several species [29, 30]. These plasmids-like DNA molecules represent a Low Molecular Weight (LMW) DNA fraction, visible on agarose gels of total DNA extracts (Figure 1D). Pioneering work revealed that these structures are single-stranded DNA (ssDNA) molecules of 1.5-3.0 kb that fold in a hairpin configuration and lack sequence similarity to the nuclear DNA [31, 32]. Some of the hairpin-like DNAs contain putatively transcribed sequences with similarity to chloroplast genes encoding subunits of Photosystems I and II (*psaB*, *psbB*, *psbC* and *psbF*) [31].

Here, we describe intriguing features of the plastid genome of Cladophorales, focusing on *Boodlea composita*. Through the integration of different DNA sequencing methods, combined with RNA sequencing, we found that chloroplast protein-coding genes are highly expressed and encoded on 1-7 kb linear single-stranded DNA molecules. Due to the wide-spread presence of inverted repeats, these molecules fold into a hairpin configuration. A chloroplast genome that is composed only of linear DNA molecules is unprecedented among eukaryotes and highlights unexpected variation in plastid genome architecture.

Results and Discussion

DNA and RNAseq data

Our reconstruction of the chloroplast genome of *Boodlea composita* is based on different high-throughput DNA sequencing methods (Figure S1, Experimental Procedures). The choice of short read DNA sequencing of isolated intact chloroplasts (chloroplast-enriched fraction) using Roche 454 technology was based on comparable sequencing approaches in other plants and algae that successfully resulted in assembly of chloroplast genomes [25]. To overcome possible assembly artefacts in a hypothetical scenario of an inflated chloroplast genome bloated by repetitive elements, long-read sequencing of the High Molecular Weight (HMW) DNA fraction using Pacific Biosciences Single-Molecule Real-Time (SMRT) method was applied, while long read sequencing of the LMW DNA fraction allowed characterization of the previously observed plasmid-like DNA in the chloroplast [29, 30]. To allow comparison of the results of *Boodlea* with other species of Cladophorales, we generated additional DNA sequence data from nine other species using Illumina HiSeq 2000 technology. Finally, two deep-coverage RNA-seq libraries, a total-RNA library and a mRNA library enriched for nuclear transcripts, were generated to confirm the transcription of genes, and to inform whether genes are nuclear versus plastid encoded.

A prodigious chloroplast genome with reduced gene set

Assembly of the chloroplast-enriched DNA reads generated using Roche 454 technology did not result in a typical circular chloroplast genome. Instead, 21 chloroplast protein-coding genes were found on 58 short contigs (1,203-5,426 bp): *atpA*, *atpB*, *atpH*, *atpI*, *petA*, *petB*, *petD*, *psaA*, *psaB*, *psaC*, *psbA*, *psbB*, *psbC*, *psbD*, *psbE*, *psbF*, *psbJ*, *psbK*, *psbL*, *psbT* and *rbcL*. All but the *rbcL* gene code for components of the major thylakoid transmembrane protein complexes (ATP synthase, cytochrome b6f, and photosystems I and II). The contigs contained inverted repeats at their termini and, despite high coverage by sequence reads, they could not be extended by iterative contig extension. Sequence similarity searches and a metagenomic binning approach (distribution analysis of 4-mers) demonstrated that the inverted repeats were also found on

contigs with no sequence similarity to known proteins, raising the number of contigs of chloroplast origin to 136. These contigs are further referred to as “chloroplast 454 contigs”. The length distribution of the chloroplast 454 contigs was consistent with the size of the LMW DNA fraction as estimated by agarose gel electrophoresis of *Boodlea* genomic DNA (Figures 1D, S6A, S6B).

The failure to assemble a circular chloroplast genome might be due to repetitive elements that impair the performance of short-read assemblers [33]. Inflated chloroplast genome bloated by repetitive elements have been documented in several green algae [34-36]. To overcome assembly artefacts and close putative gaps in the chloroplast 454 contigs, we applied Single-Molecule Real-Time (SMRT) sequencing (Pacific Biosciences) to the HMW and LMW DNA fractions. Only 22 HMW DNA reads (ca. 0.044 %) harboured protein-coding genes commonly present in chloroplast genomes of Archaeplastida (Figure 2A). All but three of these genes (*psbA*, *psbB* and *psbC*, which likely correspond to carry-over LMW DNA) contained spliceosomal introns, were absent in the chloroplast 454 contigs, and revealed a high ratio between mapped mRNA and total-RNA reads, altogether suggesting that they are encoded in the nucleus (Figure S2A). Conversely, 22 chloroplast genes (that is, the 21 protein-coding genes identified in chloroplast 454 contigs as well as the 16S rRNA gene) were found in the LMW DNA reads (Figure 2A). An orthology-guided assembly, where the chloroplast 454 contigs harbouring protein-coding genes guided the assembly of LMW DNA reads with sequence similarity to chloroplast genes, resulted in 34 contigs between 1,179 and 6,925 bp in length, henceforth referred to as “chloroplast genome” (Figure 2B, Table S1).

Four contigs of the *Boodlea* chloroplast genome (contigs 10, 19, 32 and 33) display long palindromic sequences that include full-length coding sequences (CDSs), and a less conserved tail region (Figure 2B). The remaining contigs have similar palindromic structures but appear to be not completely assembled. Such palindromes allow regions of the single-stranded LMW DNA molecules to fold into hairpin-like secondary structures. Additional smaller inverted repeats were identified in many of the contigs (Figure 2B), which may result in more complex secondary structures.

Chloroplast 454 contigs could not be scaffolded with long HMW DNA reads, nor did an hybrid assembly between chloroplast 454 contigs and long HMW DNA reads generate a circular chloroplast genome (Figure S1, Experimental Procedures). The LMW DNA reads are concordant

and consistent with the palindromic sequences of the assembled chloroplast genome, indicating that the palindromes are not a result of assembly artefacts (Figure 3, Table S1). As a consequence, we conclude that the chloroplast genome is not a single large molecule but that it is instead fragmented in several molecules in the LMW DNA.

A chloroplast genome that is entirely fragmented into hairpin chromosomes is in line with earlier observations of abundant LMW DNA in chloroplasts of several species of Cladophorales [29]. The hairpin configuration of the chromosomes derived from our sequence data corresponds with earlier data based on electron microscopy, endo- and exonuclease digestion experiments, acridine orange staining, and denaturing gel electrophoresis [29, 32]. Fluorescence *in situ* hybridization, and Southern blot hybridisation indicated that these plasmid-like DNA molecules are present within the chloroplast only [30], supporting the congruence between chloroplast 454 contigs and sequences from the LMW fraction (Figure 3, Supplementary Figure S7A).

The chloroplast genome contigs of *Boodlea* feature an exceptionally high GC-content, ranging from 54 to 60 % in the gene-containing contigs (average 57%) (Table S1). These values are concordant with the high density of the LMW fraction observed in CsCl/bisbenzimidazole gradients [29], and also with sequence data from cloned plasmids of *Ernodesmis* (51-59% GC) [31]. Plastid genomes are generally AT-rich, and in green algal species GC-content typically ranges between 26 and 43% [9, 37]. GC-rich plastid genomes are very rare, but higher values have been reported for the trebouxiophycean green algae *Coccomyxa subellipsoidea*, *Paradoxia multiseta* (both 51% GC), and Trebouxiophyceae sp. MX-AZ01 (58%) [24]. These species, however, feature standard plastid genome architectures.

The size of the *Boodlea* chloroplast genome could not be estimated by inspection of k-mer frequency distributions of the reads in the 454 library nor from those of the uncorrected and corrected LMW DNA reads [38]. Histograms of k-mer frequency distributions revealed several small peaks, indicating a heterogeneous population of molecules present in different stoichiometries, and the signal to noise ratio was too small to make a comfortable estimation of the sizes (Figure S3). The cumulative length of the 34 *Boodlea* chloroplast genome contigs is 91 kb (Table S3). However, if we would consider the large and heterogeneous population of LMW DNA reads bearing no similarity to protein-coding genes (“empty” hairpin chromosomes, see below) as part of the chloroplast genome, its size could be regarded as much larger.

The largest known circular-mapping chloroplast genomes have been documented in the red algae *Bulboplastis apyrenoidosa* (610 kb) and *Corynoplastis japonica* (1.127 Mb), where the genomes are bloated by group II introns and include transposable elements of possible bacterial origin [8]. Within green algae, expanded chloroplast genomes have been reported in two distinct clades: the Chlorophyceae and the Ulvophyceae. Inflation of the 521-kb chloroplast genome of *Floydiella terrestris* (Chlorophyceae) resulted mainly from the proliferation of dispersed, heterogeneous repeats (>30 bp) in intergenic regions, representing more than half of the genome length [36]. Intergenic regions of the *Volvox carteri* (Chlorophyceae) chloroplast genome, instead, are populated with short palindromic repeats (average size of 50 bp) that constitute ca. 64% of the predicted 525-kb genome [34]. The mechanisms by which such palindromic selfish DNA spread throughout the *Volvox* chloroplast genome is not clear, but the presence of a reverse transcriptase and endonuclease may point toward retrotranscription [34, 39]. For *Acetabularia acetabulum* (Ulvophyceae), the chloroplast genome was sequenced only partially and its size was estimated to exceed 1 Mb; it has exceptionally long intergenic regions and features long repetitive elements (>10 kb) arranged in tandem [35, 40]. The *Boodlea* chloroplast genome is rich as well in non-coding DNA, constituting 92.2% of the 136 chloroplast 454 contigs and 72.8% of the assembled chloroplast genome, comparable to that in inflated chloroplast genomes of other green algae (*Floydiella terrestris*, 82.1%; *Volvox carteri*, ca. 80%; *Acetabularia acetabulum*, ca. 87% of the sequenced chloroplast genome) [34, 35].

The non-coding DNA regions (ncDNA) of the hairpins showed high sequence similarity among one another (52.5-100% sequence similarity). Within the ncDNA, we identified six conserved motifs, 20 to 35 bp in length and with a GC-content ranging from 36 to 84%, which lack similarity to known regulatory elements (Figure 4). Motifs 1, 2 and 5 were always present upstream of the start codon of the chloroplast genes, occasionally in more than one copy. Although their distances from the start codon were variable, their orientations relative to the gene were conserved, indicating a potential function as a regulatory element of gene expression and/or replication of the hairpin chromosomes.

These motifs were also present in 1,966 (ca. 1.8 %) LMW DNA reads lacking genes. This observation supports earlier findings of abundant non-coding LMW DNA molecules in the Cladophorales [29, 31]. In contrast, a very small fraction of the HMW DNA reads (15 corrected reads) displayed the same ncDNA motifs and these were present exclusively on long terminal repeat retrotransposons (RT-LTRs) (Figures S2D, S2E). RT-LTRs were also abundant in the 454

contigs (Figure S6D). The abundance of RT-LTRs in the 454 contigs and the presence of ncDNA motifs in both the *Boodlea* chloroplast genome and nuclear RT-LTRs is suggestive of DNA transfer between the nucleus and chloroplast and may allude to the origin of the hairpin chromosomes. Hypothetically, an invasion of nuclear RT-LTRs in the chloroplast genome may have resulted in an expansion of the chloroplast genome and its subsequent fragmentation into hairpin chromosomes during replication. Chloroplast genome fragmentation could be caused by recombination between repetitive elements and displacement of the palindromic sequences from the lagging strand during the chloroplast genome replication [41, 42], and it is consistent with the expectation that recombination and cleavage of repetitive DNA will produce a heterogeneous population of molecules, as observed in dinoflagellates plastid genomes [11], and in the *Boodlea* LMW DNA.

A fragmented chloroplast genome is a common feature of Cladophorales

DNA sequence data were obtained from 9 additional Cladophorales species, representing the main lineages of the order: *Chaetomorpha aerea*, *Cladophora albida*, *C. socialis*, *C. vadorum*, *Dictyosphaeria cavernosa*, *Pithophora* sp., *Siphonocladus tropicus*, *Struvea elegans*, and *Valonia utricularis* (Tables S2 and S4). Although comparable sequencing approaches resulted in the assembly of circular chloroplast genomes for other algae, including green seaweeds [43, 44], only short chloroplast contigs (ca. 200-8,000 bp) were assembled from these libraries, similar to *Boodlea composita*. Interestingly, a similar set of chloroplast genes was identified in all sequenced Cladophorales species (Table S5). In contrast to the genes found in the *Boodlea* hairpin chromosomes, however, most of the chloroplast genes identified in the additional Cladophorales libraries were fragmented, possibly due to assembly of the shorter Illumina reads (Table S3). These findings support the idea that fragmentation of the chloroplast genome occurred before or early in the evolution of the Cladophorales.

Highly divergent chloroplast genes

The 21 chloroplast protein-coding genes of *Boodlea* and the other species of Cladophorales display extremely high sequence divergence compared to orthologous genes in other photosynthetic organisms (Figure 5). A maximum likelihood phylogenetic tree based on a

concatenated amino acid alignment of 19 chloroplast genes from Archaeplastida and Cyanobacteria species (Figure 5) shows that despite their high divergence, the Cladophorales sequences form a monophyletic group within the core Chlorophyta (Figure S4), a position that is supported by phylogenetic analyses of nuclear genes [45]. The high sequence divergence of chloroplast genes in the Cladophorales supports the notion that organellar genomes with extremely derived architectures, including those of peridinin-containing dinoflagellates, also tend to fall at the extreme ends of the range observed at the mutation rate (or gene sequence divergence) level [7, 46].

For some *Boodlea* chloroplast genes, the identification of start and stop codons was uncertain and a non-canonical genetic code was identified (Figure 6). The canonical stop codon UGA was found 11 times internally in six genes (*petA*, *psaA*, *psaB*, *psaC*, *psbC* and *rbcL*), but was also present as a genuine termination codon in several genes, *petA* and *psaA* included. At seven of these 11 positions, the corresponding amino acid residue in orthologous genes was conserved (i.e. present in more than 75% of the taxa in the alignment), but different amino acids were observed at these positions: V, Q, I, L and C (Figures 6A and 6B). The reassignment of the stop codon UGA to C has been documented in the nuclear genetic code of several species of ciliates [47]. For the remaining positions, the amino acid in the alignment was not conserved, and therefore the amino acid coded by the UGA codon could not be determined with certainty.

Deviations from the universal genetic code are widespread among mitochondrial genomes, and include loss of start and stop codons in some groups, including dinoflagellates [11, 48, 49]. In contrast, non-canonical genetic codes are much rarer in plastid genomes, and up to now have only been detected in the apicomplexans *Neospora caninum* [50], *Chromera velia* [22], and the dinoflagellate *Lepidodinium chlorophorum* [51]. In genomes of primary plastids, a non-canonical genetic code is unprecedented.

Dual meaning of UGA as both stop and sense codons has recently been reported from a number of unrelated protists [47, 52, 53]. While in *Saccharomyces cerevisiae* the tetranucleotide UGA-C allows increased incorporation of the near-cognate Cys-tRNA for the UGA premature termination codon [54], such preference was not observed in the *Boodlea* chloroplasts. Importantly, a non-canonical genetic code has also been described for Cladophorales nuclear genes, where UAG and UAA codons are reassigned to glutamine [55], which implies two independent departures from the standard genetic code in a single organism.

Unexpectedly, we found that the 16S rRNA gene in the *Boodlea* chloroplast genome is split across two distinct hairpin chromosomes and that its size is much smaller compared to its algal and bacterial homologs and (Figures 2B and 7). Fragmentation of rRNA genes has been observed in organellar genomes, including Apicomplexa, dinoflagellates, and many green algae [10]. In general, fragmentation of protein-coding and rRNA genes is more common in mitochondrial genomes than in plastid genomes [17, 22, 56]. Despite considerable effort, we could not detect the 23S rRNA gene nor the 5S rRNA gene.

The transcription of the aberrant chloroplast genes was confirmed using RNA-seq, and is concordant with previous results of Northern blots [31]. Transcripts of 21 chloroplast genes (that is, 20 protein-coding genes as well as the 16S rRNA gene) were identical to the genes encoded by the chloroplast 454 contigs (Figure 2A; 3 and S5), providing evidence for the absence of RNA editing and corroborating the use of a non-canonical genetic code (Figure S5). Lack of RNA editing was also evidenced for the 11 internal occurrences of UGA (Figure S5). This observation, in combination with conservation of the sequence after the UGA codon, serves as evidence that it is not a termination codon but an alternative code. The high total-RNA to mRNA ratio observed for reads that mapped to the chloroplast 454 contigs confirmed that these genes were not transcribed in the nucleus (Figure S6C). All coding sequences of the same protein-coding genes found on different contigs of the *Boodlea* chloroplast genome were expressed, despite minor differences in their nucleotide sequences (Table S1).

Additional transcripts of 66 genes that have been located in the chloroplast in other Archaeplastida were identified (Figure 2A). Although their subcellular origin was not determined experimentally, they are probably all nuclear-encoded, based on high mRNA to total-RNA reads ratio and their presence on High Molecular Weight (HMW) DNA reads.

Conclusions

We collected several lines of evidence indicating that *Boodlea composita* lacks a typical large circular chloroplast genome. The chloroplast genome is instead fragmented into multiple linear hairpin chromosomes, and has a highly reduced gene repertoire compared to other chloroplast genomes. Thirty-four hairpin chromosomes were identified, harbouring 21 protein-coding genes and the 16S rRNA gene, which are highly divergent in sequence compared to orthologs in other algae, and display an alternative genetic code. The exact set of *Boodlea* chloroplast genes remains elusive, but at least 19 genes coding for chloroplast products appear to be nuclear-encoded, of which nine are always chloroplast-encoded in related green algae (Figure 2A). This suggests that fragmentation of a conventional chloroplast genome in the Cladophorales has been accompanied with an elevated transfer of genes to the nucleus, similarly to the situation in peridinin-containing dinoflagellates [11], with plastid genomes encoding about 12 genes or less [11, 57]. Notably, the two distantly related algal groups have converged on a very similar gene distribution: chloroplast genes code only for the subunits of photosynthetic complexes (and also for Rubisco in *Boodlea*), whereas the expression machinery appears to be fully nucleus-encoded (Figure 2A). Other nonstandard chloroplast genome architectures have recently been observed, such as a monomeric linear chromosome in the alveolate microalga *Chromera velia* [22] and three circular chromosomes in the green alga *Koshicola spirodelophila* [12], but these represent relatively small deviations from the paradigm, when compared to the chloroplast genome of the Cladophorales. The highly fragmented chloroplast genome in the Cladophorales is wholly unprecedented and will be of significance to understanding processes driving organellar genome fragmentation and gene reduction, endosymbiotic gene transfer, and the minimal functional chloroplast gene set.

Supplemental Information

Supplementary information includes seven figures and six tables can be found with this article online at

Author Contributions

Conceptualization and methodology, F.L., M.T., C.B., H.V., K.V. and O.D.C; Resources, F.L., C.B., A.D.C., J.M.L-B., K.V. and O.D.C; Formal Analysis and investigation, A.D.C., F.L., K.A.B., J.J., K.V. and O.D.C; Writing – Original Draft and visualization, A.D.C., F.L., K.V. and O.D.C; Writing – Review & Editing, A.D.C., F.L., K.A.B., M.T., J.J., H.V., K.V. and O.D.C.

Acknowledgments

We thank Ellen Nisbett, Christopher J. Howe, Bram Verhelst, Sven Gould, Joe Zuccarello, and John W. La Claire II for help and advice. This work was supported by Ghent University BOF/01J04813, the Australian Research Council (DP150100705) to H.V., and the National Science Foundation (GRAToL 10136495) to J.L.B

References

1. Ponce-Toledo, R.I., Deschamps, P., López-García, P., Zivanovic, Y., Benzerara, K., and Moreira, D. (2017). An early-branching freshwater cyanobacterium at the origin of plastids. *Curr. Biol.* **27**, 386–391.
2. Sánchez-Baracaldo, P., Raven, J.A., Pisani, D., and Knoll, A.H. (2017). Early photosynthetic eukaryotes inhabited low-salinity habitats. *Proc. Natl. Acad. Sci. USA* **114**, E7737-E7745.
3. Keeling, P.J. (2010). The endosymbiotic origin, diversification and fate of plastids. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **365**, 729-748.
4. Green, B.R. (2011). Chloroplast genomes of photosynthetic eukaryotes. *The Plant J.* **66**, 34-44.
5. Kleine, T., Maier, U.G., and Leister, D. (2009). DNA transfer from organelles to the nucleus: the idiosyncratic genetics of endosymbiosis. *Annu. Rev. Plant Biol.* **60**, 115-138.
6. Lang, B.F., and Nedelcu, A.M. (2012). Plastid genomes of algae. In *Genomics of chloroplasts and mitochondria*, Volume 35, R. Bock and V. Knoop, eds. (Dordrecht: Springer Netherlands), pp. 59-87.
7. Simpson, C.L., and Stern, D.B. (2002). The treasure trove of algal chloroplast genomes. Surprises in architecture and gene content, and their functional implications. *Plant Physiol.* **129**, 957-966.
8. Smith, D.R., and Keeling, P.J. (2015). Mitochondrial and plastid genome architecture: Reoccurring themes, but significant differences at the extremes. *Proc. Natl. Acad. Sci. USA*, **112**, 10177-10184.
9. Muñoz-Gómez, S.A., Mejía-Franco, F.G., Durnin, K., Colp, M., Grisdale, C.J., Archibald, J.M., and Slamovits, C.H. (2017). The new red algal subphylum Proteorhodophytina comprises the largest and most divergent plastid genomes known. *Curr. Biol.* **27**, 1677-1684. e1674.
10. Barbrook, A.C., Howe, C.J., Kurniawan, D.P., and Tarr, S.J. (2010). Organization and expression of organellar genomes. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **365**, 785-797.
11. Howe, C.J., Nisbet, R.E.R., and Barbrook, A.C. (2008). The remarkable chloroplast genome of dinoflagellates. *J. Exp. Bot.* **59**, 1035-1045.
12. Watanabe, S., Fučíková, K., Lewis, L.A., and Lewis, P.O. (2016). Hiding in plain sight: *Koshicola spirodelophila* gen. et sp. nov. (Chaetopeltidales, Chlorophyceae), a novel green alga associated with the aquatic angiosperm *Spirodela polyrhiza*. *Am. J. Bot.* **103**, 865-875.
13. Zhang, Z., Green, B., and Cavalier-Smith, T. (1999). Single gene circles in dinoflagellate chloroplast genomes. *Nature* **400**, 155-159.
14. Nelson, M.J., and Green, B.R. (2005). Double hairpin elements and tandem repeats in the non-coding region of *Adenoides eludens* chloroplast gene minicircles. *Gene* **358**, 102-110.
15. Laatsch, T., Zauner, S., Stoebe-Maier, B., Kowallik, K., and Maier, U.-G. (2004). Plastid-derived single gene minicircles of the dinoflagellate *Ceratium horridum* are localized in the nucleus. *Mol. Biol. Evol.* **21**, 1318-1322.
16. Hiller, R.G. (2001). 'Empty' minicircles and *petB/atpA* and *psbD/psbE* (*cytb 559 α*) genes in tandem in *Amphidinium carterae* plastid DNA. *FEBS Lett.* **505**, 449-452.
17. Espelund, M., Minge, M.A., Gabrielsen, T.M., Nederbragt, A.J., Shalchian-Tabrizi, K., Otis, C., Turmel, M., Lemieux, C., and Jakobsen, K.S. (2012). Genome fragmentation is not confined to the peridinin plastid in dinoflagellates. *PLoS One* **7**, e38809.
18. Green, B.R. (1976). Covalently closed minicircular DNA associated with *Acetabularia* chloroplasts. *iochim. Biophys. Acta, Nucleic Acids Protein Synth.* **447**, 156-166.
19. Ebert, C., Tymms, M.J., and Schweiger, H.-G. (1985). Homology between 4.3 μm minicircular and plastomic DNA in chloroplasts of *Acetabularia cliftonii*. *Mol. Gen. Genet.* **200**, 187-192.
20. Bendich, A.J. (2007). The size and form of chromosomes are constant in the nucleus, but highly variable in bacteria, mitochondria and chloroplasts. *Bioessays* **29**, 474-483.

21. Oldenburg, D.J., and Bendich, A.J. (2016). The linear plastid chromosomes of maize: terminal sequences, structures, and implications for DNA replication. *Curr. Genet.* *62*, 431-442.
22. Janouškovec, J., Sobotka, R., Lai, D.-H., Flegontov, P., Koník, P., Komenda, J., Ali, S., Prášil, O., Pain, A., and Oborník, M. (2013). Split photosystem protein, linear-mapping topology and growth of structural complexity in the plastid genome of *Chromera velia*. *Mol. Biol. Evol.* *30*, 2447-2462.
23. Gabrielsen, T.M., Minge, M.A., Espelund, M., Tooming-Klunderud, A., Patil, V., Nederbragt, A.J., Otis, C., Turmel, M., Shalchian-Tabrizi, K., and Lemieux, C. (2011). Genome evolution of a tertiary dinoflagellate plastid. *PLoS One* *6*, e19132.
24. Turmel, M., Otis, C., and Lemieux, C. (2015). Dynamic Evolution of the Chloroplast Genome in the Green Algal Classes Pedinophyceae and Trebouxiophyceae. *Genome Biol. Evol.* *7*, 2062-2082.
25. Lemieux, C., Otis, C., and Turmel, M. (2016). Comparative chloroplast genome analyses of streptophyte green algae uncover major structural alterations in the Klebsormidiophyceae, Coleochaetophyceae and Zygnematophyceae. *Front. Plant Sci.* *7*, 697.
26. Leliaert, F., Tronholm, A., Lemieux, C., Turmel, M., DePriest, M.S., Bhattacharya, D., Karol, K.G., Fredericq, S., Zechman, F.W., and Lopez-Bautista, J.M. (2016). Chloroplast phylogenomic analyses reveal the deepest-branching lineage of the Chlorophyta, Palmophyllophyceae class. *nov. Sci. Rep.* *6*, 25367.
27. Fučíková, K., Leliaert, F., Cooper, E.D., Škaloud, P., D'hondt, S., De Clerck, O., Gurgel, F., Lewis, L.A., Lewis, P.O., Lopez-Bautista, J., et al. (2014). New phylogenetic hypotheses for the core Chlorophyta based on chloroplast sequence data. *Front. Ecol. Evol.* *2*, 63.
28. Deng, Y., Zhan, Z., Tang, X., Ding, L., and Duan, D. (2013). Molecular cloning and expression analysis of *rbcl* cDNA from the bloom-forming green alga *Chaetomorpha valida* (Cladophorales, Chlorophyta). *J. Appl. Phycol.* *26*, 1853-1861.
29. La Claire, J.W., Zuccarello, G.C., and Tong, S. (1997). Abundant plasmid-like DNA in various members of the orders Siphonocladales and Cladophorales (Chlorophyta). *J. Phycol.* *33*, 830-837.
30. La Claire, J.W., and Wang, J. (2000). Localization of plasmidlike DNA in giant-celled marine green algae. *Protoplasma* *213*, 157-164.
31. La Claire, J.W., Loudenslager, C.M., and Zuccarello, G.C. (1998). Characterization of novel extrachromosomal DNA from giant celled marine green algae. *Curr. Genet.* *34*, 204-211.
32. La Claire, J.W., and Wang, J.S. (2004). Structural characterization of the terminal domains of linear plasmid-like DNA from the green alga *Ernodesmis* (Chlorophyta). *J. Phycol.* *40*, 1089-1097.
33. Miller, J.R., Koren, S., and Sutton, G. (2010). Assembly algorithms for next-generation sequencing data. *Genomics* *95*, 315-327.
34. Smith, D.R., and Lee, R.W. (2009). The mitochondrial and plastid genomes of *Volvox carteri*: bloated molecules rich in repetitive DNA. *BMC Genomics* *10*, 132.
35. de Vries, J., Habicht, J., Woehle, C., Huang, C., Christa, G., Wägele, H., Nickelsen, J., Martin, W.F., and Gould, S.B. (2013). Is *ftsH* the key to plastid longevity in sacoglossan slugs? *Genome Biol. Evol.* *5*, 2540-2548.
36. Brouard, J.-S., Otis, C., Lemieux, C., and Turmel, M. (2010). The exceptionally large chloroplast genome of the green alga *Floydiella terrestris* illuminates the evolutionary history of the Chlorophyceae. *Genome Biol. Evol.* *2*, 240-256.
37. Leliaert, F., Smith, D.R., Moreau, H., Herron, M.D., Verbruggen, H., Delwiche, C.F., and De Clerck, O. (2012). Phylogeny and molecular evolution of the green algae. *Crit. Rev. Plant Sci.* *31*, 1-46.
38. Hozza, M., Vinař, T., and Brejová, B. (2015). How big is that genome? Estimating genome size and coverage from k-mer abundance spectra. In *String Processing and Information Retrieval*:

- 22nd International Symposium, SPIRE 2015, London, UK, September 1-4, 2015, Proceedings, C. Iliopoulos, S. Puglisi and E. Yilmaz, eds. (Cham: Springer International Publishing), pp. 199-209.
39. Burt, A.T., Robert (2006). *Genes in conflict: the biology of selfish genetic elements*, (Cambridge: Harvard University Press).
 40. Tymms, M.J., and Schweiger, H.-G. (1985). Tandemly repeated nonribosomal DNA sequences in the chloroplast genome of an *Acetabularia mediterranea* strain. *Proc. Natl. Acad. Sci. USA* *82*, 1706-1710.
 41. Bikard, D., Loot, C., Baharoglu, Z., and Mazel, D. (2010). Folded DNA in action: hairpin formation and biological functions in Prokaryotes. *Microbiol. Mol. Biol. Rev.* *74*, 570-588.
 42. Ellis, T.H.N., and Day, A. (1986). A hairpin plastid genome in barley. *EMBO J.* *5*, 2769-2774.
 43. Leliaert, F., and Lopez-Bautista, J.M. (2015). The chloroplast genomes of *Bryopsis plumosa* and *Tydemania expeditionis* (Bryopsidales, Chlorophyta): compact genomes and genes of bacterial origin. *BMC Genomics* *16*, 204.
 44. Marcelino, R.V., Cremen, M.C.M., Jackson, C.J., Larkum, A.A.W., and Verbruggen, H. (2016). Evolutionary dynamics of chloroplast genomes in low light: a case study of the endolithic green alga *Ostreobium quekettii*. *Genome Biol. Evol.* *8*, 2939-2951.
 45. Cocquyt, E., Verbruggen, H., Leliaert, F., and De Clerck, O. (2010). Evolution and cytological diversification of the green seaweeds (Ulvoophyceae). *Mol. Biol. Evol.* *27*, 2052-2061.
 46. Zhang, Z., Green, B.R., and Cavalier-Smith, T. (2000). Phylogeny of ultra-rapidly evolving dinoflagellate chloroplast genes: a possible common origin for sporozoan and dinoflagellate plastids. *J. Mol. Evol.* *51*, 26-40.
 47. Heaphy, S.M., Mariotti, M., Gladyshev, V.N., Atkins, J.F., and Baranov, P.V. (2016). Novel ciliate genetic code variants including the reassignment of all three stop codons to sense codons in *Condyllostoma magnum*. *Mol. Biol. Evol.* *33*, 1537-1719.
 48. Waller, R.F., and Jackson, C.J. (2009). Dinoflagellate mitochondrial genomes: stretching the rules of molecular biology. *Bioessays* *31*, 237-245.
 49. Slamovits, C.H., Saldarriaga, J.F., Larocque, A., and Keeling, P.J. (2007). The highly reduced and fragmented mitochondrial genome of the early-branching dinoflagellate *Oxyrrhis marina* shares characteristics with both apicomplexan and dinoflagellate mitochondrial genomes. *J. Mol. Biol.* *372*, 356-368.
 50. Lang-Unnasch, N., and Aiello, D. (1999). Sequence evidence for an altered genetic code in the *Neospora caninum* plastid. *Int. J. Parasitol.* *29*, 1557-1562.
 51. Matsumoto, T., Ishikawa, S.A., Hashimoto, T., and Inagaki, Y. (2011). A deviant genetic code in the green alga-derived plastid in the dinoflagellate *Lepidodinium chlorophorum*. *Mol. Phylogen. Evol.* *60*, 68-72.
 52. Zahonova, K., Kostygov, A.Y., Sevcikova, T., Yurchenko, V., and Elias, M. (2016). An unprecedented non-canonical nuclear genetic code with all three termination codons reassigned as sense codons. *Curr. Biol.* *26*, 1879-0445.
 53. Swart, Estienne C., Serra, V., Petroni, G., and Nowacki, M. (2016). Genetic Codes with No Dedicated Stop Codon: Context-Dependent Translation Termination. *Cell* *166*, 691-702.
 54. Beznoskova, P., Gunisova, S., and Valasek, L.S. (2016). Rules of UGA-N decoding by near-cognate tRNAs and analysis of readthrough on short uORFs in yeast. *RNA* *22*, 456-466.
 55. Cocquyt, E., Gile, G., Leliaert, F., Verbruggen, H., Keeling, P., and De Clerck, O. (2010). Complex phylogenetic distribution of a non-canonical genetic code in green algae. *BMC Evol. Biol.* *10*, 327.
 56. Smith, D.R., Lee, R.W., Cushman, J.C., Magnuson, J.K., Tran, D., and Polle, J.E.W. (2010). The *Dunaliella salina* organelle genomes: large sequences, inflated with intronic and intergenic DNA. *BMC Plant Biol.* *10*, 83.

57. Barbrook, A.C., Voolstra, C.R., and Howe, C.J. (2014). The chloroplast genome of a *Symbiodinium* sp. clade C3 isolate. *Protist* 165, 1-13.
58. Andersen, R.A. (2005). *Algal culturing techniques*, (Amsterdam: Elsevier).
59. Doyle, J.J., and Doyle, J.L. (1987). A rapid DNA isolation procedure for small quantities of fresh leaf tissue. *Phytochemical Bulletin* 19, 11-15.
60. Palmer, J.D. (1982). Physical and gene mapping of chloroplast DNA from *Atriplex triangularis* and *Cucumis sativa*. *Nucleic Acids Res.* 10, 1593-1605.
61. Chevreux, B., Pfisterer, T., Drescher, B., Driesel, A.J., Muller, W.E., Wetter, T., and Suhai, S. (2004). Using the miraEST assembler for reliable and automated mRNA transcript assembly and SNP detection in sequenced ESTs. *Genome Res.* 14, 1147-1159.
62. Boratyn, G.M., Camacho, C., Cooper, P.S., Coulouris, G., Fong, A., Ma, N., Madden, T.L., Matten, W.T., McGinnis, S.D., Merezuk, Y., et al. (2013). BLAST: a more efficient report with usability improvements. *Nucleic Acids Res.* 41, W29-W33.
63. Lin, H.-H., and Liao, Y.-C. (2016). Accurate binning of metagenomic contigs via automated clustering sequences using information of genomic signatures and marker genes. *Sci. Rep.* 6, 24175.
64. Nawrocki, E.P., and Eddy, S.R. (2013). Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics* 29, 2933-2935.
65. Ruby, J.G., Bellare, P., and Derisi, J.L. (2013). PRICE: software for the targeted assembly of components of (Meta) genomic sequence data. *G3* 3, 865-880.
66. Rice, P., Longden, I., and Bleasby, A. (2000). EMBOSS: The European Molecular Biology Open Software Suite. *Trends Genet.* 16, 276-277
67. Noé, L., and Kucherov, G. (2005). YASS: enhancing the sensitivity of DNA similarity search. *Nucleic Acids Res.* 33, W540-W543.
68. Bashir, A., Klammer, A.A., Robins, W.P., Chin, C.S., Webster, D., Paxinos, E., Hsu, D., Ashby, M., Wang, S., Peluso, P., et al. (2012). A hybrid approach for the automated finishing of bacterial genomes. *Nat. Biotechnol.* 30, 701-707.
69. Ye, C., Hill, C.M., Wu, S., Ruan, J., and Ma, Z. (2016). DBG2OLC: efficient assembly of large genomes using long erroneous reads of the third generation sequencing technologies. *Sci. Rep.* 6, 31900.
70. Hackl, T., Hedrich, R., Schultz, J., and Förster, F. (2014). proofread: large-scale high-accuracy PacBio correction through iterative short read consensus. *Bioinformatics* 30, 3004-3011.
71. Vandepoele, K., Van Bel, M., Richard, G., Van Landeghem, S., Verhelst, B., Moreau, H., Van de Peer, Y., Grimsley, N., and Piganeau, G. (2013). pico-PLAZA, a genome database of microbial photosynthetic eukaryotes. *Environ. Microbiol.* 15, 2147-2153.
72. Koren, S., Harhay, G.P., Smith, T.P., Bono, J.L., Harhay, D.M., Mcvey, S.D., Radune, D., Bergman, N.H., and Phillippy, A.M. (2013). Reducing assembly complexity of microbial genomes with single-molecule sequencing. *Genome Biol.* 14, 1-16.
73. Marçais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* 27, 764-770.
74. Berlin, K., Koren, S., Chin, C.-S., Drake, J.P., Landolin, J.M., and Phillippy, A.M. (2015). Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat. Biotechnol.* 33, 623-630.
75. Eddy, S.R. (2011). Accelerated profile HMM searches. *PLoS Comp. Biol.* 7, 10.
76. Rutherford, K., Parkhill, J., Crook, J., Horsnell, T., Rice, P., Rajandream, M.-A., and Barrell, B. (2000). Artemis: sequence visualization and annotation. *Bioinformatics* 16, 944-945.
77. Bailey, T.L., Boden, M., Buske, F.A., Frith, M., Grant, C.E., Clementi, L., Ren, J., Li, W.W., and Noble, W.S. (2009). MEME Suite: tools for motif discovery and searching. *Nucleic Acids Res.* 37, W202-W208.

78. Medina-Rivera, A., Defrance, M., Sand, O., Herrmann, C., Castro-Mondragon, J.A., Delerce, J., Jaeger, S., Blanchet, C., Vincens, P., Caron, C., et al. (2015). RSAT 2015: Regulatory Sequence Analysis Tools. *Nucleic Acids Res.* *43*, W50-56.
79. Mathelier, A., Zhao, X., Zhang, A.W., Parcy, F., Worsley-Hunt, R., Arenillas, D.J., Buchman, S., Chen, C.-y., Chou, A., Ienasescu, H., et al. (2013). JASPAR 2014: an extensively expanded and updated open-access database of transcription factor binding profiles. *Nucleic Acids Res.* *42*, D142-147.
80. Gupta, S., Stamatoyannopoulos, J.A., Bailey, T.L., and Noble, W.S. (2007). Quantifying similarity between motifs. *Genome Biology* *8*, R24.
81. Wu, T.D., Reeder, J., Lawrence, M., Becker, G., and Brauer, M.J. (2016). GMAP and GSNAP for Genomic Sequence Alignment: Enhancements to Speed, Accuracy, and Functionality. In *Statistical Genomics: Methods and Protocols*, E. Mathé and S. Davis, eds. (New York, NY: Springer New York), pp. 283-334.
82. Le Bail, A., Dittami, S.M., de Franco, P.-O., Rousvoal, S., Cock, M.J., Tonon, T., and Charrier, B. (2008). Normalisation genes for expression analyses in the brown alga model *Ectocarpus siliculosus*. *BMC Mol. Biol.* *9*, 75.
83. Grabherr, M.G., Haas, B.J., Yassour, M., Levin, J.Z., Thompson, D.A., Amit, I., Adiconis, X., Fan, L., Raychowdhury, R., Zeng, Q., et al. (2011). Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat. Biotechnol.* *29*, 644-U130.
84. Veeckman, E., Ruttink, T., and Vandepoele, K. (2016). Are we there yet? Reliably estimating the completeness of plant genome sequences. *Plant Cell* *28*, 1759-1768.
85. Bankevich, A., Nurk, S., Antipov, D., Gurevich, A.A., Dvorkin, M., Kulikov, A.S., Lesin, V.M., Nikolenko, S.I., Pham, S., Prjibelski, A.D., et al. (2012). SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.* *19*, 455-477.
86. Talavera, G., and Castresana, J. (2007). Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. *Syst. Biol.* *56*, 564-577.
87. Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* *30*, 1312-1313.
88. Miller, M.A., Pfeiffer, W., and Schwartz, T. (2011). The CIPRES science gateway: a community resource for phylogenetic analyses. In *Proceedings of the 2011 TeraGrid Conference: Extreme Digital Discovery*. (Salt Lake City, Utah: ACM), pp. 1-8.

Main-text figure legends

Figure 1. *Boodlea composita*. (A) Specimen in natural environment. (B) Detail of branching cells. (C) Detail of chloroplasts, each containing a single pyrenoid, and forming a parietal network (the white line is a calcium oxalate crystal). (D) Native agarose gel comparing genomic DNA of *Bryopsis plumosa* (Bryopsidales) and *Boodlea composita* (Cladophorales). Lane 1: 1-kb ladder, sizes in bp; lane 2: *B. plumosa*; lane 3: *B. composita*. High molecular weight (HMW) and low molecular weight (LMW) DNA bands of *B. composita* are indicated. See also Figure S1.

Figure 2. Schematic representation of *Boodlea* chloroplast genome. (A) Distribution of *Boodlea* genes having orthologs in the chloroplast of other Archaeplastida. gDNA (genomic DNA): chloroplast (cp) 454 contigs, HMW and LMW corrected reads; RNA: mRNA and total-RNA assemblies. Asterisks (*) indicate "core" chloroplast genes, i.e. protein-coding genes conserved between chloroplast genomes of Chlorophyta (see Experimental Procedures). The following nine "core" chloroplast genes were not found in any of the *Boodlea* libraries sequenced: *atpF*, *petG*, *petL*, *psaJ*, *psbM*, *psbZ*, *rpl36*, *rps2* and *ycf1*. Grey cells denote putative LMW DNA read contaminants as suggested by the ratios of HMW to LMW DNA reads and mRNA to total-RNA reads (Figures S2B and S2C). (B) Overview of the 34 contigs representing the *Boodlea* chloroplast genome. Purple arrows indicate rRNA genes, red arrows indicate CDSs of protein-coding genes, and blue arrows indicate repetitive elements. For each contig, repetitive elements with similar length indicate similar sequences. Distance between vertical grey lines in the background represents 500 bp. Oga: contig obtained by orthology-guided assembly. 454: chloroplast 454 contig. See also Figure S1, S2, S3, S6 and S7 and Table S1, S3.

Figure 3. LMW DNA reads containing chloroplast genes are expressed, enriched in the total-RNA fraction and congruent to the respective chloroplast 454 contigs. (A) Representation of *petA* LMW DNA read (3,398 bp). The red arrows indicate CDSs, the blue arrows indicate inverted repeats. (B) Corresponding Genome Browser track, from top to bottom: corrected HMW DNA coverage [0], corrected LMW DNA read coverage [range 0-541], 454 read coverage [range 0-567], mRNA library read coverage [range 0-17], assembled mRNA transcripts mapped [0], total-RNA library read coverage [range 0-7,936], and assembled total-RNA transcripts mapped [range 0-17]. (C) Dotplot showing congruence between *petA* LMW DNA read (x axis) and the corresponding *petA*-containing chloroplast 454 contig (y axis, 2,012 bp). Green lines indicate similar sequences; red lines indicate sequences similar to the respective reverse complements. See also Figure S1, S2 and S6.

Figure 4. Conserved non-coding motifs in *Boodlea* LMW DNA. (A) Sequence logos and GC contents of the conserved motifs predicted in the *Boodlea* chloroplast genome. (B) Schematic representation of the distribution of the motifs in the 1,441 bp ncdNA region from the *atpI* group A read used for the identification of additional chloroplast reads in the LMW DNA library. Motifs with conserved orientation relative to the downstream genes are represented by green arrows, while motifs without conserved orientation to the downstream genes are represented by yellow arrows. CDSs are represented by red arrows, inverted repeats are represented by blue arrows.

Figure 5. *Boodlea* chloroplast genes have large sequence divergence. (A) Maximum likelihood phylogenetic tree, with indication of relevant bootstrap values (see also Figure S4). The scale represents 0.5 substitution per amino acid position. (B) Maximum pairwise amino acid sequence distances of the concatenated amino acid alignment within and between clades (*excluding Cladophorales). See also Figure S4 and Table S4, S5 and S6.

Figure 6. Non-canonical genetic code in *Boodlea* chloroplast genes. *Boodlea* chloroplast protein-coding genes were aligned with the respective orthologs of 43 Archaeplastida and 14 Cyanobacteria. (A) Relevant parts of amino acid sequence alignment for six chloroplast genes of *Boodlea* and representatives of Archaeplastida and Cyanobacteria. Positions corresponding to UGA codons in *Boodlea* are indicated by an asterisk. Slashes represent regions of the sequence alignment that were omitted for simplicity. Dots indicate amino acid identity with the top-most sequence. For each gene, position in the alignment is indicated by the numbers shown above the sequence alignment. The numbers below the gene names indicate the eleven positions where UGA was identified as premature termination codon in the six *Boodlea* genes. (B) Sequence logo of the Position Weight Matrix reporting the relative amino acid frequencies in the alignment for each premature termination UGA position in *Boodlea*. See also Figure S5.

Figure 7. The *Boodlea* chloroplast 16S rRNA is fragmented and reduced compared to its algal and bacterial homologs. (A) *Boodlea* chloroplast 16S rRNA sequence was compared with the *E. coli* 16S rRNA secondary structure model [RF00177]. Residues shown in green and red on the *E. coli* model represent the 16S rRNA regions coded by the two hairpin chromosomes. Residues in black are absent in *Boodlea* 16S rRNA. Blue numbers indicate secondary structure helices in the 16S rRNA model. (B) Comparison between *Boodlea* and *E. coli* 16S rRNA annotated functional regions. Quality of the alignment was assessed based on the predicted posterior probability (in percentage) of each aligned region: very low < 25%; low between 25-50%; high between 50-95%; and perfect > 95%. See also Table S1.

STAR Methods

CONTACT FOR REAGENT AND RESOURCE SHARING

Further information and requests for resources and reagents should be directed to and will be fulfilled by the Lead Contact, Olivier De Clerck (Olivier.declerck@ugent.be).

EXPERIMENTAL MODEL AND SUBJECT DETAILS

Clonal cultures of *Boodlea composita* FL1110, *Chaetomorpha aerea* UTEX799, *Cladophora albida* Calb2, *Cladophora socialis* Csoc2, *Cladophora vadorum* Cvad2, *Dictyosphaeria cavernosa* FL1134, *Pithophora* sp. UTEX787, *Siphonocladus tropicus* Siph3, *Struvea elegans* Sele1, *Valonia utricularis* Vutric3 and *Valonia ventricosa* UTEX 2260 are maintained in the algal culture collection of the Phycology Research Group, Ghent University. The specimens were grown in enriched sterilized natural seawater at 22°C under 12:12 (light:dark) cool white fluorescent light at 60 $\mu\text{mol photons m}^{-2} \text{ s}^{-1}$. To prepare the enriched natural seawater, 20 mL of enriched solution is added to 980 mL of filtered and sterilized natural seawater. The enriched solution consists of: Tris base 5.0 g/L; NaNO₃ 3.5 g/L; Na₂ β -glycerophosphate · H₂O; Na₂EDTA · 2 H₂O 0.529 g/L; Fe(NH₄)₂(SO₄)₂ · 6 H₂O 0.176 g/L; FeCl₃ · 6 H₂O 12.1 mg/L; H₃BO₃ 0.286 g/L; MnSO₄ · 4 H₂O 40.6 mg/L; ZnSO₄ · 7 H₂O 5.5 mg/L; CoSO₄ · 7 H₂O 1.2 mg/L; Thiamine–HCl 0.5 mg/L; Biotin 5.0 mg/L; Cyanocobalamin 10.0 mg/L [58].

METHOD DETAILS

Genomic DNA sequencing. Total genomic DNA from fresh *Boodlea* cultures was isolated by using a modified CTAB extraction protocol [59]. Briefly, 100 mg of fresh algal material was blotted dry on paper, placed inside a 1.5 ml test tube and immediately frozen in liquid nitrogen. Samples were ground with a pestle that fits the 1.5 mL tubes and resuspended in 500 μL of CTAB isolation buffer (2% w/v cetyltrimethylammonium bromide, 1.4 M NaCl, 100 mM Tris-HCl pH 8.0, 20 mM EDTA pH 8.0, 1% w/v polyvinylpyrrolidone) with 5 μL of Proteinase K (QUIAGEN, Germany). The samples were then incubated at 60°C for 40 min. After 30 min, 5

μL of RNase A (QUIAGEN) was added to each sample. Cellular debris were spun down and the aqueous layer was extracted first with phenol:chloroform:isoamyl alcohol (25:24:1 v/v) and then with chloroform:isoamyl alcohol (24:1 v/v). Genomic DNA was precipitated with the addition of two volumes of ice-cold absolute ethanol and 0.3 M of sodium acetate pH 5.5 to each sample and overnight incubation at -20°C . The genomic DNA was washed with ice-cold 70% ethanol, air-dried and dissolved in 50 μL TE buffer (10 mM Tris-HCl pH 8.0, 1 mM $\text{Na}_2\text{-EDTA}$). HMW and LMW DNA bands were size-selected using a BluePippin™ system (Sage Science, USA). The HMW DNA band was isolated with a cut-off range of 10 kb to 50 kb, while the LMW DNA band was isolated with a cut-off range of 1.5 kb to 2.5 kb. The quantity, quality and integrity of the extracted DNA were assessed with Qubit (ThermoFisher Scientific, USA), Nanodrop spectrophotometer (ThermoFisher Scientific), and Bioanalyzer 2100 (Agilent Technologies, USA).

Intact chloroplasts were isolated from living *Boodlea* cells following the protocol of Palmer et al. [60]. In short, 200 g of *Boodlea* filaments were placed in 400 ml of ice-cold isolation buffer (0.35 M sorbitol, 50 mM Tris-HCl pH 8.0, 5 mM EDTA, 0.1 % BSA, 1.5 mM β -mercaptoethanol), homogenized in a blender at 4°C , and filtered through miracloth (Calbiochem). The filtrate was centrifuged at 1000 g for 15 min at 4°C , the supernatant was poured off, and the pellet resuspended in 8 ml of ice-cold wash buffer (0.35 M sorbitol, 50 mM Tris-HCl pH 8.0, 25 mM EDTA). The resuspended pellet was loaded on a step gradient consisting of 18 ml of 52% w/v sucrose, over-layered with 7 ml of 30% w/v sucrose, and centrifuged at 25,000 rpm for 40 min at 4°C . The chloroplast band was removed from the 30%-52% interface using a Pasteur pipette, diluted with 6 volumes of wash buffer, centrifuged at 1,500 g for 15 min at 4°C , and resuspended in wash buffer to a final volume of 10 ml. This fraction of isolated chloroplasts is further referred to as “chloroplast-enriched fraction”. DNA from the chloroplast-enriched fraction was sequenced with Roche 454 GS FLX at GATC Biotech, Germany. The HMW and LMW DNA fractions were sequenced on two SMRT cells on a PacBio RS II (VIB Nucleomics Core facilities, Leuven, Belgium) using PacBio P5 polymerase and C3 chemistry combination (P5-C3). For the HMW DNA fraction, a 20-kb SMRT-bell library was constructed, while for the LMW DNA fraction, a 2-kb SMRT-bell library was constructed.

Chloroplast DNA assembly and annotation. Quality of the reads from the 454 library was assessed with FastQC v.0.10.1 (<http://www.bioinformatics.babraham.ac.uk>, last accessed March 01, 2017) (Table S2). Low-quality reads (average Phred quality score below 20) were discarded and low-quality 3' ends of the reads were trimmed with Fastx v.0.0.13 (https://github.com/agordon/fastx_toolkit, last accessed March 01, 2017). After trimming, reads shorter than 50 bp were discarded. *De novo* assembly of the trimmed reads was performed with MIRA v. 4.0rc5 [61]. The assembly resulted in 3,735 contigs, which will be further referred to as “454 contigs” (Table S3). Length distribution of the 454 contigs is reported in Figure S6B.

After the assembly, putative chloroplast contigs were identified by comparing their translated sequences against the NCBI non-redundant protein database using BLAST 2.2.29+ [62], resulting in the identification of 58 contigs harbouring fragments or full-length chloroplast genes by sequence similarity search. These contigs had long stretches of conserved repetitive sequences at their 5' and 3' extremities. The conserved inverted repeats were used in a sequence similarity search with high stringency (high mismatch cost, high cost for gap opening and gap extension, long minimal word-size) against the 454 contigs to identify 18 additional contigs of putative chloroplast origin. This initial set of 76 contigs had a mean coverage of 84×, ranging between 11× and 191×. 17 of the 76 contigs had internal inverted repeats, with a sudden drop in read coverage. These contigs were regarded as chimeric contigs and were cleaved at the sites of coverage drop, raising the number of contigs of chloroplast origin to 89.

Additional chloroplast contigs without similarity to protein-coding genes were identified by metagenomic binning (distribution analysis of 4-mers) with MyCC [63], resulting in 21 clusters of 454 contigs (Figure S6D). The initial set of 89 chloroplast 454 contigs was present in three neighbouring clusters: Cluster 14, Cluster 17 and Cluster 21. These clusters contained 122, eight and six contigs, respectively, raising the number of identified chloroplast contigs assembled from the chloroplast-enriched fraction from 89 to 136 (“chloroplast 454 contigs” in Table S3). Of these, 71 contigs had no sequence similarity to known protein-coding genes, 36 contigs harboured full-length chloroplast genes, 36 contigs harboured fragments of chloroplast genes, and 7 contigs harboured both fragments and full-length CDSs of different chloroplast genes.

Contigs potentially coding for chloroplast tRNAs and rRNAs were identified using Infernal 1.1 [64]. The chloroplast 454 contigs served as seeds for iterative contig extension with PRICE

1.0.1 [65]. Single-end 454 reads were used as false paired-end reads with expected insert size equal to the median length of the 454 reads. 141 different combinations of parameters were tested in order to optimize the contig extension. None of the selected assemblies showed a length improvement for the initial set of chloroplast 454 contigs. The length distribution of the chloroplast 454 contigs was consistent with the size of the LMW DNA fraction as estimated by agarose gel electrophoresis of *Boodlea* genomic DNA (Figure 1D, Figure S6B).

Repetitive regions in the contigs were identified with ‘inverted’, ‘etandem’ and ‘palindrome’ from the EMBOSS 6.5.7 [66] package. Dotplots for all contigs were generated with YASS v. 1.14, using standard parameters [67]. Coverage of the chloroplast 454 contigs was evaluated by mapping the 454 reads, the mRNA and the total-RNA libraries to these contigs with CLC Genomics Workbench 7.0 (Qiagen) (Figure S6C).

The chloroplast 454 contigs were used together with HMW DNA reads for two independent hybrid assemblies. First, we tried to close hypothetical gaps between the chloroplast 454 contigs with the pbahaScaffolder.py script integrated in the smrtanalysis 2.3.0 pipeline [68]. Secondly, the pre-assembled chloroplast 454 contigs were used as anchors for HMW DNA reads in a round of hybrid assembly with dbg2olc [69]. These analyses failed to close the hypothetical gaps between the short chloroplast 454 contigs and did not yield longer contigs. These results stand in stark contrast to the mitochondrial 454 contigs, where the same approaches yielded markedly longer contigs (Figures S2F and S2G).

Since the hybrid assemblies with uncorrected reads could not reconstruct a circular chloroplast genome, HMW DNA reads were further characterized after error correction. The high-noise HMW DNA reads were corrected by applying a hybrid correction with proovread 2.12 [70] using 454 reads and reads from Illumina RNA-seq libraries (see below). Corrected reads encoding chloroplast genes were identified by aligning them against a custom protein database, named Chloroprotein_db, including genes from the pico-PLAZA protein database [71] and protein-coding genes from published green algal chloroplast genomes (Chlorophyta sensu Bremer 1985, NCBI Taxonomy id: 3041).

LMW DNA reads were self-corrected with the PBcR pipeline [72]. Since the LMW DNA size is unknown and PBcR requires an estimate of the genome size for proper read correction, six different putative genome sizes were tested (100 kb, 1 Mb, 2.24 Mb, 10 Mb, 100 Mb). The best

performance in terms of number of corrected reads was obtained by the combination of "10 Mb" for the estimated genome size and the *-sensitive* flag turned on; these corrected reads were used for the downstream analysis. After error correction, the number of reads was reduced from 154,852 to 106,428 (Table S2), with a similar length distribution as the uncorrected reads library (Figure S6A).

In order to estimate the *Boodlea* chloroplast genome size, k-mer frequency distributions were calculated with jellyfish 2.0 [73]. K-mers ranging from 11 to 47 were analysed for uncorrected and corrected LMW DNA reads, for the filtered 454 reads and for the 454 reads that could be mapped on the chloroplast 454 contigs (Figure S3).

De novo genome assembly of corrected LMW DNA reads was performed with the Celera WGS assembler version 8.3rc2 [74]. The resulting assembly, hereafter called the Celera Assembly, consisted of 558 contigs (Table S3). Corrected and uncorrected reads as well as assembled contigs potentially encoding chloroplast genes were identified by aligning them with BLAST 2.2.29+ against Chloroprotein_db. In order to identify additional short protein-coding genes, HMM profiles were generated from alignments of chloroplast genes present in Chloroprotein_db and used to search the 6-frame translations of 454 contigs and corrected and uncorrected HMW and LMW DNA reads with HMMer3 [75]. To prevent assembly artefacts caused by repetitive elements and palindromic sequences, we also performed an orthology-guided assembly, in which the LMW DNA reads harbouring chloroplast CDSs were re-assembled together with the respective chloroplast 454 contigs. First, LMW DNA corrected reads and chloroplast 454 contigs were grouped according to their best BLAST hit. The corrected reads and contigs belonging to the same group were assembled using Geneious v. 8.1.7 (Biomatters, <http://www.geneious.com/>, last accessed March 01, 2017) with parameters "High Sensitivity/Medium", and each assembly (or lack of assembly) was visually screened to exclude potential chimeric contigs (e.g. palindromic corrected subreads should be collapsed in the same locus rather than being concatenated). Where possible, LMW DNA reads and chloroplast 454 contigs were assembled as larger molecules (Figure S7). The orthology-guided assembly yielded 21 contigs, 2 belonging to group A, 15 to group B and 4 to group E (Table S1). Two groups of reads could not be assembled into longer molecules, and for them, the corresponding chloroplast 454 contigs were retained. Eleven additional chloroplast 454 contigs were retained (Group E), since they were not congruent with the LMW DNA reads and could not be included in the

assembly. This resulted in a total of 32 contigs containing chloroplast protein-coding genes, which together with the two later identified Group B contigs encoding the 16S rRNA gene, are regarded as the *Boodlea* chloroplast genome contigs (Figure 2B, Table S3).

Protein-coding genes in the *Boodlea* chloroplast genome contigs were identified with a sequence similarity search against the NCBI non-redundant protein database with BLAST 2.2.29+. Their annotation was manually refined in Geneious and Artemis 16.0.0 [76] based on the BLAST search results. rRNAs were identified using Infernal 1.1 [64]. Repetitive elements were mapped on the *Boodlea* chloroplast genome by aligning the contigs with themselves using BLAST 2.2.29+. Non-coding RNAs were identified with infernal 1.1 (cut-off value 10^{-5}). Conserved motifs were predicted with MEME suite [77], and the discovered motifs were clustered with RSAT [78]. The motifs were compared with the JASPAR-2016 [79] database using TOMTOM [80] (p-value cut-off 10^{-3}).

Boodlea chloroplast genome coverage was evaluated by mapping the 454 reads with gsnap v.2016-04-04 [81]. Corrected and uncorrected LMW DNA subreads and chloroplast 454 contigs resulting from the MIRA assembly were mapped against the *Boodlea* chloroplast genome with gmap v. 2014-12-06 [81] using the *-nospllicing* flag. Due to the high number of repetitive sequences in LMW DNA reads and 454 contigs, the resulting annotated *Boodlea* chloroplast genome was carefully inspected in order to exclude sequencing and assembly artefacts.

Completeness of the chloroplast genome was evaluated by comparing the annotated chloroplast genes to a set of 60 "core" chloroplast protein-coding genes, defined as protein-coding genes conserved among the chloroplast genomes of the following representative species of Chlorophyta: *Bryopsis plumosa*, *Chlamydomonas reinhardtii*, *Chlorella vulgaris*, *Coccomyxa subellipsoidea*, *Gonium pectorale*, *Leptosira terrestris*, *Nephroselmis olivacea*, *Oltmannsiellopsis viridis*, *Parachlorella kessleri*, and *Pseudendoclonium akinetum*, and the streptophyte *Mesostigma viride* (Figure 3A).

RNA sequencing. Total RNA was isolated using a modified CTAB extraction protocol [82]. RNA quality and quantity were assessed with Qubit and Nanodrop spectrophotomete, and RNA integrity was assessed with a Bioanalyzer 2100. Two cDNA libraries for NextSeq sequencing were generated using TruSeq™ Stranded RNA sample preparation kit (Illumina, USA): one library enriched in poly(A) mRNA due to oligo-(dT) retrotranscription and one total RNA library

depleted in rRNAs with Ribo-Zero Plant kit (Epicentre, USA). The two libraries were sequenced on one lane of Illumina NextSeq 500 Medium platform at 2x76 bp by VIB Nucleomics Core facilities (Leuven, Belgium) (Table S2).

Transcriptome assembly and annotation. Quality of the reads from the two RNA-seq libraries was assessed with FastQC. Low-quality reads (average Phred quality score below 20) were discarded and low-quality 3' ends of the reads were trimmed with Fastx. After trimming, reads shorter than 30 bp were discarded. Read normalization and *de novo* assembly of the libraries were performed with Trinity 2.0.4 [83]. The resulting contigs (hereafter, transcripts) were compared using sequence similarity searches against the NCBI non-redundant protein database using Tera-BLAST™ DeCypher (Active Motif, USA). Taxonomic profiling of the transcripts was performed using the following protocol: for each transcript, sequence similarity searches were combined with the NCBI Taxonomy information of the top ten BLAST Hits in order to discriminate between eukaryotic and bacterial transcripts or transcripts lacking similarity to known protein-coding genes (Table S3). Transcripts classified as "eukaryotic" were further examined to assess transcriptome completeness and to identify chloroplast transcripts. These transcripts were analysed using Tera-BLAST™ DeCypher against Chloroprotein_db. Transcriptome completeness was evaluated with a custom Perl script that compared gene families identified in the *Boodlea* transcriptome to a set of 1816 "core" gene families shared between Chlorophyta genomes present in pico-PLAZA 2.0 [71], following Veeckman et al. guidelines to estimate the completeness of the annotated gene space [84] (mRNA 1,741; total-RNA 1,724 out of 1,816 core gene families identified respectively).

Boodlea chloroplast genome expression and presence of potential RNA editing were evaluated by mapping the reads from the mRNA and total-RNA libraries to the chloroplast genome contigs with gsnap, and by aligning the transcripts resulting from the *de novo* assembly of the RNA-seq libraries to the chloroplast genome contigs with gmap.

Cladophorales genomic DNA sequencing. Sequence data were obtained from 9 additional Cladophorales species, representing the main lineages of the order (Tables S2 and S4). Total genomic DNA was extracted using a modified CTAB extraction protocol as described above, and sequenced using Illumina HiSeq 2000 technology (2×100 bp paired-end reads) on 1/5th of a lane by Cold Spring Harbor Laboratory (Cold Spring Harbor, NY, USA). Quality of the reads

from the sequenced libraries was assessed with FastQC 0.10.1. Low-quality reads (average Phred quality score below 20) were discarded and low-quality 3' ends of the reads were trimmed with Fastx 0.0.13 toolkit. After trimming, reads shorter than 50 bp were discarded. Trimmed reads were assembled with CLC Genomics Workbench, MIRA and SPAdes 3.6.2 [85].

The taxonomic profiling of the contigs was performed with the following protocol: for each contig, sequence similarity searches were combined with the NCBI Taxonomy ID's of the top ten BLAST hits in order to discriminate between eukaryotic and bacterial contigs and contigs with no similarity to known proteins ("NoHit"). Contigs classified as eukaryotic were further analysed to identify chloroplast contigs with a sequence similarity search using Tera-BLAST™ (DeCypher, www.timelogic.com) against Chloroprotein_db. After chloroplast contig identification, the assembly that allowed the reconstruction of the highest number of full-length chloroplast genes was retained. An overview of the assembly metrics is reported in Table S3.

Phylogenetic analysis. Phylogenetic analysis was based on a concatenated alignment of 19 chloroplast protein-coding genes (*atpA*, *atpB*, *atpH*, *atpI*, *petA*, *petB*, *petD*, *psaA*, *psaB*, *psaC*, *psbA*, *psbB*, *psbC*, *psbD*, *psbE*, *psbF*, *psbJ*, *psbL*, and *rbcL*) from *Boodlea*, nine other Cladophorales species, 41 additional species of Archaeplastida, and 14 Cyanobacteria species (Table S6). For each gene, DNA sequences were translated to amino acid sequences and aligned using ClustalW in Geneious using the BLOSUM weight matrix, with gap open penalty 10 and gap extension penalty 0.1. The 19 alignments were concatenated and poorly aligned positions were removed using Gblocks server [86], using the least stringent settings, resulting in an amino acid alignment of 5,704 positions. A maximum likelihood (ML) phylogenetic tree was inferred from the amino acid alignment using RAxML with the cpREV + Γ model [87]. Branch support was assessed by bootstrapping with 500 replicates. Phylogenetic analysis was run on the CIPRES Science Gateway v3.3 [88].

DATA AND SOFTWARE AVAILABILITY

DNA and RNA sequence data have been deposited to the NCBI Sequence Read Archive as BioProject PRJNA384503. The annotated chloroplast and mitochondrial contigs of *Boodlea composita* were deposited to GenBank under accession numbers MG257795 – MG257880. Chloroplast genes from additional Cladophorales species were made available on Mendeley Data

(<http://dx.doi.org/10.17632/7dyphg7pbk.1>). Phylogenetic data (sequence alignments, analyses and phylogenetic tree) were deposited in TreeBase under accession number 21737 (<http://purl.org/phylo/treebase/phyloids/study/TB2:S21737>).