

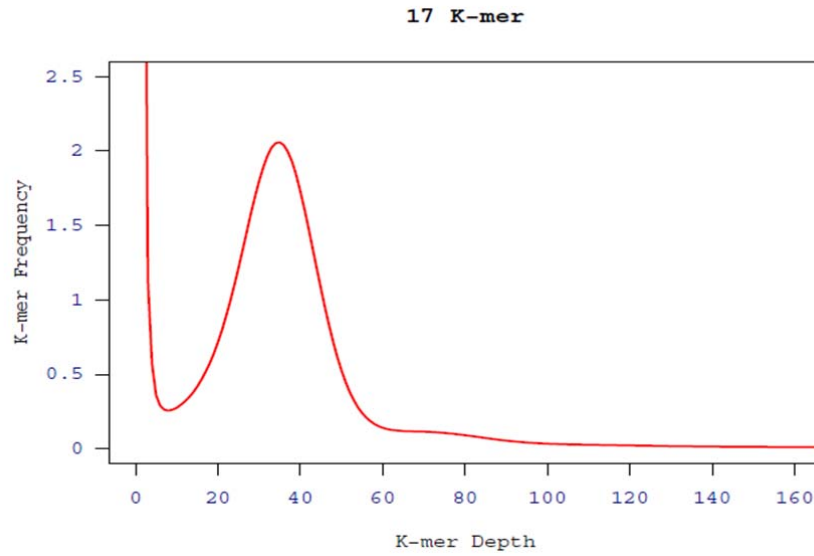
In the format provided by the authors and unedited.

Genomes of early-diverging streptophyte algae shed light on plant terrestrialization

Sibo Wang^{1,2,3,12}, Linzhou Li^{1,4,5,12}, Haoyuan Li^{1,2}, Sunil Kumar Sahu^{1,4}, Hongli Wang^{1,2}, Yan Xu^{1,6}, Wenfei Xian^{1,2}, Bo Song^{1,2}, Hongping Liang^{1,6}, Shifeng Cheng^{1,2}, Yue Chang^{1,2}, Yue Song^{1,2}, Zehra Çebi⁷, Sebastian Wittek⁷, Tanja Reder⁷, Morten Peterson³, Huanming Yang^{1,2}, Jian Wang^{1,2}, Barbara Melkonian^{7,11}, Yves Van de Peer^{8,9}, Xun Xu^{1,2}, Gane Ka-Shu Wong^{1,10*}, Michael Melkonian^{7,11*}, Huan Liu^{1,3,4*} and Xin Liu^{1,2,4*}

¹BGI-Shenzhen, Shenzhen, China. ²China National GeneBank, BGI-Shenzhen, Shenzhen, China. ³Department of Biology, University of Copenhagen, Copenhagen, Denmark. ⁴State Key Laboratory of Agricultural Genomics, BGI-Shenzhen, Shenzhen, China. ⁵Department of Biotechnology and Biomedicine, Technical University of Denmark, Lyngby, Denmark. ⁶BGI Education Center, University of Chinese Academy of Sciences, Shenzhen, China. ⁷Botanical Institute, Cologne Biocenter, University of Cologne, Cologne, Germany. ⁸Department of Plant Biotechnology and Bioinformatics, Ghent University and VIB/UGent Center for Plant Systems Biology, Ghent, Belgium. ⁹Centre for Microbial Ecology and Genomics, Department of Biochemistry, Genetics and Microbiology, University of Pretoria, Pretoria, South Africa. ¹⁰Department of Biological Sciences and Department of Medicine, University of Alberta, Edmonton, Alberta, Canada. ¹¹Present address: University of Duisburg-Essen, Campus Essen, Faculty of Biology, Essen, Germany. ¹²These authors contributed equally: Sibow Wang, Linzhou Li. *e-mail: gane@ualberta.ca; michael.melkonian@uni-koeln.de; liuhuan@genomics.cn; liuxin@genomics.cn

1



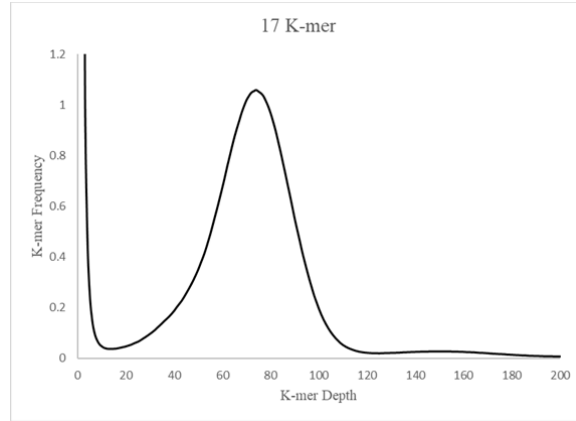
2

3 **Supplementary Figure 1. 17-mer analysis to estimate genome complexity**
4 **(genome size, repeat content estimation, heterozygosity calculation) of *M. viride*.**

5 The 17-mer frequency distribution generated from the sequencing reads was plotted.

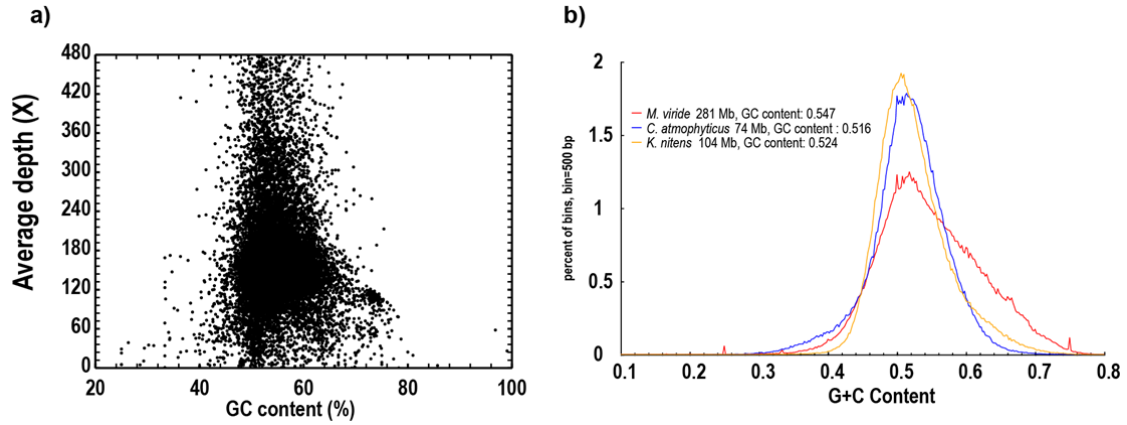
6 The peak is approximately 35 and the total K-mer count is 11,517,226,902. The
7 genome size is estimated as about 329 Mb.

8



9

10 **Supplementary Figure 2. 17-mer analysis to estimate genome complexity**
11 **(genome size, repeat content estimation, heterozygosity calculation) of *C.***
12 ***atmophyticus*.** The 17-mer frequency distribution generated from the sequencing
13 reads was plotted. The peak depth is about 74 and the total K-mer count is
14 6,340,143,223. The genome size was estimated as about 85.68 Mb., low
15 heterozygosity rate and repeat content can be observed.



17

18

19 **Supplementary Figure 3. Assessment of genomic data for assembly of *M. viride*.** (a) The

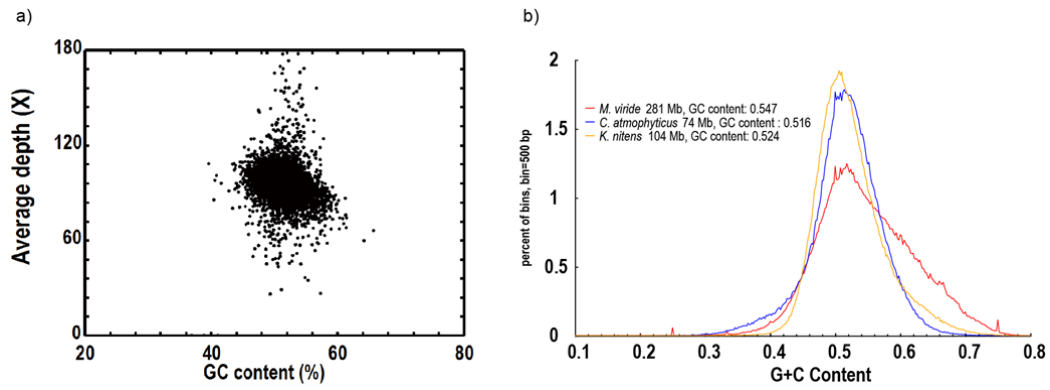
20 GC content and the average depth were calculated from 10 kb non-overlapping sliding

21 windows. The distribution pattern of the GC content indicates a relative pure single genomic

22 sample without contaminations; (b) comparison of GC content across closely related species

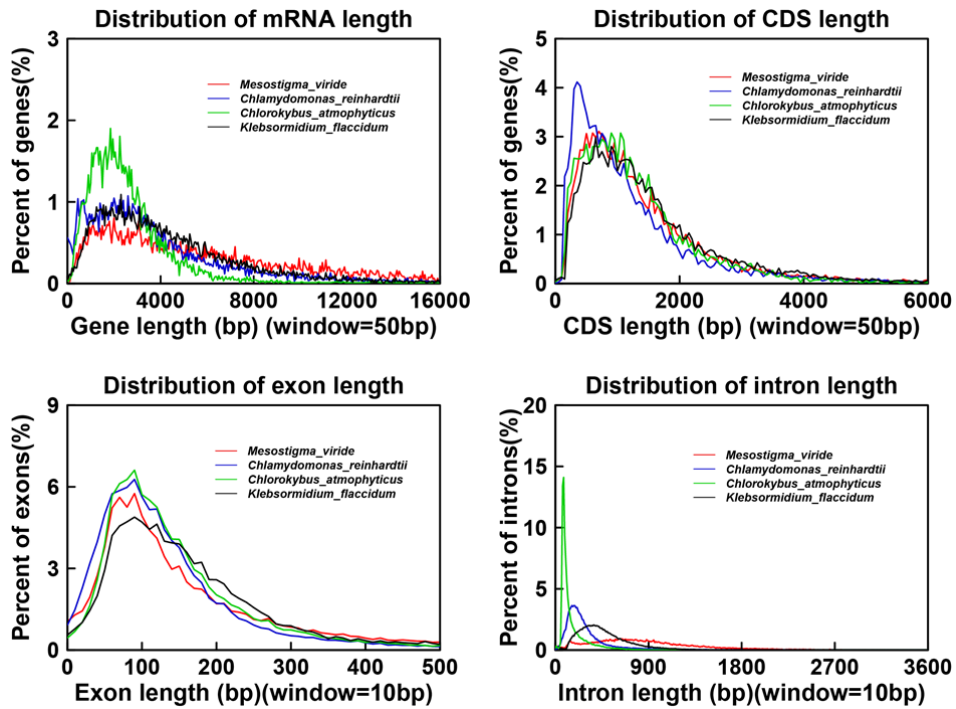
23 and model species.

23



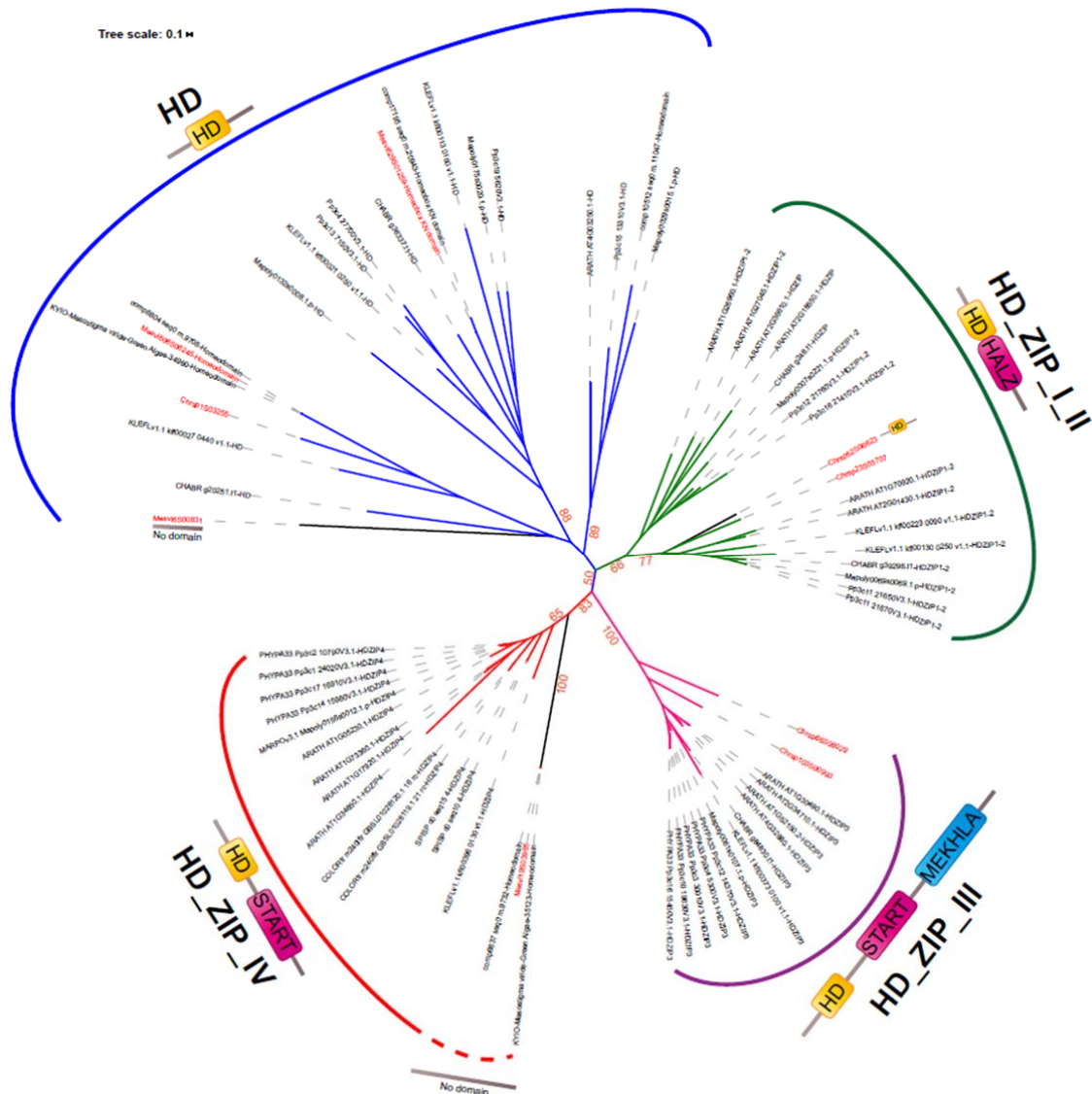
24
25
26
27
28
29

Supplementary Figure 4. Distribution of GC depth of *C. atmophyticus*. (a) The GC content and the average depth were calculated from 10 kb non-overlapping sliding windows. The distribution pattern of GC content indicates a relative pure single genomic sample without contamination and no GC bias; (b) comparison of GC content across closely related species;



30
31
32
33
34

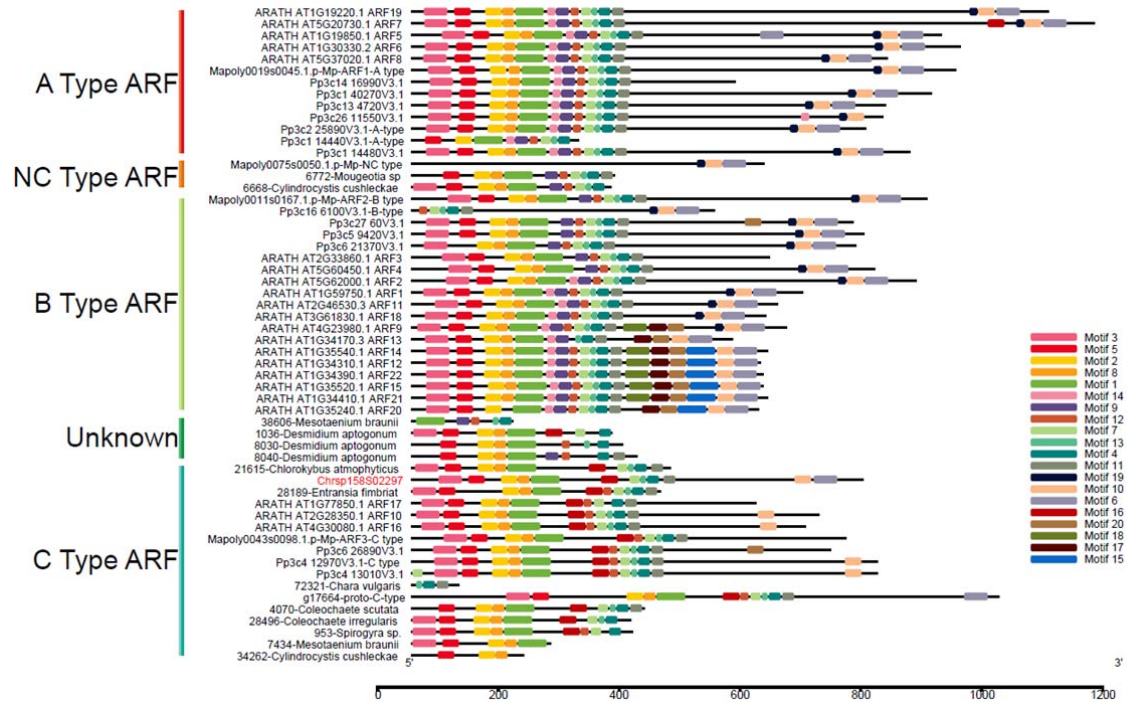
Supplementary Figure 5. Comparison of gene structural features among four green algal species. The distributions of mRNA length, CDS length, exon length and intron length for *M. viride* and *C. atmophyticus* genome were compared against other relative species.



35
 36 **Supplementary Figure 6. Phylogenetic tree of the transcription factor HD and HD-Zip.**
 37 The tree derived from a MAFFT alignment and constructed using IQ-TREE (see Methods).
 38 Bootstrap values (200 replicates) $\geq 50\%$ are shown. The sequences derived from the *M. viride*
 39 and *C. atrophyticus* genomes are highlighted in red. The tree also included the sequences
 40 derived from the transcriptomes of *M. viride* and *C. atrophyticus*. Each clade was marked
 41 with the domain structure information and HD-ZIP subfamily name. Although one of *M.*
 42 *viride* genome sequence and its corresponding transcriptome sequence clustered with HD-ZIP
 43 IV, these sequences display no domain structure (marked as dashed red line).
 44

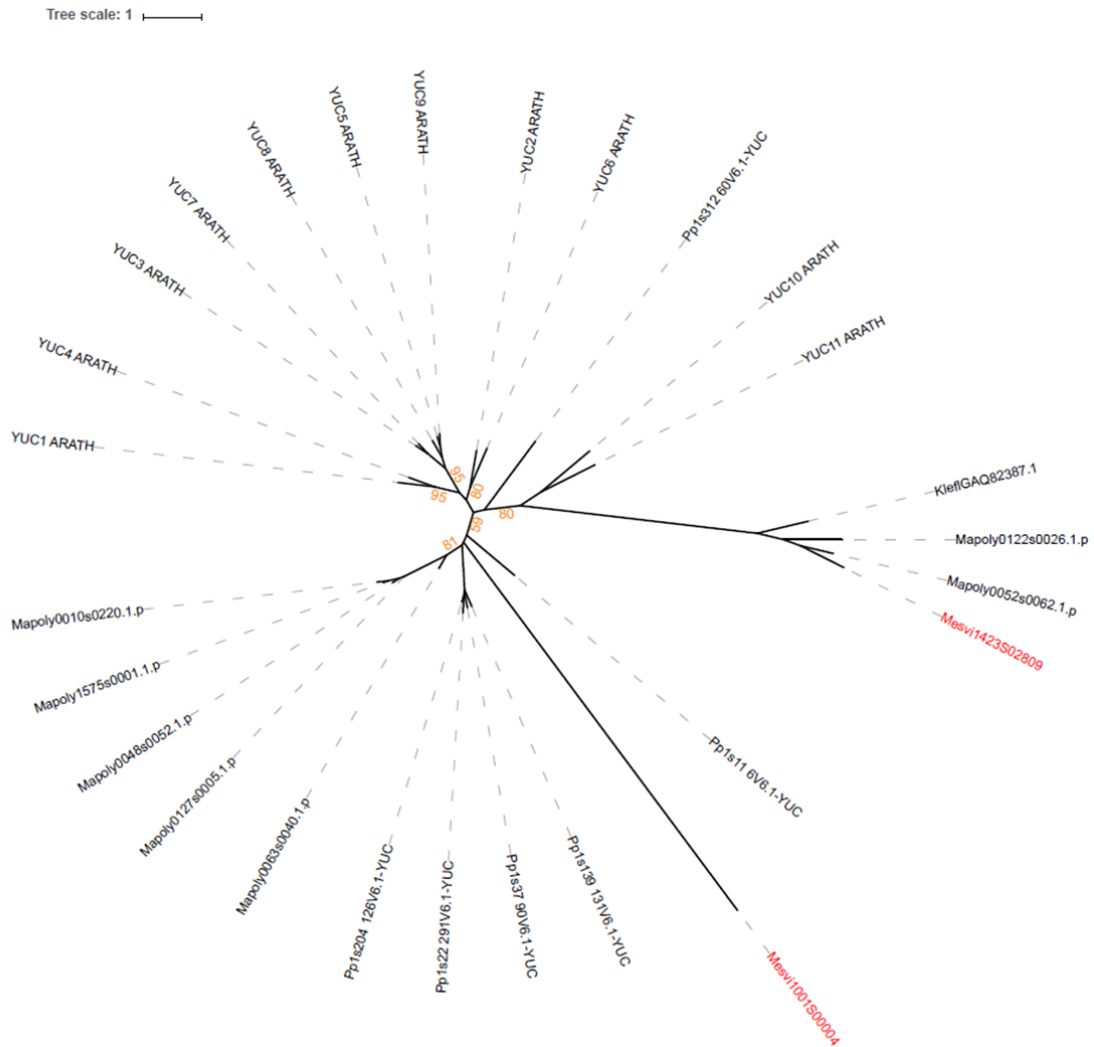


45 **Supplementary Figure 7. Phylogenetic tree of the Auxin response factor (ARF).** (a) The
 46 tree derived from a MAFFT alignment and constructed using IQ-TREE (see Methods).
 47 Bootstrap values (200 replicates) $\geq 50\%$ are shown. The sequence derived from the *C.*
 48 *atmophyticus* genome is highlighted in red. The tree also included the sequence derived from
 49 transcriptome of *C. atmophyticus* (21615-Chlorokybus_atmophyticus). (b). Conserved motifs
 50 of each sequence were identified in the respective clade through MEME analysis. 20 motifs
 51 were identified.
 52
 53



54
 55
 56
 57

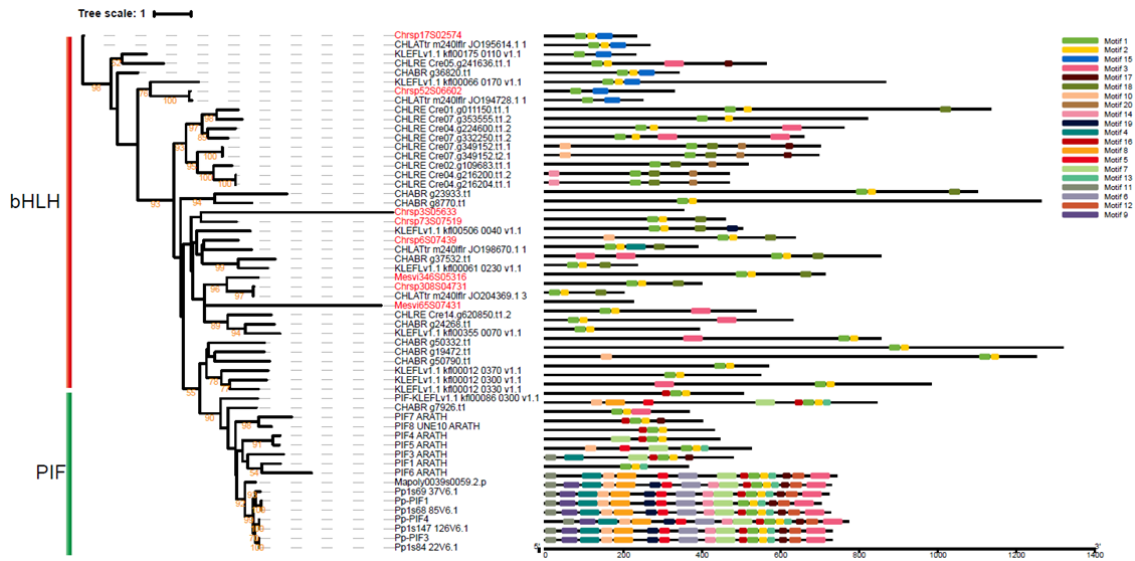
Supplementary Figure 8. Conserved motifs of each sequence were identified in the respective clade through MEME analysis of the Auxin response factor (ARF). 20 motifs were identified.



59

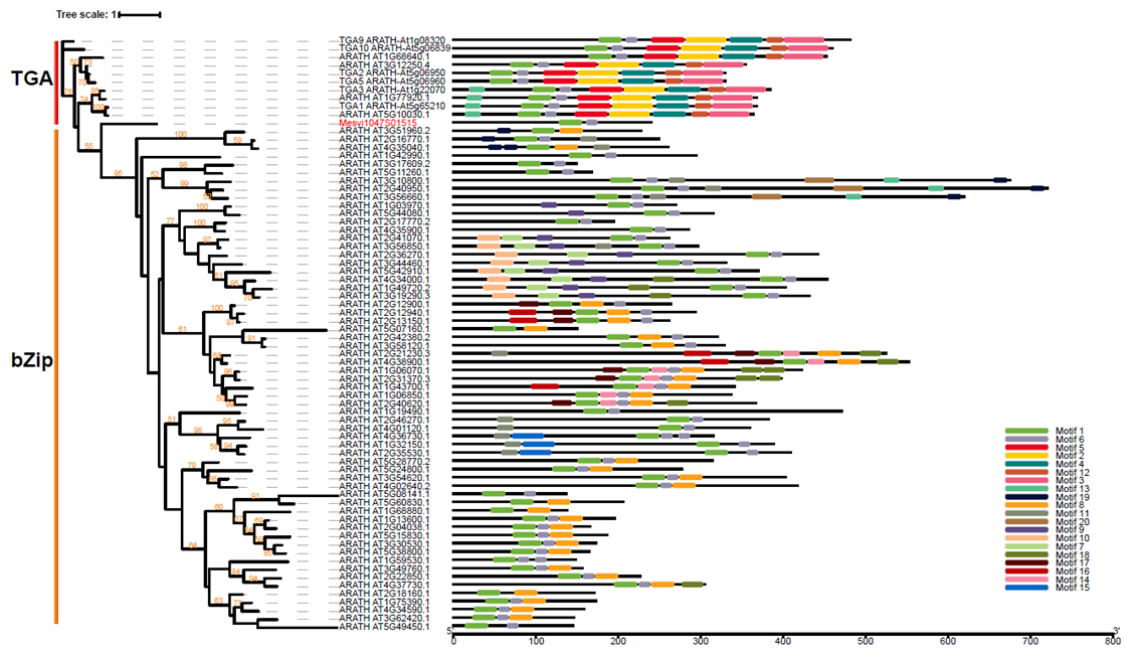
60 **Supplementary Figure 9. Phylogenetic tree of YUCCA.** The tree derived from a MAFFT
 61 alignment and constructed using IQ-TREE (see Methods). Bootstrap values (200 replicates) \geq
 62 50% are shown. The sequences derived from the *M. viride* genome are highlighted in red.
 63 The tree includes all identified YUCCA sequences of *Arabidopsis thaliana*, *Marchantia*
 64 *polymorpha*, *Physcomitrella patens* and *Klebsormidium nitens* downloaded from Uniport
 65 database (<https://www.uniprot.org>).

66



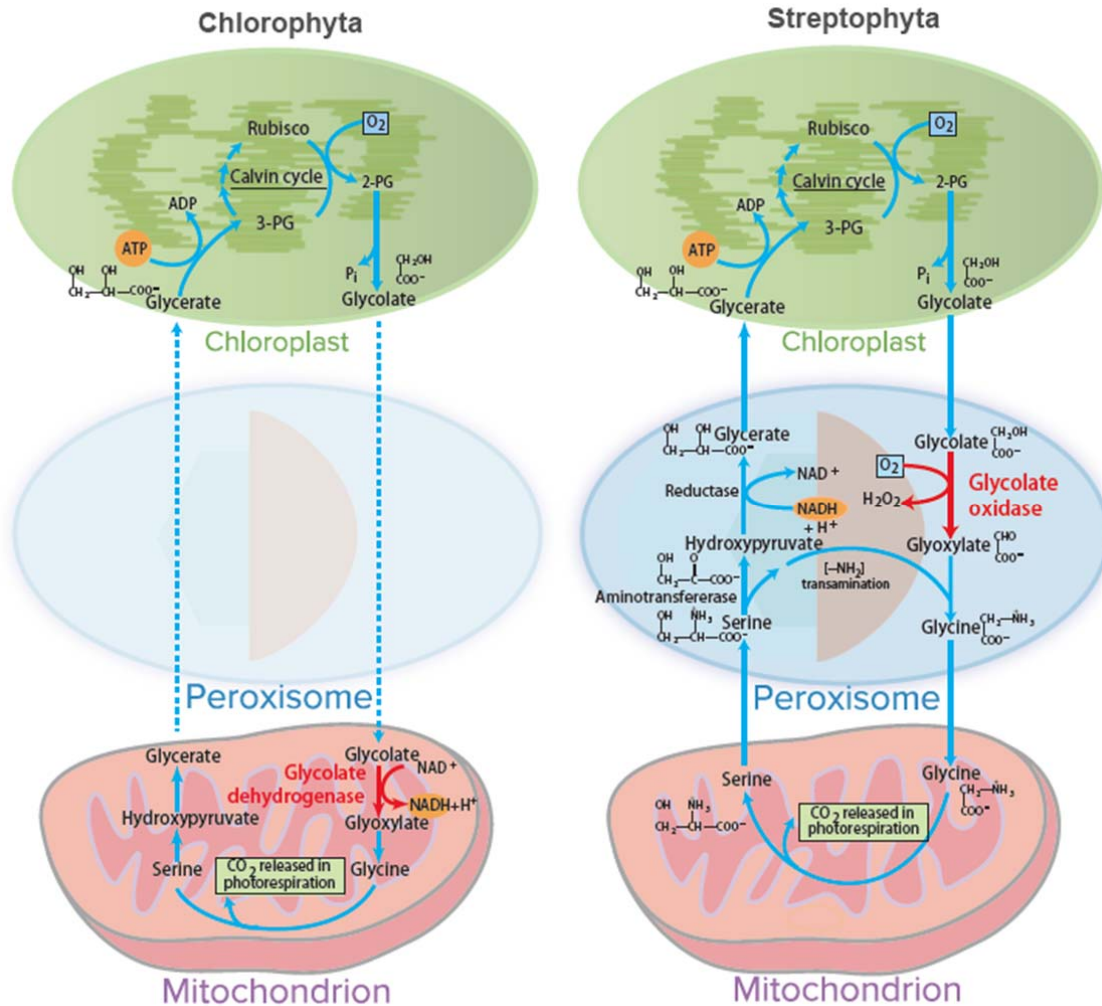
75
76
77
78

Supplementary Figure 11. Conserved motifs of each sequence were identified in the respective clade through MEME analysis of PIF and bHLH. 20 motifs were identified.



86
87
88
89
90

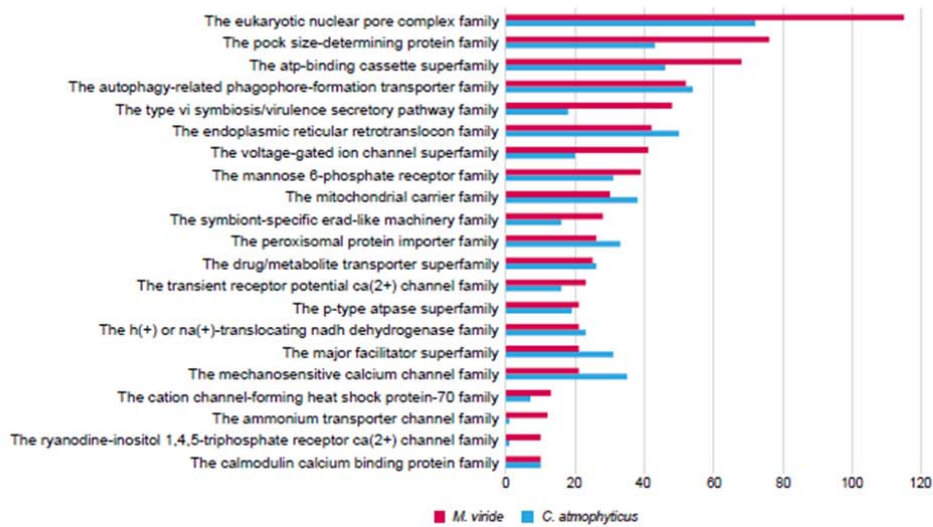
Supplementary Figure 13. Conserved motifs of each sequence were identified in the respective clade through MEME analysis of the TGA and bZIP. 20 motifs were identified.



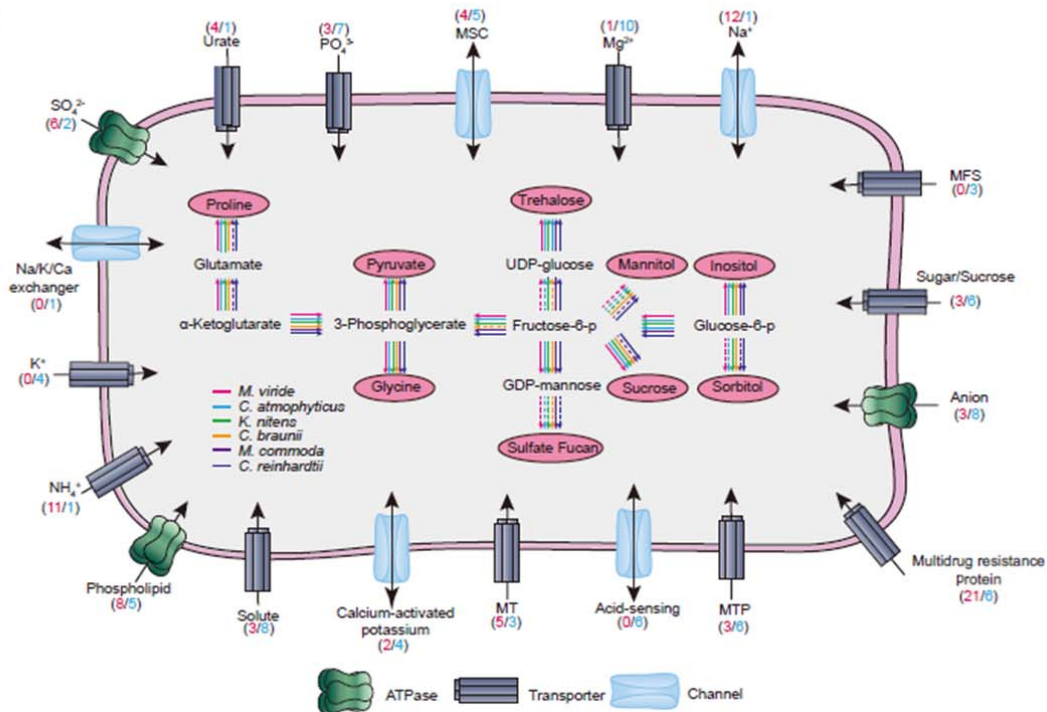
91
 92
 93
 94
 95
 96
 97
 98
 99
 100
 101
 102
 103

Supplementary Figure 14. The photorespiration pathway is different in Chlorophyta compared to Streptophyta. Analysis of targeting signal peptide of glycolate dehydrogenase and glycolate oxidase in Chlorophyta and Streptophyte algae. Glycolate oxidase of streptophyte algae is located in peroxisomes, but no peroxisome targeting signal was identified in the glycolate dehydrogenase of Chlorophyta. All glycolate dehydrogenases derived from Chlorophyta are located in mitochondria. However, glycolate dehydrogenase of streptophyte algae was found to be targeted to various organelles (chloroplast, mitochondria) except the peroxisomes (Table S). The major metabolites produced/involved during photorespiration in chlorophytes and streptophytes are also depicted for streptophytes, and the major differences among them are highlighted in red color.

a



b



104

105

Supplementary Figure 15. Overview of the key components involved in the organic and

106

ionic osmoregulation

107

representative Streptophyta and Chlorophyta are shown. Dash arrows mean incomplete

108

osmoregulation. Only the significantly different transporters, ATPase, and channels are

109

presented in the figure. The compatible osmolyte biosynthetic pathways between

110

representative Streptophyta and Chlorophyta are shown. Dash arrows mean incomplete

111

pathway, while solid arrows mean complete pathway. The copy number shows in pink and

112

blue color for *Mesostigma* and *Chlorokybus*, respectively.