# Cell

# The Chara genome: secondary complexity and implications for plant terrestrialization
## --Manuscript Draft--

| | Clemence Bonnot |
| --- | --- |
| | Holger Breuninger |
| | Aikaterini Symeonidi |
| | Guru V Radhakrishnan |
| | Filip Van Nieuwerburgh |
| | Dieter Deforce |
| | Caren Chang |
| | Kenneth G Karol |
| | Rainer Hedrich |
| | Peter Ulvskov |
| | Gernot Glöckner |
| | Charles F Delwiche |
| | Jan Petrasek |
| | Yves Van de Peer |
| | Jiri Friml |
| | Mary Beilby |
| | Liam Dolan |
| | Yuji Kohara |
| | Sumio Sugano |
| | Asao Fujiyama |
| | Pierre-Marc Delaux |
| | Marcel Quint |
| | Günter Theissen |
| | Martin Hagemann |
| | Jesper Harholt |
| | Christophe Dunand |
| | Sabine Zachgo |
| | Jane Langdale |
| | Florian Maumus |
| | Dominique Van Der Straeten |
| | Sven Gould |
| | Stefan A Rensing, Ph.D. |

**Abstract:** Land plants evolved from charophytic algae, among which Charophyceae possess the most complex body plans. We present the genome of Chara braunii; comparison of the genome to those of land plants identified evolutionary novelties for plant terrestrialization and land plant heritage genes. C. braunii employs unique xylan synthases for cell wall biosynthesis, a phragmoplast (cell separation) mechanism similar to that of land plants, and many phytohormones. C. braunii plastids are controlled via land plant-like retrograde signaling, and transcriptional regulation is more elaborate than in other algae. The morphological complexity of this organism may result from expanded gene families, with three cases of particular note: genes effecting tolerance to reactive oxygen species (ROS), LysM receptor-like kinases, and transcription factors (TFs). Transcriptomic analysis of sexual reproductive structures reveals intricate control by TFs, activity of the ROS gene network, and the ancestral

| | use of plant-like storage and stress protection proteins in the zygote. |
|---|---|
| **Opposed Reviewers:** | John Bowman<br>john.bowman@sci.monash.edu.au<br>Conflict of interest (Marchantia paper) |
| **Suggested Reviewers:** | David Domozych<br>ddomoz@skidmore.edu<br>Expert on charophytes |
| | Linda Graham<br>lkgraham@wisc.edu<br>Expert on charophytes |
| | James Leebens-Mack<br>jleebensmack@uga.edu<br>Expert on plant phylogenomics / plant genome evolution |

# The *Chara* genome: secondary complexity and implications for plant terrestrialization

Tomoaki Nishiyama[1,*], Hidetoshi Sakayama[2,*], Jan de Vries[4,5], Henrik Buschmann[3], Denis Saint-Marcoux[6,7], Kristian K. Ullrich[8,40], Fabian B. Haas[8], Lisa Vanderstraeten[9], Dirk Becker[10], Daniel Lang[38], Stanislav Vosolsobě[17], Stephane Rombauts[11], Per K.I. Wilhelmsson[8], Philipp Janitza[12], Ramona Kern[13], Alexander Heyl[14], Florian Rümpler[15], Luz Irina A. Calderón Villalobos[30], John M. Clay[16], Roman Skokan[17], Atsushi Toyoda[18], Yutaka Suzuki[19], Hiroshi Kagoshima[20], Elio Schijlen[39], Navindra Tajeshwar[14], Bruno Catarino[6], Alexander J Hetherington[6], Assia Saltykova[11,21,22], Clemence Bonnot[6,36], Holger Breuninger[6,23], Aikaterini Symeonidi[8], Guru V. Radhakrishnan[24], Filip Van Nieuwerburgh[37], Dieter Deforce[37], Caren Chang[16], Kenneth G. Karol[25], Rainer Hedrich[10], Peter Ulvskov[26], Gernot Glöckner[27], Charles F. Delwiche[16], Jan Petrášek[17], Yves Van de Peer[11,28], Jiri Friml[29], Mary Beilby[31], Liam Dolan[6], Yuji Kohara[20], Sumio Sugano[19], Asao Fujiyama[18], Pierre-Marc Delaux[32], Marcel Quint[12,30], Günter Theißen[15], Martin Hagemann[13], Jesper Harholt[33], Christophe Dunand[32], Sabine Zachgo[3], Jane Langdale[6], Florian Maumus[34], Dominique Van Der Straeten[9], Sven B. Gould[4], Stefan A. Rensing[8,35,*,+]

[1] Advanced Science Research Center, Kanazawa University, Kanazawa 920-0934, Japan

[2] Graduate School of Science, Kobe University, Kobe 657-8501, Japan

[3] Botany Department, School of Biology and Chemistry, Osnabrück University, 49076 Osnabrück, Germany.

[4] Institute for Molecular Evolution, Heinrich Heine University, 40225 Düsseldorf, Germany

[5] Department of Biochemistry and Molecular Biology, Dalhousie University, Halifax, Nova Scotia B3H 4R2, Canada

[6] Department of Plant Sciences, University of Oxford, Oxford, OX1 3RB, United Kingdom

[7] Université de Lyon, UJM-Saint-Étienne, CNRS, BVpam FRE3727, 42023 Saint-Étienne, France

[8] Plant Cell Biology, Faculty of Biology, University of Marburg, 35043 Marburg, Germany

[9] Laboratory of Functional Plant Biology, Department of Biology, Gent University, 9000 Gent, Belgium

[10] Molecular Plant Physiology & Biophysics, University of Wuerzburg, 97082 Wuerzburg, Germany

[11] Department of Plant Biotechnology and Bioinformatics, Gent University and VIB Center for Plant Systems Biology, 9052 Gent, Belgium

[12] Institute of Agricultural and Nutritional Sciences, Martin-Luther-University Halle-Wittenberg, 06120 Halle (Saale), Germany

[13] Plant Physiology, University Rostock, 18051 Rostock, Germany

[14] Department of Biology, Adelphi University, Garden City, NY 11530, USA

[15] Department of Genetics, Friedrich Schiller University Jena, 07743 Jena, Germany

[16] CBMG, University of Maryland, College Park, MD 20742, USA

[17] Department of Experimental Plant Biology, Faculty of Science, Charles University, 128 44 Prague 2, Czech Republic

[18] Comparative Genomics Laboratory and Advanced Genomics Center, National Institute of Genetics, Shizuoka 411-8540, Japan

[19] Department of Computational Biology and Medical Sciences, University of Tokyo, Kashiwa, Chiba 277-8562, Japan

[20] Genome Biology Laboratory, National Institute of Genetics, Shizuoka 411-8540, Japan

[21] Platform Biotechnology and Molecular Biology, Scientific Institute of Public Health (WIV-ISP), Brussels, Belgium

[22] Department of Information Technology, Gent University, IMinds, 9052 Gent, Belgium

[23] ZMBP, Entwicklungsgenetik, 72076 Tübingen, Germany

[24] Department of Cell and Developmental Biology, John Innes Centre, Norwich NR4 7UH, United Kingdom

[25] Lewis B. and Dorothy Cullman Program for Molecular Systematics, The New York Botanical Garden, Bronx, NY 10458, USA

[26] Department of Plant and Environmental Sciences, University of Copenhagen, DK-1871 Frederiksberg C, Denmark

[27] Biochemistry I, Medical Faculty, University of Cologne, 50931 Cologne, Germany

[28] Department of Biochemistry, Genetics and Microbiology, University of Pretoria, Pretoria, 0028, South Africa

[29] Institute of Science and Technology, 3400 Klosterneuburg, Austria

[30] Department of Molecular Signal Processing, Leibniz Institute of Plant Biochemistry, 06120 Halle (Saale), Germany

[31] School of Physics, University of NSW, Sydney, Kensington, 2052, NSW, Australia

[32] Laboratoire de Recherche en Sciences Végétales, Université de Toulouse, CNRS, UPS, Auzeville, BP42617, 31326 Castanet Tolosan, France

[33] Carlsberg Research Laboratory, 1799 Copenhagen V, Denmark

[34] URGI, INRA, Université Paris-Saclay, 78026 Versailles, France

[35] BIOSS Centre for Biological Signalling Studies, Unigersity Freiburg, Germany

[36] Present address: Labex ARBRE, UMR 1136 INRA-Université de Lorraine (IAM), INRA-Grand Est-Nancy, Champenoux, France

[37] Laboratory of Pharmaceutical Biotechnology, Gent University, 9000, Gent, Belgium

[38] PGSB, Helmholtz Center Munich, 85764 Neuherberg, Germany

[39] Wageningen University, B.U. Bioscience, 6700 AA Wageningen, The Netherlands

[40] Present address: Max Planck Institute for Evolutionary Biology, 24306, Ploen, Germany.


+ Lead contact: Stefan A. Rensing


* Authors for correspondence:

tomoakin@staff.kanazawa-u.ac.jp

hsak@port.kobe-u.ac.jp

stefan.rensing@biologie.uni-marburg.de

1    **Summary**

2    Land plants evolved from charophytic algae, among which Charophyceae possess the most
3    complex body plans. We present the genome of *Chara braunii*; comparison of the genome to
4    those of land plants identified evolutionary novelties for plant terrestrialization and land plant
5    heritage genes. *C. braunii* employs unique xylan synthases for cell wall biosynthesis, a
6    phragmoplast (cell separation) mechanism similar to that of land plants, and many
7    phytohormones. *C. braunii* plastids are controlled *via* land plant-like retrograde signaling, and
8    transcriptional regulation is more elaborate than in other algae. The morphological complexity
9    of this organism may result from expanded gene families, with three cases of particular note:
10   genes effecting tolerance to reactive oxygen species (ROS), LysM receptor-like kinases, and
11   transcription factors (TFs). Transcriptomic analysis of sexual reproductive structures reveals
12   intricate control by TFs, activity of the ROS gene network, and the ancestral use of plant-like
13   storage and stress protection proteins in the zygote.

14

15   **Keywords:** plant evolution, charophyte, phytohormones, transcriptional regulation,
16   Phragmoplastophyta, Chara, streptophyte, reactive oxygen species, phragmoplast

## Introduction

A pivotal event in the emergence of plant life was the mid-Paleozoic adaptation to land. While several algal lineages evolved to occupy terrestrial environments, only one represents the land plant ancestor; its terrestrialization event was fostered by a range of evolutionary novelties. The specific complement of traits that allowed a particular algal lineage to give rise to land plants and dominate the terrestrial environment remains under active study. Similarity of critical plant developmental, sensory, and regulatory pathways to homologous pathways in charophyte green algae has been demonstrated in several recent studies, emphasizing the close relationship among these lineages (reviewed in (Rensing, 2018)).

Charophytic algae are the closest living relatives of land plants (embryophytes), with both groups collectively referred to as streptophytes (Fig. 1). The Charophyceae, Coleochaetophyceae, and Zygnematophyceae together with the land plants represent the clade Phragmoplastophyta (Lecointre and Le Guyader, 2006), united by the presence of the phragmoplast (Pickett-Heaps, 1975), an array of microtubules perpendicular to the cell division plane that functions in the formation of the nascent cell wall. The Klebsormidiophyceae, Chlorokybophyceae, and Mesostigmatophyceae share fewer traits with land plants (Fig. 1). While Charophyceae were hypothesized to be most closely related to land plants on the basis of similar body plans (Pringsheim, 1862), recent studies indicated that the Zygnematophyceae are the land plant sister group (Wickett et al., 2014).

Extant Zygnematophyceae have simple body plans that seem to reflect secondary loss of morphological complexity. In contrast, the earlier diverging Charophyceae are morphologically more complex than all other charophytic algae: the haploid thallus body plan encompasses a shoot-like axis consisting of nodes with whorls, internodes, a simplex apical meristem, plus multicellular rhizoids (Fig. 2). Cells of the internode are large and complex, featuring endo- and ectoplasma, multiple plastids and nuclei, and communicate *via* electrical signals. The morphology of extant charophytic groups thus infers mosaic evolution and suggests that the genomes of Charophyceae, not Zygnematophyceae, will likely reveal the suite of traits that facilitated terrestrialization (Delwiche, 2016).

Here we present the genome sequence of the charophycean alga *Chara braunii*, one of the most morphologically complex extant Charophyta, shedding light on early embryophyte diversification and the colonization of land by plants.

## Results and Discussion

### The *Chara braunii* genome: assembly, annotation and comparison

*C. braunii* features a haplontic life cycle (Fig. 2), the draft sequence reported here represents a haploid genome. 1.75 Gbp of nuclear scaffold data were obtained, of which 1.43 Gbp were assembled into contigs, corresponding to ~74% of the *C. braunii* genome. RNA-seq of vegetative and reproductive stages was used together with full-length cDNA sequences to annotate the genome. 23,546 putative protein coding gene models were identified, of which 53% are supported by RNA-seq data (Table S4). At least 94% of several conserved core gene sets are encoded by the genome, indicating its suitability for genomic and comparative analyses (STAR Methods).

The observed chromosome number n=14 (Fig. S1) corresponds to the base chromosome number of *Chara* species. Indeed, synonymous substitution distance (Ks)-based analysis of *C. braunii* paralogs revealed no evidence of whole genome duplication (WGD) events (Fig. S3) and thus paralog acquisition and retention is probably due to small-scale duplications. Repetitive elements (Table S1F, S1G) collectively contribute approximately 1.1 Gbp (61%, or 75% if gaps are excluded). Unlike in most plants and green algae, there are no Copia-type long terminal repeat (LTR) retrotransposons (RT) detectable. We discovered a family of repeats with putative GIY-YIG homing endonuclease and reverse transcriptase domains, which are hallmarks of Penelope RTs and group II introns that are uncommon in plant genomes.

The density of LTR elements in the genome is intermediate between compact genomes like those of *Arabidopsis thaliana* or *Klebsormidium nitens*, and other large genomes such as maize and barley (Fig. 3). *C. braunii* introns are an order of magnitude longer than in any of the other genomes investigated here (Table S1L), although intron boundaries appear to be conserved. The high intron length coincides with a low number of introns per gene (3.82), similar to the value for the barley genome (3.89, Table S1L); intron length and number show negative correlation (r = -0.42). Repetitive elements represent 39% of the intron space (Fig. 3, Mendeley archive) which is strikingly enriched with Penelope-like elements and depleted of other types of repeats including Class 2 transposable elements (Helitrons and DNA transposons), suggesting differential integration bias and/or retention in introns as compared to intergenic space (Table S1L).

### Evolutionary novelties enabling terrestrialization and land plant heritage genes (LPHG)

The lineage harboring *C. braunii* diverged from land plants 550–750 Ma (Morris et al., 2018). By identifying features that are shared between the *C. braunii* genome and extant land plants, putative ancestral traits can be identified that have been retained over several hundred Ma. Here we refer to the genes underlying these traits as land plant heritage genes (LPHG) and similarly deduce evolutionary novelties.

*Cell division and cell wall*

C. braunii, like land plants, performs cytokinesis by assembling a cell plate using a phragmoplast microtubule array while *K. nitens* divides by an evolutionarily older cleavage (Fig. 1). Phragmoplast-mediated cell division is assumed to have enabled filament branching through a shift in the plane of cell division (Buschmann and Zachgo, 2016). Land plants also evolved another microtubule array, the preprophase band (PPB), which functions in phragmoplast and cell plate guidance. Focusing on genes for phragmoplast and PPB function, a list of 221 *A. thaliana* cytokinesis genes was compiled (Table S1C). Sequence comparisons showed that the genomes of *A. thaliana*, *P. patens*, *C. braunii* and *K. nitens* have a highly similar complement of cytokinesis-related genes while the unicellular chlorophyte *Chlamydomonas reinhardtii* is divergent. Interestingly, the *C. braunii* genome lacks the *TANGLED1* gene. In land plants, microtubule-associated TANGLED1 localizes to PPBs and is required for phragmoplast guidance (Walker et al., 2007). Since TANGLED1 homologs are found in several bryophytes, but none in any algae, this gene likely played an important role in PPB evolution (Fig. 1). To gain further insight into the evolution of the phragmoplast, we determined how many paralogs each of the cytokinesis genes has in *C. braunii* as compared to *K. nitens*. In this way we identified possible phragmoplast signature genes (Table S1C). Among others, we detected expansion of cyclins as well as EXOCYST and SNARE complex members (Table S1C, Data S1Q-S). The expansion of phragmoplast-related gene families in *C. braunii* / the Phragmoplastophyta, but not in Chlorophyta, *K. nitens* or *M. viride*, suggests their sub- and neofunctionalization to enable phragmoplast function.

Like land plant cell walls, those of *C. braunii* consist of cellulose embedded within a pectin and hemicellulose matrix (Sorensen et al., 2011), its synthesis is orchestrated by a repertoire of glycosyltransferases much like in land plants (Table S1H), with the exception of an apparently unique mechanism for xylan synthesis. The GT47 xylan synthase XYS1 has been identified in *K. nitens*, as well as IRX9 and IRX14 from GT43 (Data S1A), implicated in xylan biosynthesis despite no apparent requirement for being an active enzyme (Ren et al., 2014). Orthologs to neither XYS1 nor IRX9/14 could be identified in *C. braunii*, however, a deep branching, highly diverged form of GT43 was identified as the most likely *C. braunii* xylan synthase, providing the first hint that GT43 sequences are enzymatically involved in synthesizing xylan.

*Phytohormones*

Phytohormones enable the integration of environmental stimuli with developmental programs. As such, they are a key feature of land plants, with some apparently having origins in algae (Hori et al., 2014; Ju et al., 2015; Wang et al., 2015). Potential orthologs of phytohormone pathway genes were defined across *K. nitens*, *C. braunii*, *P. patens* and *A. thaliana* (Table 1, Fig. 4, Table S1J).

*Auxin (AUX)*

AUX is one of the major regulators of plant growth and development. Biosynthesis of AUX (Hori et al., 2014) as well as transcriptional and physiological response to high concentrations have been shown in *K. nitens* (Ohtaka et al., 2017). In contrast to *K. nitens,* genes encoding biosynthetic enzymes of the TAA and YUCCA families are absent from *C. braunii* (Table 1). In *C. australis* IAA, serotonin and melatonin accumulate in a synchronized manner during the

day/night cycle (Beilby et al., 2015). As the tryptamine IAA biosynthetic pathway intersects with the serotonin/melatonin pathway (Tivendale et al., 2014), *Chara* may synthesize and metabolize AUX *via* a different route than land plants.

Homologous genes for both *PIN*s and *ABCBs* are present in the *C. braunii* genome (Table 1, Table S1K), supporting previous data on polar AUX transport (PAT) in *K. nitens* (Hori et al., 2014) and Charales (Boot et al., 2012). Homologous sequences for AUX1/LAX-like influx carriers as well as the intracellular PIN-like (PILS) transporters, however, are absent from the *C. braunii* genome (Table 1), suggesting that AUX transport and homeostasis display an evolutionary history different from other streptophytes.

The land plant-type AUX signaling cascade, consisting of SCF$^{TIR1/AFB}$ and Aux/IAA co-receptors and ARF TFs, was suggested to be absent in *K. nitens* (Hori et al., 2014; Ohtaka et al., 2017). *K. nitens* encodes an Aux/IAA domain containing protein (Wang et al., 2015) that features an additional B3 domain, is not induced by IAA (Ohtaka et al., 2017) and thus not classified as canonical Aux/IAA (Table 1). In addition to all components of the ubiquitin-proteasome system (Table S1I), *C. braunii* features a single *ARF* (Data S1E) with land plant-like domain composition (Flores-Sandoval et al., 2018), and two Aux/IAA-like sequences (Table 1, Data S1F) clustering with the *A. thaliana* non-canonical IAA33 (lacking a TOPLESS interacting motif and degron for AUX-dependent SCF$^{TIR1/AFB}$-Aux/IAA interactions.

*C. braunii* also encodes several F-box proteins (FBPs) with sequence similarity to land plant phytohormone co-receptors (Data S1P). None of the TIR1/AFB-like FBPs cluster with the land plant AUX co-receptor gene family (Data S1G). Our structural modeling, however, reveals that the *C. braunii* sequences adopt a solenoid-fold architecture resembling TIR1 (Tan et al., 2007). Ligand binding modeling supported the potential ability to form an AUX binding pocket (Data S1P). The existence of only degron-less *C. braunii* Aux/IAAs, however, prompts to postulate that a land plant-like TIR1/AFB-Aux/IAA co-receptor pair is most likely not functional in *C. braunii*.

Consequently, while obvious candidates for canonical land plant AUX biosynthesis genes are lacking, there is a partial candidate gene set of the major land plant AUX signaling and PAT pathways in *C. braunii*. In conclusion, AUX biosynthesis, transport, and some form of signaling were already present in the last common ancestor of *C. braunii* and *K. nitens*, but AUX signaling *via* ARFs was apparently gained in the common ancestor of Phragmoplastophyta, as was ARF repression by Aux/IAAs (Table S1Q, Fig. 4).


*Cytokinin (CK)*

The CK signaling pathway consists of four protein families: the receptor, the histidine-containing phosphotransfer protein, and the type A and B response regulators (RRA and RRB) (Heyl et al., 2013). The *C. braunii* genome encodes members of the first three, but no RRBs (Table 1, Fig. 4). This is contrast to their presence in all chlorophytes and charophytes analysed (Hori et al., 2014; Wang et al., 2015). Several RR domains closely related to RRBs were found, but none contained the Myb domain essential for RRB function (Table S1J). Given the complexity of the *C. braunii* genome, it is possible that not all genes were correctly or completely predicted, but neither genome nor transcriptome data (Data S1H) provide evidence

172 for RRB genes. Their loss suggests either the rewiring of CK signaling or substitution of RRB
173 function by other genes.

174

175 *Ethylene (ETH)*

176 The *C. braunii* genome possesses one or more potential homologs of all of the core components
177 associated with ETH signaling (Table 1, Fig. 4 and Table S1J). *Chara* exhibits ETH-binding
178 activity (Wang et al., 2006), and *C. braunii* encodes several ETH receptor homologs. Notably,
179 *C. braunii* possesses a full-length homolog of *EIN2*, a central regulator in ETH signaling. This
180 is in contrast to both the *K. nitens* genome, which lacks *EIN2* (Hori et al., 2014), and the
181 *Spirogyra pratensis* transcriptome, which shows only a partial *EIN2* sequence (Ju et al., 2015).
182 Except for *EIN2*, *S. pratensis* possesses an ETH signaling pathway that is functionally
183 conserved with the pathway known in land plants (Ju et al., 2015). These findings indicate that
184 the land plant-like ETH signaling pathway was established in the common ancestor of the
185 Phragmoplastophyta, after its divergence from the lineage leading to *Klebsormidium*.

186

187 *Abscisic acid (ABA)*

188 Orthologs of the core ABA signaling components are present in bryophytes and it has been
189 suggested that all ABA biosynthesis/signaling components were gained in the common ancestor
190 of Charophyta (Ju et al., 2015; Wang et al., 2015), with the exception of PYR/PYL receptors
191 that were probably gained in the common ancestor of Zygnematophyceae and land plants (de
192 Vries et al., 2018). The *C. braunii* genome does not contain homologs of the co-receptors
193 ABI/HAB, nor the PYR/RCAR family of receptors (Park et al., 2009), but possesses homologs
194 of genes encoding enzymes that act early in the ABA synthesis pathway (from carotenoid
195 synthesis to violaxanthin; Table 1, Fig. 4 and Table S1J). Given that the presence of ABA has
196 been confirmed in *C. braunii* (Hackenberg and Pandey, 2014), it is likely that the biosynthetic
197 pathway differs from that found in land plants, with ABA possibly being synthesized directly
198 from farnesyl pyrophosphate.

199

200 *Strigolactones (SL)*

201 Orthologs of all the core SL signaling components have been identified exclusively in the
202 genomes of seed plants; however, D14-like receptor homologs are found encoded by
203 bryophytes and charophytes (Bythell-Douglas et al., 2017; Wang et al., 2015). Two SL-related
204 homologs were identified in *C. braunii*, one encoding beta-carotene isomerase D27, and one
205 encoding the candidate SL/karrikin receptor D14-like (Table 1, Fig. 4 and Table S1J). Given
206 the presence of SL in several Charales species, and the activity of the synthetic SL GR24 on
207 *Chara corallina* rhizoid growth (Delaux et al., 2012), it is likely that SL synthesis and signaling
208 differ in charophyceans and in seed plants (Bythell-Douglas et al., 2017). It has been suggested
209 that D14-like proteins might act as (the) SL receptor(s) in this group.

210 In summary, although the phytohormones AUX and CK seem to be ancestral features of
211 streptophytes, and SL and ABA of Phragmoplastophyta (Fig. 1), the respective biosynthesis
212 and/or signaling pathways differ between seed plants and *C. braunii*. Some features of these

213 four phytohormone networks, and of ETH signaling, first appeared in the Phragmoplastophyta
214 as evident by their presence in *C. braunii*. Others were either not present in the ancestor or have
215 since diverged.

216

217 *Plastid evolution: photorespiration and retrograde signaling*

218 Photorespiration, which recycles the 2-carbon compound formed when ribulose bisphosphate
219 carboxylase/oxygenase reacts with oxygen instead of $CO_2$, is crucial to photosynthesis in an
220 oxygen-rich atmosphere. The *C. braunii* genome encodes proteins necessary to carry out a
221 plant-like photorespiratory cycle, including a plant-type glycolate oxidase (GOX) (Table S1M)
222 with structural features preferring glycolate over lactate (Hackenberg et al., 2011). Plant-type
223 GOX is also present in *K. nitens*, while *C. reinhardtii* uses a mitochondrial glycolate
224 dehydrogenase for photorespiratory glycolate metabolism (Nakamura et al., 2005). Apparently,
225 plant-like photorespiration was present in the common ancestor of Streptophyta, the pathway
226 being a feature that might have aided terrestrialization.

227 The plastid to nucleus signaling network optimizes plastid function in land plants. All
228 Chloroplastida (Fig. 1) share EXECUTOR-transduced singlet oxygen and a rudimentary
229 tetrapyrrole-derived retrograde signaling, to which streptophytes recruited GUN2/3 (Fig. 5A).
230 Our data show that *C. braunii,* but not *K. nitens*, encodes GUN1, at which multiple retrograde
231 signals converge in land plants (reviewed by (Chan et al., 2016)). The only GUN1 candidate
232 protein in *K. nitens* (kfl00096_0090) clustered with streptophyte algae- and bryophyte-specific
233 PPRs, but not GUN1 (Data S1I). Hence, retrograde signaling featuring GUN1 might represent
234 an evolutionary novelty of the Phragmoplastophyta (Fig. 1).

235 Plastid-encoded RNA-polymerase (PEP) is the ancestral and for most Archaeplastida the only
236 PEP. In land plants, PEP activity is controlled through PEP-associated proteins (PAPs) (Pfalz
237 and Pfannschmidt, 2013). We detected 5, 8, 10 and 11 PAP orthologs in *C. reinhardtii*, *K. nitens*,
238 *C. braunii*, and *P. patens*, respectively. PAPs were thus already present in streptophyte algae
239 (Fig. 5A) and underwent expansion in land plants. Most of the detected PAPs are predicted to
240 be targeted to the chloroplast, the mitochondrion or both (Table S1N); dual-localization of PAPs
241 to both organelles might be an ancient and conserved character state.

242

243 *Transcriptional regulation*

244 Within the Chloroplastida, morphological complexity correlates with the number of TF (acting
245 in a sequence-specific manner, typically by binding to *cis*-regulatory elements) and
246 transcriptional regulator (TR, acting on chromatin or *via* protein-protein interaction) genes
247 (Lang et al., 2010). We identified 730 TF/TR genes in the *C. braunii* genome (Table S1Q), the
248 complement of such proteins thus being larger than in *K. nitens* (627) or *C. reinhardtii* (542),
249 coinciding with morphological complexity. *C. braunii* encodes several TFs that are not present
250 in other algae, including *K. nitens*. Based on the available data, these families first appear in the
251 Phragmoplastophyta, although they were previously thought to have been gained in the
252 common ancestor of Coleochaetophyceae, Zygnematophyceae and land plants (Wilhelmsson et
253 al., 2017). They include the single canonical ARF mentioned before, as well as TCP, HRT and

Zn cluster TFs (Fig. 1). The *C. braunii* genome encodes two TCP genes, which belong to TCP-P (class I) and TCP-C (class II). The two TCP subgroups are known to exert antagonistic functions in *A. thaliana* with regard to growth proliferation of organs and tissues (Nicolas and Cubas, 2016), implying that the appearance of two different TCP genes might have contributed to regulation of proliferation in the Phragmoplastophyta.

Two separate clades of MADS-box genes exist (Type I and II), with land plant Type II genes further subdivided into so called MIKC$^C$ and MIKC*-type genes (Gramzow and Theißen, 2010). No Type I genes were identified in the *C. braunii* genome, but three Type II genes, of which only *CbMADS1* shows a canonical MIKC-type domain structure. Phylogeny reconstructions together with exon-intron structure analysis (Fig. 5B; Fig. S5, Data S1K) suggest that MIKC$^C$ and MIKC*-type genes evolved from the duplication of an ancestral Type II gene followed by different exon duplications in both gene lineages. As such, *CbMADS1* may serve as a model for the ancestral MIKC-type gene that gave rise to MIKC$^C$- and MIKC*-type genes of land plants.

*C. braunii* encodes 11 bHLH TFs in 5 subfamilies. The Va(2) subfamily is present in chlorophytic and charophytic genomes and not present in land plants, suggesting that this subfamily was lost in the lineage leading to land plants (Table S1O, S1P). The *C. braunii* genome encodes 11 homeodomain (HD) TFs grouped into 9 subfamilies (Table S1O, S1P). Consistent with previous analyses (Catarino et al., 2016), *C. braunii* contains members of the KNOX, BEL, DDT and PINTOX subfamilies that are conserved in chlorophytes.


*Zygotes and spores as analogs to seeds*

Dormant haploid spores of mosses share features of regulation and coat biosynthesis with diploid seeds of flowering plants (Daku et al., 2016; Vesty et al., 2016). The diploid zygotes of *Chara* are dormant diaspores that presumably undergo meiosis and germinate upon suitable environmental cues (Delwiche and Cooper, 2015). Differential expression analysis shows that a number of transcripts related to seed storage proteins (cupin superfamily, oleosins) and to stress tolerance proteins found in seeds (e.g. late embryogenesis abundant), accumulate to high levels in zygotes (Fig. S4). These proteins probably enable the *C. braunii* zygotes to withstand harsh environmental conditions and represent a reservoir of nutrients to facilitate germination and growth. Homologs of these genes have apparently been adopted during land plant evolution to enable dormancy in other diaspores, namely spores and seeds.


**Evolutionary novelties of the *Chara* lineage**

*Trihelix TFs*

The number of TFs per family is lower in *C. braunii* than in land plants for most families, with the trihelix family being an exception: 302 members are encoded, while land plant genomes typically encode ca. 30 copies (Table S1Q). Trihelix TFs are involved in the regulation of development (e.g. embryogenesis, flower development), as well as responses to abiotic and biotic factors. Based on RNA-seq data, at least 28 of the *C. braunii* genes are expressed (Table S4, Fig. S4); 19 in vegetative tissue (of which 6 are expressed exclusively in vegetative tissue)

295 and 22 in reproductive tissues (antheridia, oogonia, zygotes; Fig. S4). Phylogenetic analysis
296 shows that the vast majority of *C. braunii* trihelix paralogs groups outside of the four clades
297 previously defined (Kaplan-Levy et al., 2012) (Data S1J). Similar to secondary expansion of
298 TF families in other lineages the expansion of trihelix TFs in *C. braunii* might be connected to
299 the independent evolution of morphological complexity.

300

301 *Phytohormones: PINs*

302 There are six PIN AUX transporter proteins potentially encoded by the *C. braunii* genome
303 (Table S1J). In land plants, the evolution of morphological complexity in the gametophytic
304 generation, and later in the sporophytic generation, coincides with independent radiations
305 within the *PIN* gene family (Bennett, 2015). Given its high morphological complexity, the same
306 might have occurred in *C. braunii*.

307

308 *Motor network*

309 The evolution of land plants is accompanied by increased abundance of myosin and kinesin
310 domain proteins. Because *K. nitens* has slightly more predicted kinesins than *C. braunii* (Table
311 S1S), it appears that phragmoplast evolution did not depend on the neofunctionalization of
312 kinesin paralogs. However, myosin motors use filamentous actin as tracks. The expansion of
313 the actin family in *C. braunii* (*K. nitens* and *C. reinhardtii* encode 7 actin genes, whereas *C.
314 braunii* has 16; Data S1T, U), with each gene encoding a slightly different protein, hints at
315 varying functions among the cytoskeleton. Land plants have 9 actin genes (*Marchantia
316 polymorpha*) to often 12 (*A. thaliana*, papaya, *Amborella trichocarpa*), and up to 34 in the
317 polyploid maize, while transcriptomic data of other Charales suggests high numbers of
318 underlying genes, e.g. 27 transcripts in *Nitella mirabilis*, 101 in *N. hyalina* (and 46 in the desmid
319 *Penium margaritaceum*). The high numbers of actin genes detected in the amoebal protist
320 *Naegleria gruberi* (86), and the slime mold *Dictyostelium discoideum* (39) (Joseph et al., 2008),
321 can to a large part be explained by their involvement in cell movement. Thus, the additional
322 actin genes of *Chara*, *Nitella* and *Penium* may serve the enhanced cytoplasmic streaming
323 observed in these organisms.

324

325 *Electrical excitability*

326 Inspired by the work of (Hodgkin and Huxley, 1952) on the squid axon, the large internodal
327 cells of *Chara* emerged as an excellent experimental system for electrophysiological studies on
328 plant excitability -  the "Green Axon" (Beilby, 2007). On a slower time scale (1000x), the
329 internodal cells fire action potentials (APs) in response to such as depolarization, light, heat
330 shock, injury or touch. The *C. braunii* genome encodes several putative Touch/Mechano-
331 Sensitive (MS) channels: two members of the MscS-like (MSL) family, as well as an ortholog
332 of the eukaryote specific Piezo-type channel. The negative resting potential (up to -250 mV)
333 across the plasma membrane is generated by the P-type $H^+$-ATPases, encoded in the *C. braunii*
334 genome (Table S1R). $Ca^{2+}$ and $Cl^-$ contribute to the depolarizing phase of the *Chara* AP, while
335 $K^+$ efflux shapes the AP repolarization phase as in animals. No animal-like voltage-gated $Na^+$

336  or $Ca^{2+}$ channels were identified, but a single ALMT-type anion channel gene is present in *C.*
337  *braunii*. The anion channel in *Chara* is $Ca^{2+}$-activated and voltage sensitive, so an Anoctamin-
338  like channel poses another possibility. A Shaker-type, voltage-gated $K^+$ channel in *C. braunii*
339  genome probably mediates the depolarization-activated potassium efflux of the AP
340  repolarization phase. The *C. braunii* habit of long internodal cells might require long distance
341  electrical signaling (Beilby, 2015) enabled by its peculiar set of ion channels. The similarities
342  or differences of *C. braunii* AP, as compared to flowering plants, are yet to be established.

343

344  *Sensing of biotic interaction and microbiome*

345  Land plants harbor a large number of LysM receptor-like kinases (RLK) involved in the
346  perception of chitin-based signals produced by pathogenic and beneficial microorganisms. One
347  member of this family has been described in charophytic algae suggesting either an inability to
348  discriminate microorganisms or an alternative system to do so (Delaux et al., 2015). The *C.*
349  *braunii* genome revealed the presence of seven LysM-RLKs (Fig. 5C; Data S1N) that expanded
350  independently of land plant LysM-RLKs. This expansion may reflect an adaptation of *C.*
351  *braunii* to an extended range of interacting microorganisms (co-cultured bacteria: Table S1T,
352  S1U). This is noteworthy given that many have failed to axenically cultivate Charophyceae,
353  raising the possibility that growth may be dependent on microbiotic commensalism or
354  mutualism.

355

356  *Sexual reproduction and the ROS network*

357  To analyze reproductive mechanisms, transcriptomes of antheridia, oogonia and zygotes were
358  generated (Fig. 5/6, Fig. S6, Table S2 & S3). For antheridia, the data demonstrate that cell
359  motility is up-regulated as expected (Fig. 5D; Fig. S6A). Of 949 differentially expressed genes
360  (DEGs) upregulated in antheridia, 49 encode proteins harboring dynein heavy chains. Dynein-
361  mediated transport is employed in flagellate cells such as spermatozoids and was lost during
362  land plant evolution, concomitant with the loss of motile cells (Rensing et al., 2008). 22 of 302
363  trihelix TFs are expressed in reproductive tissues. Of those, 9 are expressed in all three tissues,
364  5 specifically in antheridia, 7 in oogonia and antheridia, and 1 specifically in the zygote (Fig.
365  S4B). This expression profile may suggest a possible role for these genes in sexual
366  reproduction, in particular in antheridia. Transcripts of a HMG TR and a RWP-RK TF also
367  specifically accumulated in antheridia. Members of these families were shown to be involved
368  in mating in fungi (Barve et al., 2003) and gamete differentiation in *C. reinhardtii* (Lin and
369  Goodenough, 2007), and the single RWP-RK TF in *M. polymorpha* keeps egg cells quiescent
370  in the absence of fertilization (Rovekamp et al., 2016).

371  Zygote transcriptome profiles are characterized by transcription, microtubule-based movement
372  and protein kinase activities (Fig. S6D), processes that might be hallmarks of the diploid zygote
373  maturing and entering dormancy. 87 TFs/TRs are differentially expressed between zygotes and
374  oogonia, among them families typically linked to the regulation of development (e.g. bHLH,
375  HD, AP2/EREBP; Fig. S4C), supporting the hypothesis that transcription undergoes a switch
376  after fertilization. One of the seven LysM RLKs (g44510) is strongly induced in zygotes. In
377  line with potential commensalism mentioned above, this protein might detect the presence of

beneficial microbes as a putative factor triggering meiosis and germination of the dormant zygote.

Of particular interest is the up-regulation of oxidation reduction processes in oogonia as compared to antheridia or zygotes (Fig. 5E; Fig. S6B/C). Like all living organisms, *C. braunii* needs to deal with constitutive production of reactive oxygen species (ROS) using the ROS gene network (Fig. S7, Table S1X). In contrast to land plants, aquatic plants have the option to let ROS diffuse into the water. *C. braunii* encodes all families responsible for ROS scavenging, but with lower gene copy number in comparison to land plants. In contrast, CC-type glutaredoxins (GRX) (ROXYs in *A. thaliana*), which exert crucial functions during angiosperm reproductive development (Gutsche et al., 2015), could not be detected (Table S1X). Among redox-associated genes (Table S1X) the class III peroxidases (Prx), thioredoxins and respiratory burst oxidase homologs expanded greatly during land plant evolution. However, only Prx expanded in *C. braunii* compared to *K. nitens* (Data S1O). With both peroxidative and hydroxylic catalytic cycles, these enzymes can regulate ROS and polymerize cell wall compounds (Francoz et al., 2015). Most of the *C. braunii* Prx are predicted to be secreted, as such, they may contribute to the formation of the strikingly elaborate reproductive structures, e.g. the thick zygote wall (Fig. 2).

7 out of 12 Prx are 2-8 fold higher expressed in oogonia than in antheridia or zygotes (Fig. 6). The higher expression of the ROS gene network could be related to the ROS homeostasis regulation necessary for an optimum fecundation. Flowering plant stigmas exhibit high levels of peroxidase activity when receptive to pollen (McInnis et al., 2006) and have been discussed to be involved in pollen-pistil interaction or pollen tube growth/penetration (Beltramo et al., 2012). For *A. thaliana* root and shoot apical meristems it was shown that stem cell-specific Prx fine tune the balance between superoxide anions ($O_2._-$) and hydrogen peroxide ($H_2O_2$) and thereby affect the switch between cell maintenance and differentiation (Zeng et al., 2017). Differential regulation of ROS levels by Prx might control sexual reproduction in *C. braunii*. Potentially, this mechanism arose in the common ancestor of Phragmoplastophyta and has been recruited from the gametophyte to the sporophyte during land plant evolution.

## Conclusions

The *C. braunii* genome encodes more proteins than other algae, but less than most land plants. Both, specific gains / expansions and losses, can be attributed to the *Chara* lineage (Fig. 1). In absence of a WGD gene family expansions resulted from gene duplication and differential loss. Many of these events likely represent secondary gains in *Chara* complexity *via* sub- and neofunctionalization. We hypothesize that many gene family expansions detected in the *C. braunii* genome underpin its strikingly complex morphology.

Comparative genome analysis clearly reflects the phylogenetic placement of *C. braunii* as a close relative of land plants, with both striking similarities and important differences. It demonstrates the substantial insights into fundamental aspects of plant biology that can be gained by comparing diverse relatives. Molecular signatures across genomes reveal that AUX transport *via* PINs, trihelix TFs, MIKC-type MADS genes as well as photorespiration and

diaspore storage proteins were present prior to the divergence of *K. nitens* (Fig. 1). Other features, such as the non-motile vegetative phase and filamentous growth, evolved later.

Hence, much of what was previously considered land plant-like features clearly evolved in the common ancestor of the Phragmoplastophyta (Fig. 1). These features include polyplastidy, branching, cellulose synthase rosettes, apical cell growth, several features of phytohormone networks, potential involvement of ROS in sexual reproduction and the phragmoplast. Some features evolved after the split of Charophyceae or Coleochaetophyceae such as GRAS TFs and the PPB-like isthmus band of microtubules. Life on land meant increased exposure to UV light. RNA editing repairs UV-B induced mutations in land plants (Maier et al., 2008). Editing evolved after the divergence of Charophyceae from the lineage leading to Zygnematophyceae and land plants (Cahoon et al., 2017). Key editing factors (PPR proteins) are much less abundant in *C. braunii* (57) than in the *Spirogyra* (379) or *P. patens* (100) genomes (Table S1Y). Other features, such as the multicellular sporophyte and embryogenesis, the synthesis of a complex cuticle and the ability to associate with arbuscular mycorrhizal fungi evolved at the base of the land plants, and further during land plant evolution (Fig. 1). Among the latter features are hallmarks of plants' adaptations to land. Yet, before any of these adaptations evolved, LPHGs enabled the first steps of terrestrialization. The key to their identification lies in comparative genomics studies using streptophyte algae, as exemplified here for *C. braunii*.

**Authors' contributions**

AF, AT, ES, HK, HS, JF, JAL, LD, MB, MQ, SAR, SR, SS, TN, YK, YS, YVP provided resources and materials.

AF, AiS, AT, SAR, SR, TN generated the draft genome.

AH, AiS, AsS, BC, CC, CD, CFD, DB, DL, DS-M, DVDS, FBH, FM, FR, GG, GT, GVR, HB, HS, JdV, JH, JMC, JP, KGK, KKU, LIACV, LV, MH, NT, PJ, PKIW, PMD, PU, RH, RK, RS, SAR, SG, SH, SR, SV, SZ, TN analyzed data.

JdV, JAL, LD, SAR, SG, TN wrote the paper.

All authors helped discuss the results and write the paper.

**Declaration of Interests**

The authors declare no competing interests.

**References**

Abouelhoda, M.I., Kurtz, S., and Ohlebusch, E. (2004). Replacing suffix trees with enhanced suffix arrays. Journal of Discrete Algorithms *2*, 53-86.

Alneberg, J., Bjarnason, B.S., de Bruijn, I., Schirmer, M., Quick, J., Ijaz, U.Z., Lahti, L., Loman, N.J., Andersson, A.F., and Quince, C. (2014). Binning metagenomic contigs by coverage and composition. Nat Methods *11*, 1144-1146.

Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Zhang, Z., Miller, W., and Lipman, D.J. (1997). Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. Nucleic Acids Res *25*, 3389-3402.

476 Barve, M.P., Arie, T., Salimath, S.S., Muehlbauer, F.J., and Peever, T.L. (2003). Cloning and
477 characterization of the mating type (MAT) locus from Ascochyta rabiei (teleomorph: Didymella
478 rabiei) and a MAT phylogeny of legume-associated Ascochyta spp. Fungal Genet Biol *39*, 151-
479 167.

480 Bauwe, H., Hagemann, M., and Fernie, A.R. (2010). Photorespiration: players, partners and
481 origin. Trends Plant Sci *15*, 330-336.

482 Beilby, M.J. (2007). Action potential in charophytes. Int Rev Cytol *257*, 43-82.

483 Beilby, M.J. (2015). Salt tolerance at single cell level in giant-celled Characeae. Front Plant Sci
484 *6*, 226.

485 Beilby, M.J., Turi, C.E., Baker, T.C., Tymm, F.J., and Murch, S.J. (2015). Circadian changes in
486 endogenous concentrations of indole-3-acetic acid, melatonin, serotonin, abscisic acid and
487 jasmonic acid in Characeae (Chara australis Brown). Plant Signal Behav *10*, e1082697.

488 Beltramo, C., Torello Marinoni, D., Perrone, I., and Botta, R. (2012). Isolation of a gene
489 encoding for a class III peroxidase in female flower of Corylus avellana L. Mol Biol Rep *39*,
490 4997-5008.

491 Bennett, T. (2015). PIN proteins and the evolution of plant development. Trends Plant Sci *20*,
492 498-507.

493 Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina
494 sequence data. Bioinformatics *30*, 2114-2120.

495 Boot, K.J., Libbenga, K.R., Hille, S.C., Offringa, R., and van Duijn, B. (2012). Polar auxin
496 transport: an early invention. J Exp Bot *63*, 4213-4218.

497 Buschmann, H., and Zachgo, S. (2016). The Evolution of Cell Division: From Streptophyte
498 Algae to Land Plants. Trends Plant Sci *21*, 872-883.

499 Bythell-Douglas, R., Rothfels, C.J., Stevenson, D.W.D., Graham, S.W., Wong, G.K., Nelson,
500 D.C., and Bennett, T. (2017). Evolution of strigolactone receptors by gradual neo-
501 functionalization of KAI2 paralogues. BMC Biol *15*, 52.

502 Cahoon, A.B., Nauss, J.A., Stanley, C.D., and Qureshi, A. (2017). Deep Transcriptome
503 Sequencing of Two Green Algae, Chara vulgaris and Chlamydomonas reinhardtii, Provides No
504 Evidence of Organellar RNA Editing. Genes (Basel) *8*.

505 Catarino, B., Hetherington, A.J., Emms, D.M., Kelly, S., and Dolan, L. (2016). The Stepwise
506 Increase in the Number of Transcription Factor Families in the Precambrian Predated the
507 Diversification of Plants On Land. Mol Biol Evol *33*, 2815-2819.

508 Chan, K.X., Phua, S.Y., Crisp, P., McQuinn, R., and Pogson, B.J. (2016). Learning the
509 Languages of the Chloroplast: Retrograde Signaling and Beyond. Annu Rev Plant Biol *67*, 25-
510 53.

511 Chin, C.S., Alexander, D.H., Marks, P., Klammer, A.A., Drake, J., Heiner, C., Clum, A.,
512 Copeland, A., Huddleston, J., Eichler, E.E.*, et al.* (2013). Nonhybrid, finished microbial
513 genome assemblies from long-read SMRT sequencing data. Nat Methods *10*, 563-569.

514 Daku, R.M., Rabbi, F., Buttigieg, J., Coulson, I.M., Horne, D., Martens, G., Ashton, N.W., and
515 Suh, D.Y. (2016). PpASCL, the *Physcomitrella patens* Anther-Specific Chalcone Synthase-
516 Like Enzyme Implicated in Sporopollenin Biosynthesis, Is Needed for Integrity of the Moss
517 Spore Wall and Spore Viability. PLoS One *11*, e0146817.

518 Darriba, D., Taboada, G.L., Doallo, R., and Posada, D. (2011). ProtTest 3: fast selection of best-
519 fit models of protein evolution. Bioinformatics *27*, 1164-1165.

520 de Vries, J., Curtis, B.A., Gould, S.B., and Archibald, J.M. (2018). Embryophyte stress
521 signaling evolved in the algal progenitors of land plants. Proc Natl Acad Sci U S A.

522 Delaux, P.-M., Xie, X., Timme, R.E., Puech-Pages, V., Dunand, C., Lecompte, E., Delwiche,
523 C.F., Yoneyama, K., Bécard, G., and Séjalon-Delmas, N. (2012). Origin of strigolactones in the
524 green lineage. New Phytologist *195*, 857-871.

525 Delaux, P.M., Radhakrishnan, G.V., Jayaraman, D., Cheema, J., Malbreil, M., Volkening, J.D.,
526 Sekimoto, H., Nishiyama, T., Melkonian, M., Pokorny, L*., et al.* (2015). Algal ancestor of land
527 plants was preadapted for symbiosis. Proc Natl Acad Sci U S A *112*, 13390-13395.

528 Delcher, A.L., Kasif, S., Fleischmann, R.D., Peterson, J., White, O., and Salzberg, S.L. (1999).
529 Alignment of whole genomes. Nucleic Acids Res *27*, 2369-2376.

530 Delwiche, C.F. (2016). The genomes of charophyte green algae. Adv Bot Res *78*, 255-270.

531 Delwiche, C.F., and Cooper, E.D. (2015). The Evolutionary Origin of a Terrestrial Flora.
532 Current biology : CB *25*, R899-910.

533 Dierckxsens, N., Mardulyn, P., and Smits, G. (2017). NOVOPlasty: de novo assembly of
534 organelle genomes from whole genome data. Nucleic Acids Res *45*, e18.

535 Duong, T., Cowling, A., Koch, I., and Wand, M.P. (2008). Feature significance for multivariate
536 kernel density estimation. Computational Statistics & Data Analysis *52*, 4225-4242.

537 Edgar, R.C. (2004). MUSCLE: a multiple sequence alignment method with reduced time and
538 space complexity. BMC Bioinformatics *5*, 113.

539 Flores-Sandoval, E., Eklund, D.M., Hong, S.F., Alvarez, J.P., Fisher, T.J., Lampugnani, E.R.,
540 Golz, J.F., Vazquez-Lobo, A., Dierschke, T., Lin, S.S*., et al.* (2018). Class C ARFs evolved
541 before the origin of land plants and antagonize differentiation and developmental transitions in
542 Marchantia polymorpha. New Phytol *218*, 1612-1630.

543 Francoz, E., Ranocha, P., Nguyen-Kim, H., Jamet, E., Burlat, V., and Dunand, C. (2015). Roles
544 of cell wall peroxidases in plant development. Phytochemistry *112*, 15-21.

545 Gao, X.H., Huang, X.Z., Xiao, S.L., and Fu, X.D. (2008). Evolutionarily conserved DELLA-
546 mediated gibberellin signaling in plants. J Integr Plant Biol *50*, 825-834.

547 Garcia, M., Myouga, F., Takechi, K., Sato, H., Nabeshima, K., Nagata, N., Takio, S., Shinozaki,
548 K., and Takano, H. (2008). An Arabidopsis homolog of the bacterial peptidoglycan synthesis
549 enzyme MurE has an essential role in chloroplast development. Plant J *53*, 924-934.

550 Gnerre, S., Maccallum, I., Przybylski, D., Ribeiro, F.J., Burton, J.N., Walker, B.J., Sharpe, T.,
551 Hall, G., Shea, T.P., Sykes, S*., et al.* (2011). High-quality draft assemblies of mammalian
552 genomes from massively parallel sequence data. Proc Natl Acad Sci U S A *108*, 1513-1518.

553 Gramzow, L., and Theißen, G. (2010). A hitchhiker's guide to the MADS world of plants.
554 Genome Biology *11*, 214.

555 Guindon, S., Dufayard, J.F., Lefort, V., Anisimova, M., Hordijk, W., and Gascuel, O. (2010).
556 New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the
557 performance of PhyML 3.0. Syst Biol *59*, 307-321.

558 Gutsche, N., Thurow, C., Zachgo, S., and Gatz, C. (2015). Plant-specific CC-type
559 glutaredoxins: functions in developmental processes and stress responses. Biol Chem *396*, 495-
560 509.

561 Haas, B.J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P.D., Bowden, J., Couger,
562 M.B., Eccles, D., Li, B., Lieber, M*., et al.* (2013). De novo transcript sequence reconstruction

from RNA-seq using the Trinity platform for reference generation and analysis. Nature protocols *8*, 1494-1512.

Hackenberg, C., Kern, R., Huge, J., Stal, L.J., Tsuji, Y., Kopka, J., Shiraiwa, Y., Bauwe, H., and Hagemann, M. (2011). Cyanobacterial lactate oxidases serve as essential partners in N2 fixation and evolved into photorespiratory glycolate oxidases in plants. Plant Cell *23*, 2978-2990.

Hackenberg, D., and Pandey, S. (2014). Heterotrimeric G-proteins in green algae. An early innovation in the evolution of the plant lineage. Plant Signal Behav *9*, e28457.

Han, G.Z. (2017). Evolution of jasmonate biosynthesis and signaling mechanisms. J Exp Bot *68*, 1323-1331.

Heyl, A., Brault, M., Frugier, F., Kuderova, A., Lindner, A.C., Motyka, V., Rashotte, A.M., Schwartzenberg, K.V., Vankova, R., and Schaller, G.E. (2013). Nomenclature for members of the two-component signaling pathway of plants. Plant Physiol *161*, 1063-1065.

Hodgkin, A.L., and Huxley, A.F. (1952). A quantitative description of membrane current and its application to conduction and excitation in nerve. J Physiol *117*, 500-544.

Hori, K., Maruyama, F., Fujisawa, T., Togashi, T., Yamamoto, N., Seo, M., Sato, S., Yamada, T., Mori, H., Tajima, N.*, et al.* (2014). *Klebsormidium flaccidum* genome reveals primary factors for plant terrestrial adaptation. Nature Communications *5*, 3978.

Huang, X., and Madan, A. (1999). CAP3: A DNA sequence assembly program. Genome Res *9*, 868-877.

Huelsenbeck, J.P., and Ronquist, F. (2001). MRBAYES: Bayesian inference of phylogenetic trees. Bioinformatics *17*, 754-755.

Huson, D.H., Beier, S., Flade, I., Gorska, A., El-Hadidi, M., Mitra, S., Ruscheweyh, H.J., and Tappu, R. (2016). MEGAN Community Edition - Interactive Exploration and Analysis of Large-Scale Microbiome Sequencing Data. PLoS Comput Biol *12*, e1004957.

Inupakutika, M.A., Sengupta, S., Devireddy, A.R., Azad, R.K., and Mittler, R. (2016). The evolution of reactive oxygen species metabolism. J Exp Bot *67*, 5933-5943.

Iseli, C., Jongeneel, C.V., and Bucher, P. (1999). ESTScan: a program for detecting, evaluating, and reconstructing potential coding regions in EST sequences. In Proc Int Conf Intell Syst Mol Biol (Menlo Park, CA, USA: American Association for Artificial Intelligence ), pp. 138-148.

Joseph, J.M., Fey, P., Ramalingam, N., Liu, X.I., Rohlfs, M., Noegel, A.A., Muller-Taubenberger, A., Glockner, G., and Schleicher, M. (2008). The actinome of Dictyostelium discoideum in comparison to actins and actin-related proteins from other organisms. PLoS One *3*, e2654.

Jouffroy, O., Saha, S., Mueller, L., Quesneville, H., and Maumus, F. (2016). Comprehensive repeatome annotation reveals strong potential impact of repetitive elements on tomato ripening. BMC Genomics *17*, 624.

Ju, C., Van de Poel, B., Cooper, E.D., Thierer, J.H., Gibbons, T.R., Delwiche, C.F., and Chang, C. (2015). Conservation of ethylene as a plant hormone over 450 million years of evolution. Nat Plants *1*, 14004.

Kaplan-Levy, R.N., Brewer, P.B., Quon, T., and Smyth, D.R. (2012). The trihelix family of transcription factors--light, stress and development. Trends Plant Sci *17*, 163-171.

Katoh, K., and Standley, D.M. (2013). MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability. Molecular Biology and Evolution *30*, 772-780.

606   Keller, O., Kollmar, M., Stanke, M., and Waack, S. (2011). A novel hybrid gene prediction
607   method employing protein multiple sequence alignments. Bioinformatics *27*, 757-763.

608   Kelley, L.A., Mezulis, S., Yates, C.M., Wass, M.N., and Sternberg, M.J. (2015). The Phyre2
609   web portal for protein modeling, prediction and analysis. Nat Protoc *10*, 845-858.

610   Kent, W.J. (2002). BLAT--the BLAST-like alignment tool. Genome Res *12*, 656-664.

611   Kim, D., Pertea, G., Trapnell, C., Pimentel, H., Kelley, R., and Salzberg, S.L. (2013). TopHat2:
612   accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions.
613   Genome Biol *14*, R36.

614   Köster, J., and Rahmann, S. (2012). Snakemake--a scalable bioinformatics workflow engine.
615   Bioinformatics *28*, 2520-2522.

616   Kwantes, M., Liebsch, D., and Verelst, W. (2012). How MIKC* MADS-box genes originated
617   and evidence for their conserved function throughout the evolution of vascular plant
618   gametophytes. Mol Biol Evol *29*, 293-302.

619   Lagesen, K., Hallin, P., Rodland, E.A., Staerfeldt, H.H., Rognes, T., and Ussery, D.W. (2007).
620   RNAmmer: consistent and rapid annotation of ribosomal RNA genes. Nucleic Acids Res *35*,
621   3100-3108.

622   Lang, D., Ullrich, K.K., Murat, F., Fuchs, J., Jenkins, J., Haas, F.B., Piednoel, M., Gundlach,
623   H., Van Bel, M., Meyberg, R.*, et al.* (2018). The Physcomitrella patens chromosome-scale
624   assembly reveals moss genome structure and evolution. Plant J *93*, 515-533.

625   Lang, D., Weiche, B., Timmerhaus, G., Richardt, S., Riano-Pachon, D.M., Correa, L.G., Reski,
626   R., Mueller-Roeber, B., and Rensing, S.A. (2010). Genome-wide phylogenetic comparative
627   analysis of plant transcriptional regulation: a timeline of loss, gain, expansion, and correlation
628   with complexity. Genome Biol Evol *2*, 488-503.

629   Lawrence, M., Huber, W., Pages, H., Aboyoun, P., Carlson, M., Gentleman, R., Morgan, M.T.,
630   and Carey, V.J. (2013). Software for computing and annotating genomic ranges. PLoS Comput
631   Biol *9*, e1003118.

632   Lecointre, G., and Le Guyader, H. (2006). The Tree of Life: A Phylogenetic Classification
633   (Harvard University Press).

634   Lee, E., Helt, G.A., Reese, J.T., Munoz-Torres, M.C., Childers, C.P., Buels, R.M., Stein, L.,
635   Holmes, I.H., Elsik, C.G., and Lewis, S.E. (2013). Web Apollo: a web-based genomic
636   annotation editing platform. Genome Biol *14*, R93.

637   Leinonen, R., Sugawara, H., Shumway, M., and International Nucleotide Sequence Database,
638   C. (2011). The sequence read archive. Nucleic Acids Res *39*, D19-21.

639   Li, B., and Dewey, C.N. (2011). RSEM: accurate transcript quantification from RNA-Seq data
640   with or without a reference genome. BMC Bioinformatics *12*, 323.

641   Li, H., and Durbin, R. (2010). Fast and accurate long-read alignment with Burrows-Wheeler
642   transform. Bioinformatics *26*, 589-595.

643   Lin, H., and Goodenough, U.W. (2007). Gametogenesis in the Chlamydomonas reinhardtii
644   minus mating type is controlled by two genes, MID and MTD1. Genetics *176*, 913-925.

645   Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and
646   dispersion for RNA-seq data with DESeq2. Genome Biol *15*, 550.

647   Maier, U.G., Bozarth, A., Funk, H.T., Zauner, S., Rensing, S.A., Schmitz-Linneweber, C.,
648   Borner, T., and Tillich, M. (2008). Complex chloroplast RNA metabolism: just debugging the
649   genetic programme? BMC Biol *6*, 36.

650  Marcais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting
651  of occurrences of k-mers. Bioinformatics *27*, 764-770.

652  McInnis, S.M., Desikan, R., Hancock, J.T., and Hiscock, S.J. (2006). Production of reactive
653  oxygen species and reactive nitrogen species by angiosperm stigmas and pollen: potential
654  signalling crosstalk? New Phytol *172*, 221-228.

655  Morris, J.L., Puttick, M.N., Clark, J.W., Edwards, D., Kenrick, P., Pressel, S., Wellman, C.H.,
656  Yang, Z., Schneider, H., and Donoghue, P.C.J. (2018). The timescale of early land plant
657  evolution. Proc Natl Acad Sci U S A.

658  Nakamura, Y., Kanakagiri, S., Van, K., He, W., and Spalding, M.H. (2005). Disruption of the
659  glycolate dehydrogenase gene in the high-CO2-requiring mutant HCR89 of *Chlamydomonas*
660  *reinhardtii*. Canadian Journal of Botany *83*, 820-833.

661  Nguyen, L.T., Schmidt, H.A., von Haeseler, A., and Minh, B.Q. (2015). IQ-TREE: a fast and
662  effective stochastic algorithm for estimating maximum-likelihood phylogenies. Mol Biol Evol
663  *32*, 268-274.

664  Nicolas, M., and Cubas, P. (2016). TCP factors: new kids on the signaling block. Curr Opin
665  Plant Biol *33*, 33-41.

666  Ohtaka, K., Hori, K., Kanno, Y., Seo, M., and Ohta, H. (2017). Primitive Auxin Response
667  without TIR1 and Aux/IAA in the Charophyte Alga Klebsormidium nitens. Plant Physiol *174*,
668  1621-1632.

669  Park, S.Y., Fung, P., Nishimura, N., Jensen, D.R., Fujii, H., Zhao, Y., Lumba, S., Santiago, J.,
670  Rodrigues, A., Chow, T.F.*, et al.* (2009). Abscisic acid inhibits type 2C protein phosphatases via
671  the PYR/PYL family of START proteins. Science *324*, 1068-1071.

672  Parra, G., Bradnam, K., and Korf, I. (2007). CEGMA: a pipeline to accurately annotate core
673  genes in eukaryotic genomes. Bioinformatics *23*, 1061-1067.

674  Pfalz, J., and Pfannschmidt, T. (2013). Essential nucleoid proteins in early chloroplast
675  development. Trends Plant Sci *18*, 186-194.

676  Pickett-Heaps, J.D. (1975). Green Algae: Structure. Reproduction and Evolution in Selected
677  Genera (Sinauer).

678  Pringsheim, M. (1862). On the Pro-Embryos of the Charae. The Annals and Magazine of
679  Natural History *59*, 321-326.

680  Pruesse, E., Peplies, J., and Glockner, F.O. (2012). SINA: accurate high-throughput multiple
681  sequence alignment of ribosomal RNA genes. Bioinformatics *28*, 1823-1829.

682  Puranik, S., Acajjaoui, S., Conn, S., Costa, L., Conn, V., Vial, A., Marcellin, R., Melzer, R.,
683  Brown, E., Hart, D.*, et al.* (2014). Structural basis for the oligomerization of the MADS domain
684  transcription factor SEPALLATA3 in Arabidopsis. Plant Cell *26*, 3603-3615.

685  Quinlan, A.R. (2014). BEDTools: The Swiss-Army Tool for Genome Feature Analysis. Curr
686  Protoc Bioinformatics *47*, 11 12 11-34.

687  Ren, Y., Hansen, S.F., Ebert, B., Lau, J., and Scheller, H.V. (2014). Site-directed mutagenesis
688  of IRX9, IRX9L and IRX14 proteins involved in xylan biosynthesis: glycosyltransferase
689  activity is not required for IRX9 function in Arabidopsis. PLoS One *9*, e105014.

690  Rensing, S.A. (2018). Great moments in evolution: the conquest of land by plants. Curr Opin
691  Plant Biol *42*, 49-54.

692 Rensing, S.A., Ick, J., Fawcett, J.A., Lang, D., Zimmer, A., Van de Peer, Y., and Reski, R.
693 (2007). An ancient genome duplication contributed to the abundance of metabolic genes in the
694 moss *Physcomitrella patens*. BMC Evol Biol *7*, 130.

695 Rensing, S.A., Lang, D., Zimmer, A.D., Terry, A., Salamov, A., Shapiro, H., Nishiyama, T.,
696 Perroud, P.-F., Lindquist, E.A., Kamisugi, Y.*, et al.* (2008). The *Physcomitrella* genome reveals
697 evolutionary insights into the conquest of land by plants. Science *319*, 64-69.

698 Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D.L., Darling, A., Hohna, S., Larget, B.,
699 Liu, L., Suchard, M.A., and Huelsenbeck, J.P. (2012). MrBayes 3.2: Efficient Bayesian
700 Phylogenetic Inference and Model Choice Across a Large Model Space. Syst Biol *61*, 539-542.

701 Rost, B. (1999). Twilight zone of protein sequence alignments. Protein Eng *12*, 85-94.

702 Rovekamp, M., Bowman, J.L., and Grossniklaus, U. (2016). *Marchantia* MpRKD Regulates
703 the Gametophyte-Sporophyte Transition by Keeping Egg Cells Quiescent in the Absence of
704 Fertilization. Curr Biol *26*, 1782-1789.

705 Saier, M.H., Jr., Reddy, V.S., Tsu, B.V., Ahmed, M.S., Li, C., and Moreno-Hagelsieb, G. (2016).
706 The Transporter Classification Database (TCDB): recent advances. Nucleic Acids Res *44*,
707 D372-379.

708 Sakayama, H., Kasai, F., Nozaki, H., Watanabe, M.M., Kawachi, M., Shigyo, M., Nishihiro, J.,
709 Washitani, I., Krienitz, L., and Ito, M. (2009). Taxonomic reexamination of *Chara globularis*
710 (Charales, Charophyceae) from Japan based on oospore morphology and rbcL gene sequences,
711 and the description of *C. leptospora* sp. nov. J Phycol *45*, 917-927.

712 Schwacke, R., Schneider, A., van der Graaff, E., Fischer, K., Catoni, E., Desimone, M.,
713 Frommer, W.B., Flugge, U.I., and Kunze, R. (2003). ARAMEMNON, a novel database for
714 Arabidopsis integral membrane proteins. Plant Physiol *131*, 16-26.

715 Scrucca, L., Fop, M., Murphy, T.B., and Raftery, A.E. (2016). mclust 5: Clustering,
716 Classification and Density Estimation Using Gaussian Finite Mixture Models. R J *8*, 289-317.

717 Simao, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., and Zdobnov, E.M. (2015).
718 BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs.
719 Bioinformatics *31*, 3210-3212.

720 Sorensen, I., Pettolino, F.A., Bacic, A., Ralph, J., Lu, F., O'Neill, M.A., Fei, Z., Rose, J.K.,
721 Domozych, D.S., and Willats, W.G. (2011). The charophycean green algae provide insights into
722 the early origins of plant cell walls. Plant J *68*, 201-211.

723 Steinbiss, S., Willhoeft, U., Gremme, G., and Kurtz, S. (2009). Fine-grained annotation and
724 classification of de novo predicted LTR retrotransposons. Nucleic Acids Res *37*, 7002-7013.

725 Tan, X., Calderon-Villalobos, L.I., Sharon, M., Zheng, C., Robinson, C.V., Estelle, M., and
726 Zheng, N. (2007). Mechanism of auxin perception by the TIR1 ubiquitin ligase. Nature *446*,
727 640-645.

728 Theißen, G., Melzer, R., and Rümpler, F. (2016). MADS-domain transcription factors and the
729 floral quartet model of flower development: linking plant development and evolution.
730 Development *143*, 3259-3271.

731 Timme, R.E., Bachvaroff, T.R., and Delwiche, C.F. (2012). Broad phylogenomic sampling and
732 the sister lineage of land plants. PLoS ONE *7*, e29696.

733 Tivendale, N.D., Ross, J.J., and Cohen, J.D. (2014). The shifting paradigms of auxin
734 biosynthesis. Trends Plant Sci *19*, 44-51.

735  Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J., Salzberg,
736  S.L., Wold, B.J., and Pachter, L. (2010). Transcript assembly and quantification by RNA-Seq
737  reveals unannotated transcripts and isoform switching during cell differentiation. Nat
738  Biotechnol *28*, 511-515.

739  Van Bel, M., Proost, S., Wischnitzki, E., Movahedi, S., Scheerlinck, C., Van de Peer, Y., and
740  Vandepoele, K. (2012). Dissecting plant genomes with the PLAZA comparative genomics
741  platform. Plant Physiol *158*, 590-600.

742  Vesty, E.F., Saidi, Y., Moody, L.A., Holloway, D., Whitbread, A., Needs, S., Choudhary, A.,
743  Burns, B., McLeod, D., Bradshaw, S.J*., et al.* (2016). The decision to germinate is regulated by
744  divergent molecular networks in spores and seeds. New Phytol *211*, 952-966.

745  Vriet, C., Lemmens, K., Vandepoele, K., Reuzeau, C., and Russinova, E. (2015). Evolutionary
746  trails of plant steroid genes. Trends Plant Sci *20*, 301-308.

747  Walker, K.L., Muller, S., Moss, D., Ehrhardt, D.W., and Smith, L.G. (2007). Arabidopsis
748  TANGLED identifies the division plane throughout mitosis and cytokinesis. Curr Biol *17*,
749  1827-1836.

750  Wang, C., Liu, Y., Li, S.S., and Han, G.Z. (2015). Insights into the origin and evolution of the
751  plant hormone signaling machinery. Plant Physiol *167*, 872-886.

752  Wang, W., Esch, J.J., Shiu, S.H., Agula, H., Binder, B.M., Chang, C., Patterson, S.E., and
753  Bleecker, A.B. (2006). Identification of important regions for ethylene binding and signaling in
754  the transmembrane domain of the ETR1 ethylene receptor of Arabidopsis. Plant Cell *18*, 3429-
755  3442.

756  Wass, M.N., Kelley, L.A., and Sternberg, M.J. (2010). 3DLigandSite: predicting ligand-binding
757  sites using similar structures. Nucleic Acids Res *38*, W469-473.

758  Wickett, N.J., Mirarab, S., Nguyen, N., Warnow, T., Carpenter, E., Matasci, N., Ayyampalayam,
759  S., Barker, M.S., Burleigh, J.G., Gitzendanner, M.A*., et al.* (2014). Phylotranscriptomic
760  analysis of the origin and early diversification of land plants. Proc Natl Acad Sci U S A *111*,
761  E4859-4868.

762  Wilhelmsson, P.K.I., Muhlich, C., Ullrich, K.K., and Rensing, S.A. (2017). Comprehensive
763  Genome-Wide Classification Reveals That Many Plant-Specific Transcription Factors Evolved
764  in Streptophyte Algae. Genome Biol Evol *9*, 3384-3397.

765  Xiong, W., He, L., Lai, J., Dooner, H.K., and Du, C. (2014). HelitronScanner uncovers a large
766  overlooked cache of Helitron transposons in many plant genomes. Proc Natl Acad Sci U S A
767  *111*, 10263-10268.

768  Yang, Z. (2007). PAML 4: phylogenetic analysis by maximum likelihood. Mol Biol Evol *24*,
769  1586-1591.

770  Zeng, J., Dong, Z., Wu, H., Tian, Z., and Zhao, Z. (2017). Redox regulation of plant stem cell
771  fate. EMBO J *36*, 2844-2855.

772

773

774

**Figure titles and legends**

776

777 **Figure 1: Evolution of charophytic algae and land plant features**

778 Cladogram symbolizing streptophytic evolution shows gain/expansion (green lines) and loss
779 (red lines) of features; topology as in (Wickett et al., 2014) with phytohormone-related terms
780 in blue and transcription factors (TF) and transcriptional regulators (TR) in brown. Expansions
781 (and gains/losses) detected in the *Chara* lineage are shown by asterisk. See text for
782 abbreviations. Modes of cytokinesis: a cleavage furrow with persistent telophase spindle as
783 seen in *Klebsormidium*, and a phragmoplast seen in *Chara* that differs from that of land plants
784 as the cell plate in *Chara* shows little centrifugal growth but is formed simultaneously across
785 the cell's equator.

786

787 **Figure 2: Life cycle and habit of *Chara braunii***

788 Meiosis occurs just prior to germination. At germination, a positively gravitropic rhizoid and a
789 protonema that develops into the thallus are formed. The shoot-like thallus (phototropic and
790 negatively gravitropic) comprises stem-like structures (axes) and whorls of branchlets (lateral
791 organs appended to the main axes having adaxial-abaxial differentiation) at axial nodes. Growth
792 of the axis/stem is axial from the terminal (apical) cell. Internodal cells, up to 5 cm long, are
793 multinucleate. Internodal cells and branchlets are connected *via* specialized nodes or central
794 cells connecting the internodes. Nodal cells can serve asexual propagation as they can form
795 apical cells *de novo*. Female (oogonia) and male (antheridia) gametangia are borne on branchlet
796 nodes of the monoicous thalli and generate female (egg cell) and male (sperm cell) gametes.
797 The oogonial complex is comprised of egg cell and associated corona, jacket (five spiral tube
798 cells), and basal cells. Sperm cells arise from filaments produced on the inner surfaces of
799 antheridial shield cells. Upon fertilization the only diploid cell of the life cycle, the dormant
800 zygote or oospore, is formed. Charasomes are plasmamembrane invaginations that allow carbon
801 concentration *via* local acidification. Cells are connected by plasmodesmata. Actin-myosin
802 based cytoplasmic streaming provides a fast transport mechanism. *C. braunii* is ecorticate, other
803 species develop cortical cells (filaments with spine cells) from the nodes that cover the axis and
804 branchlet internodal cells. LS: longitudinal section.

805

806 **Figure 3: Gene and transposon length and density in selected plant and algal genomes.**

807 Comparative box and whisker plots depicting distributions of feature lengths (A) and densities
808 in 100 kbp windows (B). Organisms are ordered top-down by decreasing genome size; x-axes
809 are logarithmic scale. Features are color-coded (legend on the right) and comprise predicted
810 genes, helitrons, intact full-length long terminal repeat elements (flLTRE) and potentially
811 fragmented copies (LTREs).

812

813 **Figure 4: Overview of predicted presence of factors in phytohormone biosynthesis and**
814 **signaling pathways of *C. braunii*.**

Shown are biosynthesis enzymes (rectangles), receptors (pentagons), signal transduction components (hexagons), and TFs (ovals). Elements for which no orthologs were found (light green dashed boxes) and for which putative orthologs were identified (dark green boxes) (*cf.* Table 1, S10/11). Abbreviations as in Table 1.


**Figure 5: Land plant heritage genes present in the *C. braunii* genome**

 (A) Growing repertoire of retrograde signaling components as well as PAPs along the streptophyte trajectory. Potential retrograde signaling orthologs are marked with colored dots (see species key). PEP-associated proteins (PAPs) are shown in the bottom inset. XRN2/XRN3 were not distinguished due to paralogy; faded dots mark the paralogy of *Chlamydomonas* FSD2 and the detection of *P. patens* PTAC7 ortholog with $E<10^{-4}$; mosses encode the cyanobacterial (i.e. non-PAP) version of MurE (Garcia et al., 2008), potentially applying for algal MurEs, too. (B) Bayesian inference phylogenetic tree of plant MADS-box genes. Posterior probabilities ($\geq 0.6$) of main branches are depicted next to the tree. Insert shows the exon-intron structures of representatives of MIKC$^C$-type genes together with the *Chara* MIKC-type genes. (C) Condensed ML tree of the LysM-RLK family. The Charales sequences form a single clade (blue branches) encompassing 7 *C. braunii* sequences. Duplication (red circle) leading to the LYK (orange) and LYR (green) subclades occurred at the base of the embryophytes. The moss and liverwort clades are clustered. (D, E) GO enrichment word clouds (biological process). Word clouds of genes down-regulated (D) or up-regulated (E) in oogonia as compared to antheridia. Font size correlates with significance; red terms are depleted, green terms enriched; top three terms each are shown. See also Fig. S5, S6, related to Table S2-S4.


**Figure 6: Expression of the ROS gene network during sexual reproduction.**

ROS-related gene abundance expressed in transcripts per million (TPM) was transformed to log scale and represented as heatmap in zygotes, oogonia and antheridia. Gene distance was calculated using the Euclidean method and genes were clustered using complete linkage. DEGs ($p < 0.01$) between zygotes and oogonia / oogonia and antheridia are depicted: green up arrow, log2(fold-change) > 0; red down arrow, log2(fold-change) < 0. The expanded family of class III peroxidases is shown in bold. See also Fig. S4, S6, S7, related to Table S2-S4.

**Tables**

| Gene/<br>Gene family | K. nitens | C. braunii | P. patens | A. thaliana |
|---|---|---|---|---|
| AUX biosynthesis | | | | |
| Tryptophan aminotransferase-related proteins (TAA/TAR) | 1 | 0 | 6 | 5 |
| YUCCA (YUC) | 1 | 0 | 8 | 11 |
| AUX signaling | | | | |
| Transport inhibitor response 1/AUX signaling F-box (TIR1/AFB) | 0 | 0 | 5 | 5 |
| AUX response factor (ARF) | 0 | 1 | 15 | 22 |
| Indole-3-acetic acid inducible (Aux/IAA) | 1/0[a] | 2 | 4 | 29 |
| AUX metabolism | | | | |
| Gretchenhagen (GH) | 4 | 1 | 2 | 20 |
| AUX transport | | | | |
| ATP-binding cassette B (ABCB) | 7 | 5 | 10 | 22 |
| AUX resistance 1 (AUX1/LAX) | 1 | 0 | 9 | 4 |
| PIN-formed 1 (PIN) | 1 | 6 | 4 | 8 |
| PIN-likes 1 (PILS) | 3 | 0 | 3 | 7 |
| CK Signaling | | | | |
| CHASE domain containing histidine kinase (CHK) | 6 | 2 | 11 | 3 |
| Histidine-containing phosphotransfer proteins (HPT) | 1 | 1 | 2 | 5 |
| Response regulator type B (RRB) | 1 | 0 | 5 | 11 |
| Response regulator type A (RRA) | 1 | 2 | 7 | 10 |
| ETH biosynthesis | | | | |
| 1-aminocyclopropane-1-carboxylate synthase (ACS) | 1 | 2 | 2 | 12 |
| 1-aminocyclopropane-1-carboxylate oxidase (ACO) | 0 | 0 | 0 | 5 |
| ETH signaling | | | | |
| ETH response/ETH response sensor (ETR/ERS) | 5 | 4 | 8 | 5 |
| Constitutive triple response1 (CTR1) | 1 | 2 | 1 | 1 |
| ETH insensitive2 (EIN2) | 0 | 1 | 2 | 1 |
| ETH insensitive3 (EIN3) | 1 | 4 | 2 | 6 |
| EIN3 binding F-box protein (EBF1) | 1 | 1 | 2 | 2 |

| | | | | |
|---|---|---|---|---|
| ABA biosynthesis | | | | |
| Phytoene synthase1 (PSY1) | 1 | 1 | 3 | 1 |
| Phytoene desaturase (PDS) | 2 | 1 | 2 | 1 |
| Lutein deficient (LUT) | 1 | 1 | 1 | 3 |
| Zeaxanthin epoxidase (ZEP/ABA1) | 1 | 1 | 1 | 1 |
| 9-Cis-epoxycarotenoid dioxygenase (NCED) | 0 | 0 | 2 | 5 |
| Abscisic aldehyde oxidase3 (AAO3) | 1 | 0 | 2 | 1 |
| ABA signaling | | | | |
| Pyrabactin resistance (PYR) | 0 | 0 | 4 | 14 |
| Protein phosphatase 2C (PP2C - Group A) | 1 | 0 | 2 | 9 |
| SNF related kinase (SnRK) | 1 | 1 | 4 | 5 |
| CBL-interacting protein kinase (CIPK) | 1 | 0 | 7 | 25 |
| Calcium-dependent protein kinase (CPK) | 1 | 2 | 30 | 34 |
| SL synthesis | | | | |
| Beta-carotene isomerase (D27) | 2 | 1 | 1 | 1 |
| Carotenoid cleavage dioxygenase (CCD7) | 2 | 0 | 1 | 1 |
| Carotenoid cleavage dioxygenase (CCD8) | 2 | 0 | 1 | 1 |
| SL signaling | | | | |
| Alfa beta hydrolase (D14) | 0 | 0 | 0 | 1 |
| D14-like/ Karrikin insensitive2 (KAI2) | 2 | 1 | 11 | 2 |
| More axillary branching 2 (MAX2) | 0 | 0 | 1 | 1 |

846

**Table 1: Comparison of gene families operating in the biosynthesis and signaling networks of phytohormones.**

A specific set of individual genes or gene families encoding steps in the phytohormone biosynthesis/signaling/metabolism/transport networks have been analysed in *K. nitens*, *C. braunii*, *P. patens* and *A. thaliana* (Table S1J).

a, kfl00094_0070 features Aux/IAA domains but also a B3 domain (see text for details).

853

854

**Supplemental Figure Titles and Legends**

**Figure S1, related to STAR methods: Chromosomes in an antheridial filament of *C. braunii* (n=14, strain S276).**

The chromosomes during cell division in young antheridial filaments of strain S276 were observed after Feulgen staining. The chromosome number n=14 was confirmed by counts made on chromosomes during metaphase or anaphase. Most Chara species have either n=14 or n=28 chromosomes, Nitella and the other genera have different base numbers. There are numerous examples of monoecious/dioecious species pairs in the family, with the dioecious species always displaying half the number of chromosomes than their monoecious counterpart. For Chara typically dioecious=14, monoecious=28 (or other multiples of 14). *C. braunii* is monoecious, but is unique in having the dioecious chromosome number of 14. There are no known dioecious sister taxa to *C. braunii*, perhaps due to the already reduced genome. Scale bar = 2 μm.

**Figure S2, related to STAR methods: Assembly characteristics and decontamination**

A) k-mer frequency analysis of the S276 paired end read data with k=25. Number of 25-mers at frequency 3 to 200 are shown with the solid line. Circles shows the points from 16 to 80 as what was recognized the major peak, presumably representing the single copy region in *C. braunii*. B) Scatter plot of mapped reads of two *C. braunii* strains on each scaffold. Blue and light blue points are scaffolds with GC content of at least 55% and less than 55%, respectively. C) Frequency distribution of scaffold wise GC content compared between putative *C. braunii* derived scaffolds (blue) and other scaffolds (green).

**Figure S3, related to STAR methods: Ks-based analysis of *C. braunii* paralogs**

Paranome-based WGD signature prediction. (A) Ks frequency plot highlighting mixture model components mean and standard-deviation (top: #component, bottom: mean Ks) based on raw Ks value classification. (B) Ks frequency plot highlighting mixture model components mean and standard-deviation (top: #component, bottom: mean Ks) based on log-transformed Ks value classification. (C) Ks group assignment for raw Ks classification. (D) Ks group assignment for log-transformed Ks classification. (E) Significant zero crossing (SiZer) plot. (F) Significant convexity (SiCon) plot. (G-J) Significant features of kernel density estimates using indicated bandwidths, highlighting significant gradient regions in blue and significant curvature regions in green using a significance level of 0.05. Red vertical lines represent Ks value of 0.1 and 2.0, dotted red vertical line represents Ks value of 0.235 corresponding to 12.5 Ma ago (these events might be no WGDs but only more or less recent local duplication events). For *C. braunii* no single predicted WGD signature was supported by three different bandwidth kernel densities (cf. STAR Methods).

**Figure S4, related to Figure 6: Expression profiles during sexual reproduction.**

Expression profile of trihelix TF genes based on RNA-seq evidence (Table S4) was visualized as A) a Venn diagram using venny (http://bioinfogp.cnb.csic.es/tools/venny/) and B) as a heatmap showing gene expression and DEGs from reproductive organs with RPKM > 1 in minimum two samples. C) Shows expression of differentially expressed TFs/TRs during sexual reproduction. D) Expression of DEGs associated with seeds during sexual reproduction. Transcripts per million (TPM) were transformed to log2 scale and clustered using the euclidean distance method and the complete clustering method (B, C, D).

**Figure S5, related to Figure 5: Exon-intron structure comparison of MIKC$^C$-type, MIKC*-type and charophyte MIKC-type genes.**

(A) Exon-intron structures of representatives of MIKC$^C$-type and MIKC*-type genes together with the charophyte MIKC-type genes *CbMADS1*, *CbMADS2* and *KnMADS1*. The exons encoding MADS-, I- , K- and C-domains are color coded in black, red, blue and green, respectively. Among the three Type II genes that were identified in the *C. braunii* genome only *CbMADS1* shows a canonical MIKC-type gene sequence. In contrast *CbMADS2* lacks most (but not all) introns and thus probably evolved via a retrotransposition and recombination event. *CbMADS3* lacks the conserved K-box that encodes for the protein-protein interacting K-domain (data not shown). (B and C) Analysis of exon-intron structures suggest that *CbMADS1* directly descends from an ancestral MIKC-type gene that was a common ancestor of MIKC$^C$- and MIKC*-type genes. (B) It was previously suggested that the N-terminal part of the K-domain of MIKC*-type proteins evolved through a duplication of two K-domain exons of an ancestral MIKC-type gene (Kwantes et al., 2012). The aligned amino acid sequences encoded by exon 2 of *CbMADS1*, and by the first K-domain exons of *KnMADS1*, *MpMADS1*, *PPM3*, *SmMADS4* and *AGL30* indeed strongly support this hypothesis. (C) In addition, striking similarities between the aligned amino acid sequences encoded by exon 5 of *CbMADS1*, exon 6 of *KnMADS1* and exons 5 and 6 of *MpMADS2*, *PPM1*, *SmMADS3* and *SEP3*, respectively, suggest that also the K-domain of MIKC$^C$-type proteins evolved through an exon duplication of an ancestral MIKC-type gene. This is especially intriguing considering the fact that, based on structural data, the last two K-domain exons of most if not all MIKC$^C$-type genes encode for a protein-protein interaction interface that facilitates tetramer formation of MIKC$^C$-type proteins (Puranik et al., 2014). It has already been suggested that the ability of MIKC$^C$-type proteins to tetramerize was an important precondition to evolve and diversify efficient developmental switches that facilitated the transition to land and the evolution of complex body plans of land plants (Theißen et al., 2016). Thus it is tempting to speculate that an exon duplication of an ancestral MIKC$^C$-type gene in the MRCA of extant land plants created the molecular prerequisites for this evolutionary novelty.

**Figure S6, related to Figures 5/6: Transcriptome analyses of reproduction and early development.**

GO enrichment word clouds (category biological process); genes down-regulated (A) or up-regulated (B) in oogonia as compared to antheridia, genes down-regulated (C) or up-regulated (D) in zygotes as compared to oogonia. Antheridia are strongly enriched with the GO category GO:0015074 "DNA integration" (A). 349 gene models expressed in antheridia were classified in this category; of these, 324 genes were found to be overlapping with a TE to at least 50 % (Table S4). Most of these genes were annotated as "integrase", "ribonuclease H-like", "reverse transcriptase", and "aspartyl protease" by homology-based approach, terms typical of Ty3/Gypsy pol gene composition (Havecker et al., 2004). Ty3/Gypsy elements represent 20 % of the *C. braunii* genome. These results might indicate mobilization of retrotransposons and other mobile elements during male gametogenesis. This could be a consequence of genome rearrangement during male gamete formation. One could also imagine that mobilization and integration of retrotransposons might enhance genomic diversity during sexual reproduction.

**Figure S7, related to Figure 6: Major reactive oxygen species scavenging pathway in plants.**

Proteins associated with ROS scavenging are in bold. Number of genes found for *A. thaliana* and *C. braunii* (in green) are indicated in brackets. APx: Ascorbate peroxidase, Asn: ascorbate, DHA: Dehydroascorbate, DHAR: Dehydroascorbate reductase, GPx: Plant glutathione peroxidase, GR: Glutathione reductase, Grx: Glutaredoxins superfamily, GSH: reduced glutathione, GSSH: oxidized glutathione. Kat: Catalase, MDAR: Monodehydroascorbate reductase, PrxR: Peroxiredoxins family, RBOH: Respiratory burst oxidase homolog also called NADPH oxidase, SOD: Superoxide dismutase, Trx: Thioredoxins, MDA: Monodehydroascorbate, adapted from (Inupakutika et al., 2016).

963 **STAR Methods**

964

965 **CONTACT FOR REAGENT AND RESOURCE SHARING**

966 Further information and requests for resources and reagents should be directed to and will be
967 fulfilled by the lead contact Stefan A. Rensing (stefan.rensing@biologie.uni-marburg.de).

968

969 **EXPERIMENTAL MODEL AND SUBJECT DETAILS**

970 Two strains of *C. braunii* (S276 and S277) were used. The strain S276 was isolated from the
971 thallus, which germinated from the bottom soil of Lake Kasumigaura (Ibaraki, Japan) and was
972 maintained at Kobe University. The unialgal isolation of this strain was achieved as follows.
973 First, collected oospores were surface sterilized for 5 to 8 min in 20% (v/v) NaClO (aq) with
974 0.05% (v/v) Tween20. The sterilized oospores were then transferred into autoclaved soil-water
975 medium for the Charales (SWC-3), containing distilled water and two layers of substrate: a
976 mixture of black soil and river sand on top of a layer of leaf mould. In the present study, strain
977 S277 was newly collected from a pond at Saijo (Ehime, Japan) on October 18, 2011. Newly
978 collected specimens of *C. braunii* were identified based on their *rbc*L DNA sequences. The
979 methods employed for field collection and DNA barcoding followed (Sakayama et al., 2009).
980 The thalli were collected using a handmade anchor. Total DNA was extracted from field-
981 collected samples using the Qiagen DNeasy Plant Mini Kit. Partial *rbc*L DNA sequences were
982 amplified using the primers CHAR-RF-1 (5'-ATGTCACCACAGACAGAAACTAA-3') and
983 CHAR-RR-4 (5'-GCTCCTGGAGCATTTCCCCAAG-3'). PCR conditions were 95 °C for
984 5min; 32 cycles at 95 °C for 40s, 55 °C for 40s, and 72 °C for 1.5min; and 72 °C for 7 min
985 using Ex Taq (Takara Bio). PCR products were sequenced using the primers CHAR-RF-1,
986 CHAR-RR-4, CHAR-RF-2 (5'-GAGCTGTATATGAATGTCTTCG-3') and CHAR-RR-3 (5'-
987 GTTTCTGCTTGAGATTTATA-3'). The sequences obtained were aligned with published *rbc*L
988 DNA sequences of the genus *Chara* downloaded from GenBank. Sequence alignment was
989 performed using MUSCLE (Edgar, 2004) with default options. The aligned dataset of the *rbc*L
990 DNA sequences was subjected to the Neighbour-Joining (NJ) method with Jukes-Cantor
991 distances and 1,000 bootstrap replicates, using MEGA 6.0. Based on NJ trees, field-collected
992 samples were identified at the species level. The unialgal culture of S277 was established
993 following the same procedure as outlined for S276. The pressed specimens of S276 and S277
994 (TNS-AL 209137 and 209138) were deposited at the Herbarium, Department of Botany,
995 National Science Museum (TNS), Tsukuba, Japan. Routine culture was essentially performed
996 at 23 °C with a 16-h light: 8-h dark cycle with 24.5 μmol photons $m^{-2}$ $s^{-1}$ illumination provided
997 by fluorescent lamps using soil-water medium for the Charales (SWC-3).

998

999 **METHOD DETAILS**

1000 **DNA extraction**

1001 Thalli of strain S276 were harvested in SWC-3 medium, washed with distilled water, frozen in
1002 liquid nitrogen, and stored at -80 °C until DNA extraction. High molecular weight DNA was
1003 prepared by the CTAB method followed by purification with Qiagen Genomic Tip. The frozen
1004 powder was weighed and poured on 6 volumes of 2X CTAB buffer (2%

hexadecyltrimethylammonium bromide [CTAB], 1.4M NaCl, 100 mM Tris-Cl pH 8, 20 mM EDTA, 1% Polyvinylpyrrolidone, 1% 2-mercaptoethanol) on a hotplate stirrer at 60 °C. After two rounds of Chloroform:IAA 25:1 extraction, the supernatant was mixed with 3 col of CTAB precipitation buffer (1% CTAB, 50 mM Tris-Cl pH 8, 10 mM EDTA). The precipitate was recovered by centrifugation and dissolved in NaCl solution (1 M NaCl, 10 mM Tris-Cl pH 8, 1 mM EDTA), then precipitated with 0.6 vol of 2-propanol. The precipitate was dissolved in TE and further purified with a Qiagen Genomic Tip according to the manufacturer's instruction. The integrity of the DNA was confirmed with pulsed field electrophoresis using CHEF DR-II (Bio-Rad). Alternatively, genomic DNA from harvested thalli was isolated by grinding the flash frozen material, adding 15 mL extraction buffer (100mM Tris, 50mM EDTA, 500mM NaCl, 10mM 2-mercaptoethanol; pH8) and 2 mL 10% SDS, and incubating for 10 min at 65 °C with mild agitation. Subsequently, 5.4 mL 5M potassium acetate were added and incubated 20 min on ice. After centrifugation at 13,000 g for 20 min at 4 °C the DNA is precipitated by adding 14 mL 2-propanol, incubation for 30 min at -20 °C and centrifugation at 13,000 g for 15 min at 4 °C. After the isopropanol precipitation the air dried pellet was dissolved in 700 µl 1x TE buffer (pH 8), 1-3µl RNaseA (10mg/ml) was added and incubated for 10 min at 37 °C. To purify the DNA 600 µl phenol/chloroform 1:1 were added, mixed, centrifuged at 10,000 g for one minute and the aqueous phase extracted. To this phase 600 µl chloroform/isoamylalcohol 24:1 were added, mixed, centrifuged at 10,000 g for one minute and the aqueous phase extracted. To precipitate the DNA 70 µl 3M Na-acetate and 500 µl isopropanol were added, mixed and centrifuged at 10,000 g for ten minutes. The pellet was washed with one ml 70% ethanol, dried and afterwards was dissolved in deionized water. Quality was controlled using Nanodrop, Qubit measurement and agarose gel electrophoresis.

**Chromosome observation**

The thalli with young antheridia were collected within the first hour of the dark period and fixed in ethanol:acetic acid (3:1). Fixed material was stored at 4 °C until used. Chromosome preparations were made using the Feulgen squash method (Fig. S1). Fixed samples were rehydrated by passing through a graded series of ethanols and rinsed gently in distilled water. The samples were treated with 1N HCl for 5 min at room temperature, then treated with 1N HCl for 8 min in a water bath at 60 °C, and rinsed gently in distilled water. Afterwards, the samples were transferred into Schiff's reagent (Merck Millipore) for 60 min at room temperature. After rinsing the samples in distilled water, antheridia were removed from the thallus and dissected to remove the shield cells. The antheridial filaments were placed on a glass slide and covered with a glass cover-slip. Then, they were squashed to spread the cells and observed.

**Genome sequencing and assembly**

Genomic DNA of the uni-algal strain S276 isolated from Lake Kasumigaura (Ibaraki, Japan) was sequenced as the reference genome using Illumina technology and sequences were compared with those of the strain S277 that was isolated from the pond at Ehime (Japan). Approximately 0.25 Gbp of scaffolds were present in only one of the datasets and found to be of bacterial origin. After removal of these prokaryotic sequences, 1.75 Gbp of scaffold data (N50 size of 2.26 Mbp at #234) were obtained, of which 1.43 Gbp were assembled into contigs. This corresponds to ~74% of the *C. braunii* genome as measured by flow cytometry (1.89-1.96

1048  Gbp) and to ~61% of the 2.35 Gbp estimated by k-mer analysis. The plastid and mitochondrial
1049  genome were assembled separately to recover 187,771 and 67,059 bp circular genomes,
1050  respectively.

**Genome sequencing of *C. braunii* strain S276**

1052  A paired-end library with insert size of 250 bp was constructed using an S2 ultrasonicator
1053  (Covaris) and a TruSeq DNA PCR-Free LT Sample Prep Kit (Illumina) according to the
1054  manufacturer's protocols. The products were size-selected on an agarose gel and purified using
1055  the Qiagen MinElute Gel Extraction Kit. Nucleotide sequences were determined for 150 bp
1056  from both ends with an Illumina HiSeq 2500. Sixteen Mate-pair libraries were constructed using
1057  a Nextera Mate-pair library construction kit with standard and modified input DNA of 5.6, 8,
1058  and 20 μg in the reaction. The first set, four libraries were constructed using the standard
1059  protocol, a gel-free method starting with 1 μg DNA (one library), and gel-excision starting with
1060  4 μg DNA (three libraries). In the Gel-free protocol tagmented DNA was purified with AMPure
1061  XP resulting in a broad size with a peak at 2.7 kbp. In the Gel (+) protocol, the size range was
1062  3-5 kbp, 5-8 kbp, and larger, resulting in a peak of 4.5, 5.8, and 9 kbp, respectively, as measured
1063  with a Bioanalyzer after purification with a Zymoclean Large Fragment DNA Recovery Kit.
1064  After circularization, fragmentation with Covaris S2, end-repair, A-tailing and adapter ligation,
1065  gel-free and 4.5 kbp library were amplified for 10 cycles, whereas 5.8 kbp and 9 kbp libraries
1066  were amplified for 15 cycles. After purification and quantification, the libraries were further
1067  subjected to 8, 6, 6, and 8 cycles of PCR, for gel-free, 4.5, 5.8 and 9 kbp libraries, respectively
1068  (Table S1A).

1069  In the second set, two libraries were constructed using 20 μg DNA instead of the standard 4 μg
1070  DNA to obtain larger fragment size distribution after tagmentation. In this sample, though the
1071  large molecules were not well separated on the agarose gel, three fractions, thick band at high
1072  molecular weight above all marker bands, below the band to 12 kbp, and a 8-12kbp fraction
1073  were recovered. The size of the recovered DNA could not be measured accurately using a
1074  Bioanalyzer, though the peak was around the 17kbp marker. The final amplification was done
1075  for 21 cycles and additional 8 cycles. The lowest 8-12 kbp fraction did not amplify well and
1076  was not used in further analysis.

1077  In the third set, five libraries were constructed using 5.6 μg of starting DNA (1.4-fold of
1078  standard) and an additional five libraries using 8.0 μg of starting DNA (2-fold of standard);
1079  pulsed field electrophoresis on a CHEF-DRII (Bio-Rad) was used for the separation after the
1080  tagmentation. The electrophoretic conditions were 6 V/cm, 11 hours, switch time 1-6 s, on 1%
1081  agarose gel, in 0.5 X TBE buffer. The gel was stained with SYBR Gold and the gel slices were
1082  recovered in five fractions each. The lower limit of each slice was 5.0, 7.5, 10.0, 15.0, and 23.5
1083  kbp. After purification, the DNA was immediately subjected to circularization without
1084  measuring its size. The final amplification was conducted for 15 cycles. Of these (Table S1A),
1085  15 had good insert size distribution when mapped to a preliminary version of the assembly, but
1086  one (S276MP3 xk) had not and thus excluded for further analysis.

1087  Another two mate pair libraries were constructed by GATC (3-4 kbp fragment size) and
1088  sequenced on an Illumina HiSeq 2000. One library was constructed using Crelox with an insert

1089 size of 3 kbp. DNA was fragmented using the Covaris S2 AFA instrument and sequencing was
1090 performed on an Illumina HiSeq 2000 at 2 x 100 bp.

1091 **K-mer frequency analysis**

1092 *K*-mer frequency with *k* = 25 in the paired end reads were counted with JELLYFISH (Marcais
1093 and Kingsford, 2011), applying the min-quality=20 option. A clear peak at 51 was observed
1094 with a valley at 16 (Fig. S2A). The peak at 51 was interpreted as the single copy genomic
1095 sequence and those less than 16 were mostly k-mers containing sequencing errors. The
1096 cumulative *k*-mer count from 16 upto 10,000 (which was the default upper limit of JELLYFISH)
1097 divided by 51 suggested the genome size be 2.355 Gbp. Note that this number includes *k*-mers
1098 derived from organellar and bacterial sequences and supposed to be overestimate for the nuclear
1099 genome size. With the peak at 51, the amount of paired-end reads are supposed to be sufficient
1100 for the assembly. The region from 16 to 80 as the putative single copy region comprised 0.95
1101 Gbp.

1102 **Assembly**

1103 The raw sequences were assembled with ALLPATHS-LG (Gnerre et al., 2011). Initially the
1104 assembly started with R48517 on a machine having 768 GB of memory and 32 CPU cores.
1105 After running a month this process stopped at UnipathPatcher phase. Continuation was tried
1106 with the settings: PATCH_UNIPATHS=False FIX_LOCAL=False
1107 PATCH_SCAFFOLDS=False FIX_SOME_INDELS=False; unfortunately this failed again.
1108 The run directory was copied to a machine having 2 TB of memory and 80 cores and the
1109 assembly was continued with R48777 and completed after another twenty days (with 48
1110 slots=threads), with reported peak memory usage of 1,756 GB. The assembly resulted in 28,091
1111 scaffolds with a total length of 1.99 Gbp, comprised of 250,979 contigs with a total length of
1112 1.65 Gbp. The library information is summarized in Table S1B.

1113 **Genome sequencing of *C. braunii* strain S277**

1114 Thalli of strain S277 were harvested in SWC-3 medium, washed with distilled water, frozen in
1115 liquid nitrogen, and stored at -80 °C until DNA extraction. Total DNA was extracted as
1116 described above. A paired end library was constructed using a TruSeq DNA PCR-free library
1117 preparation kit (Illumina) and sequenced with HiSEQ (DRA accession: DRR054048). 1.1 μg
1118 of DNA was fragmented with Covaris S2, using micro tube, duty cycle 10%, intensity 4, 200
1119 cycles/burst and total time of 80 s. The fragments were size selected using a bead-based method
1120 following the 350-bp protocol.

1121 **PacBio sequencing of fosmid clones for quality control**

1122 *C. braunii* S276 genomic DNA was cloned into the pNGS fosmid vector using the aNxSeq 40
1123 kbp Mate-Pair Cloning Kit (Lucigen). Six fosmid clones with verified end sequence and one
1124 96 well plate of undetermined clones were pooled and shotgun sequenced on a PacBio SMRT
1125 cell (608 Mbp, 63,768 reads post-filtering). The resulting reads were assembled into contigs
1126 using HGAP (Chin et al., 2013) in smrtanalysis (PacificBiosciences). The contig sequences
1127 were further polished with two rounds of Quiver. Bacterial contamination was removed using
1128 MEGAN, and comparative mapping of S276 and S277 reads, resulting in 22 probable *C.*
1129 *braunii* contigs. All but one of those could be BLAST-mapped to the assembly. One clone

appeared to be chimeric based on mapping Illumina mate-pair library data on the clone. Of the remaining 20, 14 were mapping to single scaffolds, the other 6 to 2-4 scaffolds. 10 of the 22 contigs were found to map with >=95% identity and >= 90% coverage to the assembly, the remaining 12 did not meet these parameters, probably due to assembly gaps. In summary, 45% of the assembled fosmid clones had high quality representations in the assembly, and 91% could be mapped, demonstrating the good quality of the assembly.

**Distinction of bacterial sequences**

Paired end sequences of S276 and S277 were mapped to the assembly with bwa mem (Burrows-Wheeler Aligner) (Li and Durbin, 2010) and the number of mapped sequences were counted on each scaffold (Fig. S2). Number of tags of both samples on each scaffold was plotted and we found two groups. The two groups were separated by a line in which S277 had 1/100 of S276 tags (Fig. S2B). The GC content of each scaffold was calculated and compared between the two groups. The group showing less tags in S277 had a higher GC content distribution (Fig. S2C). Thus, these scaffolds were presumed to be derived of different organisms, which were probably bacteria that survived autoclaving. In addition, scaffold_64 was found to be of bacterial origin in manual inspection during gene prediction. Further, the genomic scaffolds were split into 1 kbp fragments. Using tera-BLASTn 9.0.0 on DeCypher 9.0.0.25 (http://www.timelogic.com/catalog/757/tera-blast) each fragment was BLASTed against the NCBI nt database. The BLAST output was analysed by MEGAN 6 (Huson et al., 2016) and bacterial hits assigned to the 1 kbp fragments. All scaffolds containing more than 50 % of bacterial hit fragments were extracted. If no non-bacterial hits were contained on the scaffold and the bit score of the bacterial contamination exceeded 50 per hit the scaffold was removed as contamination. This affected 153 scaffolds with a total length of 312 kbp (Table S5), containing 120 gene models (marked in Table S4). Thus, 11,655 scaffolds totaling 1,751,225,565 bp, comprised of 234,221 contigs totaling 1,429,911,168 bp were recovered as representing the *C. braunii* nuclear genome. N50 scaffold size, and N50 contig size were 2,261,426 bp (at #234) and 10,124 bp (at #41,610), respectively.

**Microbiome analysis**

The diversity of microorganisms is expected to be low due to lab-culturing conditions and DNA sequence extraction protocols. To isolate the microorganisms remaining in the bulk of data, we mapped reads to the eukaryotic genome and only analyzed reads left unmapped. Given that S276 and S277 were reared at different geographical locations, analyzes were done on both sets separately. The two separate sets of remaining reads were assembled into contigs and analyzed from a meta-genomics point of view. Two separate assemblies have been generated using CLC-assembly cell using the larger word-size (kmer) of 50 nt to force more specificity (CLC bio, Aarhus, Denmark). These assemblies resulted in respectively 322685 contigs with a total size of 76.7 Mbp (N50 242 bp, max size 167358 bp, min size 100 bp) and 325720 contigs with a total size of 90.1 Mbp (N50 373 bp, max size 172440 bp, min size 100 bp). The obtained contigs represent a mixture of microorganisms that where clustered using CONCOCT (Alneberg et al., 2014) according to the manual, using BEDtools (Quinlan, 2014), Picard-tools and R, to create and format the needed input files. Several runs were done, aiming at providing the minimal number of differentiated clusters. In some cases large clusters were isolated and submitted again for a new round of clustering. The clusters (or bins) were calculated based on read coverage

and sequence tetramer composition of the contigs following an iterative fitting of mixture-of-Gaussian models on the available data; each group is supposed to represent an organism that was further characterized to establish the species. Taxonomic assignment of the bins was performed using a similarity-based labeling of the fragments with MEGAN5. A first assessment of the quality and completeness of the bins was done by monitoring the presence of 36 COG single copy genes. 16S rRNA genes were isolated from the sequences using online RNAmmer 1.2 Server (Lagesen et al., 2007) and provided to SINA Alignment Service within Silva database for classification (Pruesse et al., 2012). Not all clusters could be identified up to species level, but for those for which we could find a reference genome, we show also a level of completeness by comparing to the respective reference genomes using nucmer from the MUMmer (Delcher et al., 1999) v3.23  package (Table S1T, S1U).

**Transcriptome sequencing**

Thalli of strain S276 were harvested in SWC-3 medium under controlled laboratory conditions at 23 °C with a 16-h light: 8-h dark cycle with 24.5 μmol photons $m^{-2} s^{-1}$ illumination provided by fluorescent lamps. Two and seven different samples, for full-length cDNA and RNA-seq analyses, respectively, were collected, frozen in liquid nitrogen, and stored at −80 °C until further processing. Frozen samples were ground in liquid nitrogen. Total RNAs were then extracted with ISOGEN (Nippon Gene, Tokyo, Japan), and purified using the Qiagen RNeasy Plant Mini Kit. For the extraction of total RNA in oospores and rhizoids, Fruit-mate (Takara Bio, Shiga, Japan) was used prior to the extraction by ISOGEN. Full-length cDNA libraries were constructed using the oligo-capping method. Total RNA was treated with bacterial alkaline phosphatase (BAP; Takara) at 37°C for 40 min with RNasin (Promega). After extraction with phenol:chloroform (1:1) twice and ethanol precipitation, the RNA was treated with  tobacco acid pyrophosphatase (TAP; in house purified) with RNasin at 37°C for 45 min. The BAP-TAP treated  RNA were ligated with 5′-oligo (5′-AGC AUC GAG UCG GCC UUG UUG GCC UAC UGG-3′) using T4 RNA ligase (Takara).The first strand cDNAs were amplified using 5′ (5′-AGC ATC GAG TCG GCC TTG TTG-3′) and 3′ (5′-GCG GCT GAA GAC GGC CTA TGT-3′) PCR primers. The amplified cDNAs were digested with SfiI and cloned into DraIII-digested pME18S-FL3-3 (AB009864).  Clones were picked and sequenced with ABI sequencers at National Institute of Genetics, Japan. After filtering for vector, synthetic oligonucleotides, and low-quality sequences 73,388 reads were left in total (Table S1D). RNA-seq libraries were constructed via the Illumina mRNA-Seq Sample Prep Kit using RNA extracted from various tissues (Table S1E). 76 or 101 bp paired end sequencing was performed on an Illumina HiSEQ 2000. Additionally, a late reproductive phase thalli (harvested 2-3 weeks after appearing of the gametangial primordia) library was constructed as RNA-ligation based stranded library using the combined method of mRNA-Seq Sample Prep Kit and Small RNA Sample Preparation Kit (Illumina), following the manufacturer's instructions. This library was sequenced by 76 bp single end sequencing performed on a GAIIx (Illumina).

**Quantitative transcriptome comparison of antheridia, oogonia, and zygotes**

Antheridia and oogonia were hand-dissected in Qiagen RNA*later* from *C. braunii* thalli (strain S276) grown under a 14:10 hours light:dark cycle at 22 °C. Zygotes were collected once detached from mother plants grown in identical conditions. Samples were flash frozen in liquid nitrogen then kept at -80 °C until further processing. Approximately 20 mg of starting material

was ground in liquid nitrogen then total RNA was extracted using Ambion mirVana kit following manufacturer's recommendations. DNA was digested from RNA extracts using Promega RQ1 DNase and RNA was cleaned using a Qiagen RNeasy MinElute Cleanup Kit. RNA was then amplified using an Ovation RNA-Seq System V2 (NuGEN) amplification kit following manufacturer's protocol. Final amplified cDNAs were cleaned using the Qiagen PCR cleanup kit. Three biological replicates were obtained for antheridia, oogonia and zygotes. One sample containing vegetative and reproductive tissues was similarly prepared, except for the amplification step. 20 µg of RNA from each replicate was paired-end sequenced on an Illumina HiSeq 2000 platform at the Beijing Genomics Institute in China; at least 2 x 10 million reads were obtained per sample. Reads were processed to remove low quality sequences, PCR adapters, foreign sequences introduced by the amplification procedure and any detectable bias using Trimmomatic v0.36 (Bolger et al., 2014) and Perl scripts. Transcript were inferred from the reads pooled and aligned to the *C. braunii* genome sequence using Tophat v2.1.0 (Kim et al., 2013) and Cufflinks v2.0.2 (Trapnell et al., 2010). Both programs were given the *C. braunii* genomic structure as a guide. A custom Perl script was then used to clean Cufflinks predictions from spurious gene fusions and other detectable problems. Unaligned reads were further normalised, assembled and scaffolded into transcripts. Both reference guided and *de novo* assemblies were merged. Coding sequences were predicted, and sequence annotation and GO terms were obtained from transcripts using a pipeline based on BLAST v2.2.29 (Altschul et al., 1997) and TransDecoder v2.0.1 (Haas et al., 2013). A summary of assembly and read mapping statistics is presented in Table S1W. Read counts were obtained by mapping reads onto the inferred transcriptome with RSEM v1.2.11 (Li and Dewey, 2011). Differential expression was tested between zygotes and oogonia samples and between oogonia and antheridia samples and was conducted in R using DESeq2 v1.14.1 (Love et al., 2014). Genes were considered differentially expressed between two conditions with an adjusted p-value < 0.01 and a log2 fold-change (logFC) > 2. Differentially expressed genes are listed in Table S2. GO terms enrichment analysis was conducted in R using topGO v2.22.0. Enriched GO terms and associated genes are listed in Table S3. Heatmaps were generated using R and the package pheatmap v1.0.8. Visualization of the GO terms was implemented using word clouds via the http://www.wordle.net application. The weight of the given terms was defined as the -log10(q-values) and the color scheme used for the visualization was red for down-regulated GO terms and green for those up-regulated. See Table S2 for DEGs and Table S3 for GO analyses.

**Identification of repeat sequences with RepeatModeler/RepeatMasker**

A species-specific repeat model was constructed using RepeatModeler Version open-1.0.7 with ncbi engine. Repeats were identified using RepeatMasker version open-4.0.5 with Search Engine: NCBI/RMBLAST [2.2.27+] and RepeatMaskerLib.embl (Complete Database: 20140131), resulting in masking 46% of the genome. The breakdown is shown in Table S1F.

**Gene prediction**

High throughput cDNA sequencing (RNA-seq) was conducted on several libraries representing vegetative and reproductive stages, including thallus, gametangia and zygotes. These data were used together with full-length cDNA sequences to annotate the genome with AUGUSTUS. 35,445 putatively protein-coding genes were identified, of which 63% could be annotated using similarity-based approaches. A total of 13,331 gene models overlap to at least 50% with TE

evidence and thus might not represent canonical protein-coding genes, bringing the number of protein-encoding genes down to 23,546. In total, the expression of 12,388 of those (53%) was supported by RNA-seq data (Table S4). Reciprocal best BLAST (Altschul et al., 1997) hit analysis of the *C. braunii* protein set revealed a high percentage presence of core gene sets: 96.43% of eukaryotic benchmarking universal single-copy orthologs (BUSCO, (Simao et al., 2015)), 98.65% CEGMA core eukaryotic genes (Parra et al., 2007), and 93.96% core gene families for green plants (Van Bel et al., 2012).

Gene prediction with Augustus (Keller et al., 2011) was performed following https://computationalbiologysite.wordpress.com/2013/07/25/incorporating-rnaseq-tophat-to-augustus/. Initial models were created based on the CEGMA output. RNA-seq data was mapped to the RepeatMasker masked *C. braunii* genome. Each accepted_hits.bam was sorted and processed with filterBam --uniq (--paired for paired data). Evidence of introns was extracted using bam2hints –intronsonly to obtain intron_hints.gff. The first round of Augustus was run with this as hints. An exon-exon junction database was constructed based on this output and bowtie was used to map the reads to the junctions. These mappings were further merged to the first intron hints and the second round of augustus was run. Gene prediction at this phase was manually investigated and confirmed genes on scaffold_0 and scaffold_2 were chosen and adjusted for the 5' and 3' ends of UTR based on RNA-seq mapping on Web Apollo (Lee et al., 2013). Thus, 120 manually inspected gene models were used to retrain Augustus. Construction of exon-part hints through wig file were performed according to http://augustus.gobics.de/binaries/readme.rnaseq.html. For the stranded RNA-seq data, forward and reverse mapped reads were separated with samtools and assigned the strand accordingly. Repeat hints were prepared by processing the gff file created by the RepeatMasker with "sed –e s/similarity/nonexonpart/ -e 's/Target.*/src=RM/'". Amino acid sequence of *A. thaliana* (TAIR10_pep_20110103_representative_gene_model_updated) and *P. patens* (P.patens.V6_filtered_cosmoss_proteins.fas) were mapped to the genome using exonerate and converted as hint data. The full-length EST sequences were mapped using blat (Kent, 2002) with -minIdentity=92 -extendThroughN parameters and converted to EST hints. All these hints were merged to a single hints file and the final run of Augustus was run with --gff3=on --UTR=on --alternatives-from-evidence=true --allow_hinted_splicesites=atac with a merged hints file. The output was collected and gene models predicted on the 11,808 scaffolds that were treated as *C. braunii* genome. Thus, we obtained 36,877 transcripts from 35,883 loci. For annotation see Table S4.

**Assembly of organellar genomes**

Organellar genomes were assembled using NOVOPlasty (Dierckxsens et al., 2017) v2.5.3. For chloroplast genome, two lanes of paired end data were processed using the *Chara vulgaris* chloroplast genome (NC_008097.1) as seed. This resulted in 4 possible reconstructions, two in 187 kbp and the remaining two in 200 kbp, i.e. contig arrangement 01+02+03+04+06, 01+04+05, 01+02+03+04+05, or 01+04+06. The differences are on whether 02 and 03 are inserted and whether the end is 05 or 06. 02 and 03 is contained in 01 and seems to represent an inverted repeat region and insertion of them would be excess. The 05 and 06 contain 27,447-bp common sequence, which is the small single copy region. Given there are about equal number of molecules that is flipped at the inverted repeat region, both reconstructions are

1302 equally valid and one is arbitrarily chosen. The mitochondrial genome was assembled using the
1303 *C. vulgaris* mitochondrial genome (NC_005255.1) as seed input and specifying the chloroplast
1304 genome obtained as above. This resulted in a single circularized assembly of 67,059 bp, which
1305 is very close to 67,737 of the *C. vulgaris* mitochondrial genome.

**Repeat/TE annotation**

1307 Repetitive elements collectively contribute approximately 1.1 Gbp of the genome assembly.
1308 This estimate is probably low, given that highly similar repeats are challenging to assemble and
1309 that there is ~0.5 Gbp size difference between the ungapped (1.43 Gbp) assembly and C-value
1310 estimates (1.9 Gbp). Transposable elements (TEs) and unclassified repeats are abundant (61%
1311 and 37% of repeat annotation, respectively), with Gypsy-type LTR retrotransposons
1312 representing 24% (343 Mbp) of the ungapped assembly (Table S1G).

1313 We have used the REPET package v2.4 to perform *de novo* identification, classification and
1314 annotation of repetitive elements in the *C. braunii* assembly as decribed in (Jouffroy et al.,
1315 2016). We first launched the TEdenovo pipeline on a sub-genome comprising contigs of size
1316 above 20 kb and representing a total of 362 Mb (12,655 contigs). We used default settings
1317 except that the minimum number of copies per group was set to 5 (minNbSeqPerGroup: 5),
1318 resulting in a library of 3,140 consensus sequences. This library was subsequently filtered by
1319 using the TEannot pipeline against the whole assembly and discarding consensus sequences
1320 without a single full length match, resulting in a library of 2,161 sequences. This filtered library
1321 was used to annotate the whole genome assembly using the TEannot pipeline. Threshold
1322 annotation scores were determined for each consensus as the 99th percentile of the scores
1323 obtained against a randomized sequence (whole genome reversed, not complemented and
1324 masked with TRF). Consensus sequences were then classified using the features detected with
1325 PASTEC followed by semi-manual curation. In addition to the HMM comparison against
1326 PFAM implemented in PASTEC, we have also used RPS-BLAST (-F T -e 1e-2) to search for
1327 more remote homologies against a library of CDD domains identified in the repbase library.

1328 Several unclassified consensus sequences have been classified in putative retrotransposons
1329 because they contain at least one of the following domains: cd00024 Chromatin organization
1330 modifier, cd00303 Retropepsins, cd01650 RT nLTR, cd01651 RT G2 intron, cd05482
1331 Retropepsins, cd06095 Retropepsin, cd06222 RNase H, pfam00385 Chromo, pfam00552
1332 Integrase, pfam00665 Integrase, pfam02093 Gag P30, pfam03708 Avian retrovirus envelope
1333 protein, pfam03732 Retrotransposon gag protein, pfam07727 Reverse transcriptase,
1334 pfam10536 Plant mobile domain, pfam13966 zinc-binding in reverse transcriptase, pfam13975
1335 gag-polyprotein putative aspartyl protease, pfam13976 GAG-pre-integrase domain, and
1336 smart00298 Chromatin organization modifier domain.

1337 Based on the REPET results, percentage overlap of protein coding gene models with TEs was
1338 assessed and added to Table S4. Gene models overlapping to 100% with TE evidence are
1339 considered true TE genes, while those overlapping to at least 50% (but less than 100%) might
1340 be protein-coding genes present in TE regions, or might encode TE-based proteins. All genes
1341 were kept in the gene catalog so that individual evaluation (e.g. based on the homology-based
1342 annotation) is possible.

**Screening for whole genome duplication events**

1344 To identify whole genome duplication (WGD) events we employed the KeyS software (Rensing
1345 et al., 2007) to obtain Ks (synonymous substitution) distributions of paralogous genes for *C.*
1346 *braunii.* In brief, paralogous genes were defined by a self-BLAST retaining only BLAST hits
1347 that showed at least 50% query and subject coverage and an alignment length according to the
1348 twilight zone *sensu* (Rost, 1999). Gene pairs with a BLAST identity of 98% or higher were
1349 further tested at the nucleic acid level to remove nearly identical sequences using optimal global
1350 alignments and a threshold of 98% identity. For nearly identical gene pairs only the longer
1351 sequence was kept and all gene pairs containing the shorter sequence were discarded (Rensing
1352 et al., 2007). The paralogous genes were further clustered using a minimal connectivity
1353 threshold of 50% (half linkage) and Ks values were calculated at the cluster nodes (representing
1354 duplication events rather than gene pairs) using the maximum likelihood method of CODEML
1355 implemented in PAML v4.7 (Yang, 2007).

1356 The following procedure has been described recently (Lang et al., 2018), please see there for
1357 related citations. Briefly, we employed mixture modeling to find WGD signatures using the
1358 *mclust* v5.1 R package (Scrucca et al., 2016) to fit a mixture model of Gaussian distributions to
1359 the raw Ks and log-transformed Ks distributions. All Ks values $\leq 0.1$ were excluded for analysis
1360 to avoid the incorporation of allelic and/or splice variants and to prevent the fitting of a
1361 component to infinity, while Ks values $> 5.0$ were removed because of Ks saturation. Further,
1362 only WGD signatures were evaluated between the Ks range of 0.235 (12.5 Ma ago) to account
1363 for recently duplicated gene pairs to Ks of 2.0 to account for misleading mixture modeling
1364 above this upper limit. Because model selection criteria used to identify the optimal number of
1365 components in the mixture model are prone to over fitting we also used SiZer and SiCon as
1366 implemented in the *feature* v1.2.13 R package (Duong et al., 2008) to distinguish components
1367 corresponding to WGD features at a bandwidth of 0.0188, 0.047, 0.094 and 0.188
1368 (corresponding to 1 Ma, 2.5 Ma, 5 Ma and 10 Ma ago) and a significance level of 0.05.

1369 Deconvolution of the overlapping distributions that can be derived from paranome-based Ks
1370 values without structural information shows that using mixture model estimation based on log-
1371 transformed Ks values mimics structure-based WGD predictions better than using raw Ks
1372 values, and can predict young WGD signatures and can pin point older WGD signatures (Lang
1373 et al., 2018). Since WGD signature prediction based on paranome-based Ks values can be
1374 misleading and is prone to over prediction we only considered Ks distribution peaks in a range
1375 of 0.235 to 2.0 as possible WGD signatures, thus excluding young paralogs potentially derived
1376 from tandem or segmental duplication and those for which accurate dating cannot be achieved
1377 due to high age (Fig. S3).

1378 **Genome comparison**

1379 *C. braunii* was compared with eight further Viridiplantae genomes. In addition to the genome
1380 length, GC content and the number of annotated genes, the mean intergenic and the mean intron
1381 length were calculated. The intergenic length was performed by extracting the genome regions
1382 not covered by the gff3 annotation file with bedtools complement (Quinlan, 2014) version
1383 2.25.0. The intron length was calculated by extracting the distance between the annotated CDS
1384 regions. Both mean length and the corresponding standard deviation were calculated using awk
1385 (Table S1L). The gene density of the *C. braunii* genome is relatively sparse as compared to e.g.
1386 *A. thaliana*, *O. sativa* (rice) or two algae (*K. nitens* and *Chlamydomonas reinhardtii*), but similar

1387 to other Gbp-sized genomes like *Z. mays* or *H. vulgare* (Fig. 3 and Table S1L); the distance
1388 between genes is comparable to the approximately equal-sized *Z. mays* genome.

**Comparative analysis of gene and transposons in selected plant and algae species**

1390 The genome sequences and annotations of *K. nitens*, *C. reinhardtii*, *A. thaliana*, *M. polymorpha*,
1391 *Oryza sativa*, *P. patens*, *C. braunii*, *Z. mays*, *H. vulgare* were downloaded and processed with
1392 GAG and the genome tools gff3 validator, to obtain consistent annotation files. For each
1393 annotated gene, intronic regions were inferred using the GenomeTools gff3 program. The *K.*
1394 *nitens* annotation file was manually curated for consistency with the other annotations and the
1395 GFF3 data standard.

1396 Subsequently, intact full-length long terminal repeat transposon elements (LTREs) were
1397 predicted using the GenomeTools LTRharvest and LTRdigest software (Steinbiss et al., 2009)
1398 utilizing a set of TE-associated PFAM domains and a compilation of eukaryotic tRNAs. The
1399 pipeline was implemented as a BASH/PBS shell script (run_LTR_harvest_digest.sh). The
1400 resulting set of candidate LTREs was filtered to contain 2 LTRs, >=1 protein domain match and
1401 2 target site duplications. These filtered elements were considered to represent intact full-length
1402 LTREs whose nucleotide sequences were extracted and searched against the genome using
1403 Vmatch requiring >=80% sequence identity and 100 bp alignment length. Depending on the
1404 repeat content and genome size, genomes where either split at gap boundaries into preferably
1405 100 Mbp stretches using the UCSC toolkit faSplit (A: Snakemake workflow: split_approach),
1406 or directly processed as a whole FASTA file (B: Snakemake workflow: vmatch_mask) (Köster
1407 and Rahmann, 2012). Resulting putative LTRE fragments were merged into non-redundant,
1408 non-overlapping regions using the reduce function implemented in the R/Bioconductor package
1409 GenomicRanges (A) (Lawrence et al., 2013) or the bedtools merge program (B).

1410 Helitrons were predicted using the HelitronScanner software using the parameters reported for
1411 element inference and copy number prediction in plant genomes reported in the initial
1412 manuscript (Xiong et al., 2014). Additional fragments were inferred by matching 50 bp from
1413 the 3' terminus of each full-length helitrons against the respective genome utilizing Vmatch
1414 (Abouelhoda et al., 2004) following the same approach as described for LTREs. Resulting
1415 matches and full-length helitrons were merged into non-redundant, non-overlapping regions
1416 using the bedtools merge program. The pipeline was implemented in the Snakemake workflow
1417 in folder helitrons/.

1418 Gene-to-gene, gene-to-LTRE, LTRE-to-gene and LTRE-to-LTRE distances were inferred using
1419 an R script utilizing the distanceToNearest function from the R/Bioconductor GenomicRanges
1420 package (get_distances.R/get_distances.sh). Subsequent data analysis and plotting was carried
1421 out and documented in the R Jupyter Notebooks: folder analysis/: analyseWindows.ipynb,
1422 Distances.ipynb, Introns.ipynb, Lengths.ipynb. All described, generated materials and software
1423 needed to reproduce this analysis are available from the accompanying Mendeley Data
1424 repository (doi:10.17632/9hzzf9m4kh.1), arranged as an archive
1425 ("ComparativeTE_and_genes.Lang.tar.gz") that contains input, output and scripts.

**In-depth analyses of specific gene families**

*Cell wall biosynthesis*

1428 Glycosyltransferases in the *C. braunii* genome assembly were initially identified via BLAST,
1429 using the Carbohydrate Acting enZYme database (CAZY) as of 2016-06-01 as query and a cut-
1430 off value of $10^{-25}$. The sequences were manually verified by alignment with known cell wall
1431 biosynthetic glucosyltransferases and deposited in Table S1H. Phylogenetic trees were
1432 constructed using Phylogeny.fr with standard settings, starting with muscle alignment, curation
1433 of alignment by deletion of positions with gaps, and finally PhyML maximum likehood tree
1434 construction (Guindon et al., 2010). The phylogenetic trees (Data S1A, B) were statistically
1435 supported by approximate likelihood-ratio tests using default settings and values between 0 and
1436 1 were obtained, as with bootstrap values. Approximate likelihood-ratio-test (aLRT) values
1437 were included when values were under 0.7 where *C. braunii* sequences are present.

## *Cell division*

1439 In order to compare the mode of cell division of algae and land plants we compiled a list of 221
1440 *Arabidopsis* genes involved in cytokinesis (Table S1C), focusing on genes required for
1441 phragmoplast and PPB function. With these 221 *A. thaliana* proteins, a BLASTp (version
1442 2.6.0+) search was performed against published plant and algal genomic/transcriptomic
1443 datasets (key resource table), including *C. braunii* and *K. nitens*. The e-value cutoff was set to
1444 1E-4 and the number of database sequences to show alignment for was set to 3,000. The BLAST
1445 result was filtered according to (Rost, 1999) to keep homologous sequences only. Mutiple
1446 sequence alignments for phylogenetic trees of protein families were conducted using MAFFT
1447 (Katoh and Standley, 2013) in the automatic mode, and manually curated. The best fitting
1448 evolutionary model based was determined using ProtTest (Darriba et al., 2011) and applied in
1449 Bayesian phylogenetic inference using MrBayes (Ronquist et al., 2012) with two hot and two
1450 cold chains (Data S1Q-U) until the standard deviation of split frequencies dropped below 0.01
1451 or for 6 mio generations (actin and cyclin).

1452 Using the amplification score that shows potential gene expansion between *K. nitens* and *C.*
1453 *braunii* (Table S1C) we performed phylogenetic analyses as outlined above and found cyclin
1454 genes to be amplified in *C. braunii*, suggesting a more intricate regulation of the cell cycle as
1455 compared to *K. nitens*. While there is a single A1-type cyclin in both algae, the *C. braunii*
1456 genome encodes three B1-type cyclins (like *A. thaliana*), whereas *K. nitens* encodes only one
1457 (Table S1C, Data S1Q). We also found evidence that membrane trafficking is more elaborate;
1458 there are three genes coding for EXOCYST 70A in *A. thaliana*, two in *C. braunii* (and in the
1459 transcriptomes of several Zygnematophyceae), and a single gene in *K. nitens* (as in *Mesostigma*
1460 *viride* and Chlorophyta; Data S1R). With regard to the SNARE complex, we find that the *A.*
1461 *thaliana* NOVEL PLANT SNARE (NPSN) 11/12/13 clade contains two *C. braunii* (and two
1462 *Nitella mirabilis*) and a single *K. nitens* (and *M. viride*) protein (Data S1S).

## *Phytohormones: ETH*

1464 For the identification of putative homologs for ETH biosynthesis and signaling genes,
1465 BLASTp/tBLASTn searches were carried out against the *C. braunii* gene models and genome
1466 assembly using representative *A. thaliana* protein sequences as queries [ACS1 (AT3G61510),
1467 ACO1 (AT2G19590), ETR1 (AT1G66340), CTR1 (AT5G03730), EIN2 (AT5G03280), EIN3
1468 (AT3G20770); Table S1J]. Translated sequences of putative ETH biosynthesis/signaling genes
1469 from *C. braunii* were then used as queries in reciprocal BLASTp searches to the *A. thaliana*

protein database. Multiple ACO homologs were found in the *C. braunii* genome, however, the reciprocal BLASTp search suggests that these homologs are likely to be other oxidases. The other candidate *C. braunii* ETH biosynthesis/signaling protein sequences were manually verified and screened for essential protein domains [ACS (PR00753), ETR/ERS (ETH Binding Domain), CTR1 (PF14381 and CD13999), EIN3 (PF04873 and C-terminal Signaling Domain), EBF (IPR001810)]. An additional search with BLASTP 2.8.0+ using the representative *A. thaliana* proteins as queries and the putative homologs as the subjects was performed.

### *Phytohormones: ABA*

For the identification of putative homologs for ABA biosynthesis and signaling genes, BLASTn/BLASTp searches were carried out against the *C. braunii* gene models and genome assembly using representative *A. thaliana* genomic/protein sequences as queries (Table S1J). An additional search with BLASTP 2.8.0+ using the representative *A. thaliana* proteins as queries and the putative homologs as the subjects was performed. The obtained *C. braunii* protein sequences were manually verified and screened for essential protein domains [PSY (PF00494), PDS (PF01593), GTG1 (PF12537), SnRK/CPK (PF00069)].

### *Phytohormones: SL*

For the identification of putative homologs for SL biosynthesis and signaling genes, BLASTn/BLASTp searches were carried out against the *C. braunii* gene models and genome assembly using representative *A. thaliana* genomic/protein sequences as queries (Table S1J). An additional search with BLASTP 2.8.0+ using the representative *A. thaliana* proteins as queries and the putative homologs as the subjects was performed. The obtained *C. braunii* protein sequences were manually verified and screened for essential protein domains [CCD (PF03055)].

### *Phytohormones: Jasmonates (JA), Salicylates (SA), Gibberellins (GA), Brassinosteroids (BR)*

For the identification of putative homologs for JA, SA, GA and BR biosynthesis and signaling genes, BLASTn/BLASTp searches were carried out against the *C. braunii* gene models and genome assembly using representative *A. thaliana* genomic/protein sequences as queries (Table S1J). Canonical (land-plant like) signaling pathways for JA, SA, GA and BR have been shown to have arisen in land plants [JA - (Han, 2017); SA - (Wang et al., 2015)], vascular plants [GA - (Gao et al., 2008; Wang et al., 2015)] and seed plants [BR - (Vriet et al., 2015)] respectively. Consistent with these findings, none of the genes encoding steps in the biosynthesis or signaling pathways for GA, JA, SA or BR appear to be present in the *C. braunii* genome (Table S1J). However, JA was found in *C. australis* (Beilby et al., 2015), JA and SA were detected in *K. nitens* (Hori et al., 2014), and GA was detected in *Chara tomentosa*, suggesting a different synthesis than known in land plants as in the case of AUX and ABA (Table 1, Fig. 4).

### *Phytohormones: AUX transport*

For the identification of putative homologs for AUX transporter genes, tBLASTn/BLASTp searches were carried out against the *C. braunii* gene models and genome assembly using representative *A. thaliana* genomic/protein sequences as queries (Table S1J and S11).

1510 Predicted coding sequences of PIN proteins were manually aligned with representative PIN
1511 sequences from previously published alignments, PIN sequences from charophyte algae were
1512 obtained from the NCBI database. The PIN sequence of *K. nitens* (GAQ81096.1) originated
1513 from the complete genome assembly, other algal sequences were obtained from the SRA
1514 database (Leinonen et al., 2011) of individual sequencing project by using the BLASTn
1515 algorithm, using the sequence from *K. nitens* as a query. The resulting hits were assembled with
1516 CAP3 (Huang and Madan, 1999) and repeatedly BLASTed against respective SRA databases
1517 to increase sequence length. Maximum-likelihood phylogenetic analysis was performed in
1518 MEGA 7.0 software using amino acid representation of highly conserved N- and C-terminal
1519 part of PIN sequence, LG+G+I substitution model and 500 bootstrap replicates (Data S1C, D).

### *Phytohormones: AUX signaling*

1521 For charophyte algae, mRNA sequences were downloaded and protein sequences were
1522 predicted with ESTScan v3.0.3 (Iseli et al., 1999) using the *A. thaliana* matrix [-M
1523 Arabidopsis_thaliana.smat]. Subsequently all proteins were screened with *hmmsearch* of the
1524 HMMer software suite (v3.1b2) for the abundance of the PFAM v30.0 domains: Auxin_resp
1525 (PF06507), AUX_IAA (PF02309), B3 (PF02362), F-box (PF00646) and F-box-like (PF12937)
1526 using either the gathering threshold [--cut_ga] option or an E-value of 0.1 for the complete
1527 sequence [-E 0.1] and an E-value of 0.1 for the domain [--domE 0.1] to account for possible
1528 sampling bias and cutoff bias of the curated PFAM model.

1529 The obtained results were used to classify the proteins into possible AUX gene families: ARFs
1530 [mandatory domains: Auxin_resp + B3; optional: AUX_IAA], Aux/IAA [mandatory:
1531 AUX_IAA - Auxin_resp] and TIR1/AFB [mandatory: F-box or F-box-like]. For the AUX gene
1532 familiy TIR1/AFB an additional BLAST search with BLAST+ (v2.5.0) [-matrix BLOSUM45
1533 -evalue 1e-5] using representative *A. thaliana* genes as queries [AT3G62980.1 (TIR1),
1534 AT4G03190.1 (AFB1), AT3G26810.1 (AFB2), AT1G12820.1 (AFB3), AT4G24390.2 (AFB4),
1535 AT5G49980.1 (AFB5)] and the domain containing proteins as the subjects was performed. Only
1536 BLAST hits with a query coverage (alignment length / query length) of at least 50% and a
1537 minimal protein identity according to formula (2) of (Rost, 1999) were retained as possible
1538 AUX gene family candidates. Maximum-likelihood phylogenetic analysis for each AUX gene
1539 family was performed on manual curated multiple sequence alignments obtained via MAFFT
1540 (v7.305b) and the E-INS-i algorithm. *IQ-TREE* (Nguyen et al., 2015) v1.5.3 was applied using
1541 the standard non-parametric bootstrap option with 1,000 replicates and the best model selected
1542 by *IQ-TREE* (Table S1K, Data S1E-G).

### *Phytohormones: AUX,* in silico *modeling of* C. braunii *LRR FBPs.*

1544 Leucine-RichRepeat (LRR)-containing F-Box Proteins (FBPs) from *C. braunii* with sequence
1545 similarity to land plant LRR FBPs were *in silico* modeled using "intensive" modeling mode in
1546 Protein Homology/analogY Recognition Engine V 2.0 (Phyre2) (Kelley et al., 2015). Various
1547 PDB molecule templates (coronatine-insensitive protein 1: Chain B (c3ogmB) and Chain D
1548 (c3oglD); transport inhibitor response 1: Chain E (c2p1nE); f-box/lrr-repeat max2 homolog:
1549 Chain A (c5hywA), skp2: Chain C (c1fs2C) and Chain K (c1fqvk); and protein toll: Chain A
1550 (c4lxrA)) were sele-cted to model *C. braunii* LRR FBPs based on heuristics to maximize
1551 confidence, percentage identity and alignment coverage. Structural prediction from regions

1552  modeled *ab initio* are highly unreliable. The final models (color-coded by the confidence of the
1553  match to the templates overall) were submitted to 3DLigandSite server (Wass et al., 2010) to
1554  predict potential binding sites (gray structures cartoon depiction); see Data S1P.

### *Phytohormones: CK*

1556  In order to identify putative CK receptors, BLAST searches were carried out against the *C.*
1557  *braunii* gene models and genome assembly, using PpCHK4 and AHK4 as queries. The detected
1558  sequences were run against the Interpro and PFAM databases to detect the domains (histidine
1559  kinase and response regulators) which are found in CK receptors. Two sequences were
1560  identified containing the domain architecture of CK receptors (CHBRA123g00790 and
1561  CHBRA19g00270). In order to identify putative histidine phosphor transfer protein (HPT), a
1562  search with the HPT domain (Interpro IPR008207) was conducted and retrieved one sequence
1563  (CbHPT1, CHBRA650g00040) (Table S1J). For identification of the response regulators (type-
1564  A and type-B) we used the PFAM domains Response_reg (PF00072) and Myb_DNA-binding
1565  (PF00249) in an hmmsearch and did not find any gene models. In order to make sure that this
1566  result is not due to a missing or fragmentary gene model we also screened the available
1567  transcriptome data (transcripts were translated in all possible frames). While two A-type
1568  response regulators (RRA) could be detected in the transcriptome (comp31700c0seq1num3,
1569  comp64895c0seq1/2 rc num2, Table S1J/S1K, Data S1H), no combination of the two domains
1570  and thus no B-type (RRB) could be detected. All sequences harboring Response_reg domains
1571  were aligned with the response regulator domains of the *Arabidopsis* response regulators ARR1
1572  and ARR14 (RRB) as well as ARR4 and ARR9 (RRA) and ARR 22 (RRC – not known to be
1573  involved in CK signaling) using the muscle implementation of the MEGA 7.0 suite. Using the
1574  alignment, a maximum likelihood tree was calculated with the pairwise distances estimated by
1575  a JTT model and 100 bootstrap samples. Again, two sequences were determined as RRAs. Of
1576  the *Chara* sequences in the RRB clade, again none contained a MYB domain (Data S1H).

### *Photorespiration*

1578  In land plants, the canonical photorespiratory pathway employs 8 enzymes, namely 2PG-
1579  phosphatase (PGPase), glycolate oxidase (GOX), glutamate:glyoxylate aminotransferase
1580  (GGT), glycine decarboxylase (GDC), serine hydroxymethyltransferase (SHMT),
1581  serine/alanine:glyoxylate aminotransferase (SGT), hydroxypyruvate reductase (HPR) and
1582  glycerate 3-kinase (GLYK) (Bauwe et al., 2010). Particularly, the glycolate oxidation step,
1583  which is performed by GOX in the plant peroxisomes, is catalysed by glycolate dehydrogenase
1584  in the mitochondrium of the green algae *C. reinhardtii* (Nakamura et al., 2005) and in the
1585  cytosol of cyanobacteria. To analyze the photorespiration in the Charophyte algae *C. braunii*,
1586  the protein sequences of enzymes from *A. thaliana* were used to identify homologue proteins
1587  in *C. braunii* by a BLASTp similarity search against the Chbra.pep.20151207.orcae database
1588  (Table S1M). To verify, if *C. braunii* also possess genes to oxidize glycolate via a glycolate
1589  dehydrogenase like Chlorophytes and cyanobacteria do, the polyphyletic proteins from *C.*
1590  *reinhardtii* (ABG36932.1) and *Synechocystis sp.* PCC 6803 (Sll0404 and Slr0806) were used
1591  as templates in similarity searches. To verify, if a putative glycolate oxidase prefers the substrate
1592  glycolate over lactate, three amino acids in the active site that were shown to be responsible for
1593  the substrate preference (Hackenberg et al., 2011) were analyzed. To this end, the putative
1594  glycolate oxidase from *C. braunii* and verified glycolate oxidase proteins of the land plants *A.*

1595 *thaliana* and *Spinacia oleracea*, the red alga *Cyanidioschyzon merolae* and characterized L-
1596 lactate oxidase proteins from the cyanobacterium *Nostoc* sp. PCC 7120 and the bacterium
1597 *Aerococcus viridans* were aligned and the corresponding amino acids in the active sites of the
1598 proteins compared.

### *Retrograde signaling and PAPs*

1600 Protein data from the genomes of *C. reinhardtii, K. nitens, C. braunii,* and *P. patens* was
1601 screened for orthologs of the flowering plant-type retrograde signaling pathway or PAPs *via* a
1602 reciprocal best BLASTp approach using *A. thaliana* sequenes as query. For GUN1, the
1603 BLASTp analyses were repeated using reciprocal pHMMER surveys. To further pinpoint the
1604 relation of *Cb*GUN1 to other PPRs, the high similarity *K. nitens* protein GAQ81958.1 was used
1605 as a query in BLASTP (2.2.26) search to a database comprising the NCBI nr dataset as of
1606 January 2015 supplemented with *K. nitens*, *Pinus taeda* 1.01, and *P. patens* v3.3 Ppav3.3
1607 datasets and 912 hit sequences were retrieved through (http://moss.nibb.ac.jp/cgi-bin/blast-nr-
1608 Kfl). Two *C. braunii* proteins Cbr_g9159.t1 (GUN1) and Cbr_g31394.t1, and a *M. polymorpha*
1609 protein Mapoly0154s0039.1 were added to this set. From this set, top 500 hits with
1610 GAQ81958.1 were retrieved and aligned with mafft version 6.811b and converted to nexus
1611 format file through (http://moss.nibb.ac.jp/cgi-bin/selectNalign). The alignment was edited to
1612 retain 242 aa (others were excluded; further 47 proteins that showed low conservation in the
1613 retained regions were deleted). The nexus file was subjected to http://moss.nibb.ac.jp/cgi-
1614 bin/makenjtree to construct a NJ tree based on JTT distance with 1,000 bootstraps using
1615 PHYLIP 3.695. Sequences identical within the retained 242 aa sites were treated as a single
1616 OTUs and 381 OTUs remained in the final tree. The organism name the sequence originated
1617 was recovered using NCBI taxonomydb
1618 (ftp://ftp.ncbi.nih.gov/pub/taxonomy/accession2taxid/prot.accession2taxid.gz,
1619 ftp://ftp.ncbi.nih.gov/pub/taxonomy/taxdump.tar.gz). The subcellular localization of PAPs was
1620 predicted using three online tools (Table S1N).

### *Transcription factors and transcriptional regulators*

1622 Transcription associated proteins (TAPs) comprise transcription factors (TFs, acting in
1623 sequence-specific manner, typically by binding to *cis*-regulatory elements) and transcriptional
1624 regulators (TRs, acting on chromatin or via protein-protein interaction. We classified all *C.*
1625 *braunii* proteins into 122 families and sub families of TAPs by first screening the proteins for
1626 domains and then applying a domain-based rule set to distinguish the TAPs (Lang et al., 2010;
1627 Wilhelmsson et al., 2017). We compared this genome-wide classification with genomic protein
1628 sets from *Cyanidioschyzon merolae*, *C. reinhardtii*, *Cyanophora paradoxa*, *K. nitens* and
1629 several land plants, as well as with transcriptomic data of Charophyta (Timme et al., 2012), *M.*
1630 *polymorpha* and ferns (Table S1Q, S1Z). The phylogenetic tree for the trihelix family (Data
1631 S1J) was inferred as mentioned above for the cell division related families.

1632 For the Homeodomain (HD) and basic Helix-Loop-Helix (bHLH) phylogenetic analyses (Table
1633 S1O, S1P), the *C. braunii* genome was searched using a BLASTp query that was assembled
1634 from the previously characterized bHLH and HD protein sequences (Catarino et al., 2016) in
1635 *At, A. thaliana; Os, O. sativa; Sm, Selaginella moellendorffi; Pp, P. patens; Mp, M.*
1636 *polymorpha; Kf, K. nitens; Cr, C. reinhardtii; Ot, Ostreococcus tauri; Vc, Volvox carteri; Cm,*

1637  *C. Merolae* with the addition of bHLH proteins sequences from *Cv, Coccomyxa subellipsoidea*
1638  (previously *Chlorella vulgaris*. The results of the BLASTp search were analyzed manually to
1639  ensure the presence of the HD or the bHLH conserved domain using SMART and PFAM. All
1640  protein sequences were aligned using MAFFT (Katoh and Standley, 2013) and further manually
1641  aligned independently for HD and bHLH. The Maximum likelihood analysis was carried out
1642  using PhyML (Guindon et al., 2010) 3.0, using the JTT amino acid substitution model and a
1643  predicted gamma distribution. Branch support was tested using a Shimodaira-Hasegawa-like
1644  approximate likelihood ratio test (SH-like aLRT). The generated unrooted trees were visualised
1645  using MEGA 6.0.

1646  MADS box sequences were identified using the aforementioned domain-based rule set to
1647  distinguish the TAPs (Lang et al., 2010). Phylogenies were calculated with MrBayes
1648  (Huelsenbeck and Ronquist, 2001) applying mixed AA model for 50,000,000 generations based
1649  on an amino acid alignment of Type I and Type II MADS-domain proteins from a broad set of
1650  land plants together with MADS-domain proteins from charophytes. Sequences were aligned
1651  with MAFFT (Katoh and Standley, 2013) applying E-INS-i mode. Intron structure was
1652  determined by using the transcript sequence as query for BLAST searches against the genome
1653  scaffolds. Subsequently, the genomic region that harbors the gene was extracted and aligned to
1654  the transcript sequence.

1655  *Motor proteins*

1656  PFAM domains related to the three classes of motor proteins were retrieved from the whole
1657  predicted proteomes of *C. braunii*, *C. reinhardtii*, *P. patens*, and *A. thaliana* using Interproscan
1658  (Table S1S). These selected domain signatures not only include the true motors but also
1659  domains associated with the tasks the motors have to fulfill in a cell. Since motor proteins are
1660  comparably long gene prediction on draft genomes can lead to a slight overestimation of domain
1661  numbers. Thus, retrieved predicted gene structures were examined, whether they reside adjacent
1662  to another predicted gene encoding for a motor protein part. If the domain structures from
1663  known complete proteins conformed with a fusion of two or more adjacent gene models in *C.*
1664  *braunii*, we used this fused gene model for further analysis.

1665  *Action potential related ion channels and transport proteins*

1666  Ion channels, transporters and pumps predicted to be involved in electrical signaling in plants
1667  were identified in the *C. braunii* genome *via* a tBLASTn/BLASTp approach using *A. thaliana*
1668  sequences as bait as well as on the basis of PFAM domains. Subsequent BLASTp searches of
1669  retrieved sequences against TAIR10 (https://www.arabidopsis.org) and SWISSPROT were
1670  employed to identify closest homologs. Finally, sequences were were classified into respective
1671  transporter families according to TCDB (Saier et al., 2016) and ARAMEMNON (Schwacke et
1672  al., 2003) (Table S1R). When partially split models were found, they were manually annotated
1673  with reference to RNA-seq evidence through a genome browser at https://chara.asrc.kanazawa-
1674  u.ac.jp/Cbr1/jbrowse/.

1675
1676  *LysM-RLKs*

1677  The *C. braunii* genome was screened for LysM-RLK genes via tBLASTn using Medicago NFP
1678  and Rice CERK1 as bait sequences (Table S1V). Hits with E-value $< 10^{-30}$ were collected and

1679  deduplicated. These sequences were aligned using MAFFT (Katoh and Standley, 2013) with
1680  LysM-RLKs from embryophytes and *Nitella mirabilis*. Using MEGA 6.0 the best substitution
1681  model (JTT+G) was determined and a maximum likelihood tree was inferred using all sites and
1682  100 bootstrap resamplings (Fig. 5C, Data S1L-N).

1683  ***PPR proteins***

1684  Genomic protein sets were scanned for presence of the PFAM domain PPR
1685  (http://pfam.xfam.org/family/PF01535) using HMMscan. The number of proteins harboring
1686  two or more PPR domains were considered PPR proteins putatively involved in organellar RNA
1687  editing (Maier et al., 2008) and are shown in Table S1Y.

1688  ***ROS-associated genes***

1689  21 families belonging to the well-known reactive oxygen species (ROS) gene network were
1690  searched using as a first screen the follwing PFAM. PF00141 for Class III Prx (CIII) and
1691  Ascorbate Prx (APx and APx-R), PF00199 and PF06628 for catalases (Kat), PF00255 for
1692  glutathione Prx (GPx), PF00578 and PF08534 for peroxiredoxin family, PF03098 for
1693  dioxygenase (DiOx), PF08022, PF01794, PF08030 and PF08414 for NADPH Oxidase (RBOH)
1694  and Ferric reduction oxidase (FRO), PF02777 and PF00080 for superoxide dismutase family
1695  (MnSOD, FeSOD, Cu/ZnSOD), PF00462 for Glutaredoxins superfamily, PF01786 for
1696  Alternative Oxidase (AOX and PTOX), PF02298 for Blue-copper-binding protein superfamily,
1697  PF00210 for ferritin (FER), PF13417 for dehydroascorbate reductase (DHAR), PF07992 and
1698  PF02852 for Monodehydroascorbate reductase (MDAR) and Glutathione reductase (GR),
1699  PF07992, PF02943 and PF00085 for thioredoxin superfamily and PF01070 Glycolate Oxidases
1700  (GOx). *Arabidopsis* sequences belonging to the "ROS gene network" have been used to confirm
1701  the *C. braunii* families affiliation.

1702  Only alpha-DiOxygenase (DiOx) and APx-R were not detected in the *C. braunii* assembly. The
1703  19 other families have been found in *C. braunii* with various conservation rates (Table S1X).
1704  Among these families, Class III peroxidases (Prx), described as secreted peroxidases, are
1705  usually members of a large family. The *C. braunii* genome contained 14 homologous sequences
1706  (Table S1X), which is much lower as compared with flowering plants (73 in *A. thaliana*) but
1707  higher than in *K. nitens* (3). All the 14 sequences are derived from a single gene in an ancestor
1708  of *C. braunii* as they form a presumably monophyletic clade (Data S1O). Before these
1709  duplication events only one or a few initial sequences may have existed, implied by the single
1710  sequence detected in *Chlorokybus atmophyticus* transcriptome data (Timme et al., 2012) and
1711  the low number of three sequences found in *K. nitens*. The CIII Prx protein sequences from *K.*
1712  *nitens* (3 sequences), *C. braunii* (14 sequences), *P. patens* (57 sequences) and *A. thaliana* (73
1713  sequences) were aligned using MAFFT and the tree constructed using Maximum Likelihood
1714  implemented in MEGA (Data S1O).

1715  ***UBQ proteasome system (UPS)***

1716  *Arabidopsis* genes encoding components of the plant Ubiquitin proteasome system (UPS) were
1717  manually selected and used as query sequences in a tBLASTn analysis to identify respective
1718  orthologous genes in the *C. braunii* genome. Hits with E-values $< 10^{-10}$ were collected and
1719  annotated following a reciprocal best BLASTp approach using TAIR10 (Table S1I).

## QUANTIFICATION AND STATISTICAL ANALYSES

All details of the applied statistics (*e.g.* for RNAseq-based differential gene expression analysis) are provided alongside the respective analysis in the Methods Details section. For the differential gene expression analysis between antheridia, oogonia, and zygotes, three true biological replicates were sequenced and used for the statistical analysis (computed using DESeq2). No sequencing points, i.e. samples, were removed during the analysis.

## DATA AND SOFTWARE AVAILABILITY

Raw Illumina (DRA004353, DRA006568) and PacBio (DRA006569) genomic sequence data have been deposited in the DDBJ Sequence Read Archive (DRA) at the DNA Data Bank of Japan (DDBJ) under BioProject PRJDB3348. The main scaffolds are available as entries BFEA01000001-BFEA01011654, the accompanying organisms scaffolds as BFBZ01000001-BFBZ01016437. The chloroplast genome is available as AP018555, the mitochondrial as AP018556. Raw Illumina RNA-seq data used for annotation (DRA006080, DRA002641) have been deposited in the DRA at the DDBJ under BioProject PRJDB3228. Raw Illumina RNA-seq data of reproductive stages have been deposited to NCBI SRA (PRJNA445548). The genome and its annotation is available for human curation *via* the ORCAE interface at the URL: http://bioinformatics.psb.ugent.be/orcae/. The data is freely available for browsing as well as for bulk downloads and blast searches. Persons who would like to contribute and edit the data using the web interface will have to request an account by sending an email. Any change made to gene structures will be processed automatically by adding protein domains (running interpro) and best-blast hits. These changes will be shared with the community immediately. 69,969 ABI reads of a cDNA library (minimum length of 100 bp) have been deposited at the DDBJ under the accession numbers LU106825 to LU176793 (Table S1D). Alignments that are the basis for the phylogenetic trees as well as the genome comparison datasets resulting in Fig. 3 have been deposited as Mendeley Datasets (doi:10.17632/9hzzf9m4kh.1).

## Supplemental Tables and Files

The supplemental/supporting information is arranged into:

- a PDF containing Figures S1-S7;
- a PDF containing phylogenetic trees and alignments Data S1A-U;
- five Excel spreadsheets containing Tables S1-5 (with indexing of sheets);
- alignments and supporting data for the genome comparisons (related to Fig. 3) in Mendeley (doi:10.17632/9hzzf9m4kh.1).

## Supplemental Tables

**Table S1, related to STAR methods: details of assembly, annotation and comparative analyses, with index in first sheet.**

Table S1A: Genome libraries and accession numbers

Table S1B: Libraries used for assembly

Table S1C: Cell division

Table S1D: EST data deposited in DDBJ

Table S1E: RNA-seq used for annotation

Table S1F: Repeatmasker results

Table S1G: Repetitive elements

Table S1H: Cell wall biosynthesis

Table S1I: UBQ proteasome system

Table S1J: Phytohormones

Table S1K: Auxin signaling and transport

Table S1L: Genome comparison

Table S1M: Photorespiratory pathway

Table S1N: PAP localization prediction

Table S1O: bHLH and HD TFs comparison

Table S1P: bHLH and HD TFs *C. braunii*

Table S1Q: Transcription factors and transcriptional regulators

Table S1R: Ion channels

Table S1S: Motor proteins

Table S1T: Assembled bacterial genomes

Table S1U: Most abundant bacterial genera

1783    Table S1V: LysM RLKs

1784    Table S1W: Reproductive transcriptome

1785    Table S1X: ROS network

1786    Table S1Y: PPR proteins

1787    Table S1Z: Transcription factors and transcriptional regulators gene Ids

1788    Table S1AA: genome and transcriptome datasets used for comparative studies

1789

1790    **Table S2, related to Fig. 5/6 and STAR methods: Differential gene expression analyses of**
1791    **rerproductive stages, with index in first sheet.**

1792    Table S2A: zygote versus oogonia

1793    Table S2B: oogonia versus antheridia

1794

1795    **Table S3, related to Fig. 5/6 and STAR methods: Gene Ontology analyses of the**
1796    **differential expression data in Table S2, with index in first sheet.**

1797    Table S3A: zygote versus oogonia up GO enrichment

1798    Table S3B: zygote versus oogonia up genes

1799    Table S3C: zygote versus oogonia down GO enrichment

1800    Table S3D: zygote versus oogonia down genes

1801    Table S3E: oogonia versus antheridia up GO enrichment

1802    Table S3F: oogonia versus antheridia up genes

1803    Table S3G: oogonia versus antheridia down GO enrichment

1804    Table S3H: oogonia versus antheridia down genes

1805

1806    **Table S4, related to Fig. 5/6 and STAR methods: *C. braunii* protein coding gene annotation**
1807    **(Gene Ontology, best blast hits, trihelix TFs, expression data, overlap with TE evidence,**
1808    **decontamination).**

1809

1810    **Table S5, related to STAR methods: Decontamination analyses.**

1811    S5A    summary

1812    S5B    underlying data

1813

1814    **Data S1, related to STAR methods: phylogenetic trees and alignments.**

1815

**KEY RESOURCES TABLE**

| REAGENT or RESOURCE | SOURCE | IDENTIFIER |
|---|---|---|
| Critical Commercial Assays | | |
| DNeasy Plant Mini Kit | Qiagen | Cat# 69106 |
| Ex Taq | Takara Bio, Shiga, Japan | Cat# RR001A |
| BAP | Takara Bio, Shiga, Japan | Cat# 2120A |
| Fruit-mate for RNA Purification | Takara Bio, Shiga, Japan | Cat# 9192 |
| T4 RNA ligase | Takara Bio, Shiga, Japan | Cat #2050A |
| Genomic Tip | Qiagen | Cat# 10243 |
| ISOGEN | Nippon Gene, Tokyo, Japan | Cat# 311-02501 |
| RNasin | Promega | Cat # N2111 |
| MinElute Gel Extraction Kit | Qiagen | Cat# 28604 |
| mirVana | Ambion | Cat# AM1560 |
| mRNA-Seq Sample Prep Kit | Illumina | Cat# RS-100-0801 |
| Nextera Mate-pair library construction kit | Illumina | Cat# FC-132-1001 |
| NxSeq 40 kb Mate-Pair Cloning Kit | Lucigen | Cat# 42028-1 |
| Ovation RNA-Seq System V2 | NuGEN | Cat# 7102-32 |
| RNeasy Plant Mini Kit | Qiagen | Cat# 74904 |
| Schiff's reagent | Merck Millipore | Cat# 1.09033.0500 |
| Small RNA Sample Preparation Kit | Illumina | Cat# FC-102-1009 |
| TruSeq DNA PCR-Free LT Sample Prep Kit | Illumina | Cat# FC-121-3001 |
| Zymoclean Large Fragment DNA Recovery Kit | Zymo Research | Cat# D4045 |
| Deposited Data | | |
| ARAMEMNON (plant membrane protein database) | Schwacke et al. 2003 | http://aramemnon.uni-koeln.de |
| Carbohydrate-Active enZYmes Database (CAZy) | Lombard et al., 2014 | http://www.cazy.org |
| *Chara braunii* ABI reads cDNA libraries | This study | DDBJ accessions LU106825 to LU176793 |
| *Chara braunii* Illumina RNA-seq data of reproductive stages | This study | BioProject PRJNA445548 |
| *Chara braunii* Illumina RNA-seq data used for annotation | This study | BioProject PRJDB3228 |
| *Chara braunii* PacBio and Illumina genomic DNA sequencing data | This study | BioProject PRJDB3348 |
| Genomic and transcriptomic data used for comparative analysis, see Table SAA | This study | n/a |
| Glycosyltransferase repertoire of *S. moellendorffii* and *P. patens* | PMID: 22567114 | https://doi.org/10.1371/journal.pone.0035846.s017 |
| Phylogenetic trees and alignments, data for Fig. 3 | This study | Mendeley doi:10.17632/9hzzf9m4kh.1 |
| Transporter Classification Database (TCDB) | Saier et al. 2016 | http://tcdb.org |

| Experimental Models: Organisms/Strains | | |
|---|---|---|
| *Chara braunii* S276 | isolated from soil of Lake Kasumigaura (Ibaraki, Japan) | maintained at Kobe University; Herbarium press TNS-AL 209137 available at the National Science Museum (TNS), Tsukuba, Japan |
| *Chara braunii* S277 | collected from a pond at Saijo (Ehime, Japan) for this study | maintained at Kobe University; Herbarium press TNS-AL 209138 available at the National Science Museum (TNS), Tsukuba, Japan |
| Software and Algorithms | | |
| 3DLigandSite | Wass et al., 2010 | http://www.sbg.bio.ic.ac.uk/3dligandsite/ |
| ALLPATHS-LG | Gnerre et al., 2011 | http://software.broadinstitute.org/allpaths-lg/blog/?page_id=12 |
| Augustus | Keller et al., 2011 | http://bioinf.uni-greifswald.de/augustus/ |
| BEDtools v2.25.0 | Quinlan, 2014 | http://bedtools.readthedocs.io/en/latest/ |
| Bioconductor Package GenomicRanges | Lawrence et al., 2013 | https://bioconductor.org/packages/release/bioc/html/GenomicRanges.html |
| Burrows-Wheeler Aligner (bwa mem v.0.7.8-r455) | Li and Durbin, 2009 | https://sourceforge.net/projects/bio-bwa/files/ |
| CEGMA | Parra et al., 2007 | http://korflab.ucdavis.edu/datasets/cegma/#SCT8 |
| CLC Assembly Cell | QIAGEN Bioinformatics | https://www.qiagenbioinformatics.com/products/clc-assembly-cell/ |
| CONCOCT | Alneberg et al., 2014 | https://github.com/BinPro/CONCOCT |
| Cufflinks v2.0.2 | Trapnell et al., 2010 | https://github.com/cole-trapnell-lab/cufflinks |
| DESeq2 v1.14.1 | Love et al., 2014 | https://bioconductor.org/packages/release/bioc/html/DESeq2.html |
| ESTScan | Iseli et al., 1999 | http://estscan.sourceforge.net |

| feature v1.2.13 | Duong et al., 2008 | https://cran.r-project.org/web/packages/feature/index.html |
|---|---|---|
| GAG - Genome Annotation Generator V1.0 | released under an MIT License (MIT), copyright © 2014 GAG Developers | http://genomeannotation.github.io/GAG/ |
| GenomeTools [gff3/LTRdigest/LTRharvest] V1.5.9 | Steinbiss et al., 2009 | http://genometools.org/ |
| HelitronScanner V1.0 | Xiong et al., 2014 | https://sourceforge.net/projects/helitronscanner/files/ |
| HGAP & Quiver | Chin et al., 2013 | https://www.pacb.com/support/software-downloads/ |
| IQ-Tree v1.5.3 | Nguyen et al., 2015 | http://www.iqtree.org |
| JELLYFISH | Marçais and Kingsford, 2011 | http://www.cbcb.umd.edu/software/jellyfish/ |
| Jupyter Notebook and IRKernel | Thomas Kluyver, Philipp Angerer, Jan Schulz | http://jupyter.org/ https://github.com/IRkernel/IRkernel |
| KeyS | Rensing et al., 2007 | http://plantco.de/research.html |
| MAFFT v6.811b / v7.305b | Katoh and Standley, 2013 | https://mafft.cbrc.jp/alignment/software/ |
| mclust v5.1 | Scrucca et al., 2016 | https://CRAN.R-project.org/package=mclust |
| MEGAN5 v5.11.3 / v6 | Huson et al., 2016 | http://ab.inf.uni-tuebingen.de/software/megan6/ |
| MrBayes v3.2.6 | Huelsenbeck and Ronquist, 2001 | http://mrbayes.sourceforge.net/ |
| MUMmer v3.23 | Delcher et al., 1999 | http://mummer.sourceforge.net |
| MUSCLE v3.8.31 | Edgar et al., 2004 | https://www.drive5.com/muscle/ |
| NOVOPlasty ver 2.5.3 | Dierckxsens et al., 2017 | https://github.com/ndierckx/NOVOPlasty |
| PAML v4.7 | Yang, 2007 | http://abacus.gene.ucl.ac.uk/software/paml.html |
| PASTEC | Hoede et al., 2014 | https://urgi.versailles.inra.fr/Tools/PASTEClassifier |
| pheatmap v1.0.8 | Raivo Kolde | https://cran.r-project.org/web/packages/pheatmap/index.html |
| PHYLIP 3.695 | Jerry Shurman, Mark Moehring, Joe Felsenstein | http://evolution.genetics.washington.edu/phylip.html |

| | | |
|---|---|---|
| PhyML 3.0 | Guindon et al., 2010 | http://www.atgc-montpellier.fr/phyml/ |
| Phyre2 | Kelley et al., 2015 | http://www.sbg.bio.ic.ac.uk/phyre2/html/page.cgi?id=index |
| Picard-tools 1.129 | Broad Institute | http://broadinstitute.github.io/picard. |
| Prottest | Darriba et al. 2011 | https://github.com/ddarriba/prottest3 |
| RepeatMasker version open-4.0.5 | Arian F.A. Smit, Robert Hubley & Phil Green | http://www.repeatmasker.org |
| RepeatModeler Version open-1.0.7 | Arian F.A. Smit and Robert Hubley | http://www.repeatmasker.org/RepeatModeler/ |
| REPET package v2.4 | Flutre et al., 2011 | https://urgi.versailles.inra.fr/Tools/REPET |
| RSEM v1.2.11 | Li and Dewey, 2011 | https://github.com/deweylab/RSEM |
| RNAmmer 1.2 | Lagesen et al, 2007 | http://www.cbs.dtu.dk/services/RNAmmer/ |
| SINA Alignment Service (Silva database for classification) | Pruesse et al, 2012 | https://www.arb-silva.de/aligner/ |
| Smrtanalysis 2.0.1 | PacificBiosciences | http://programs.pacificbiosciences.com/l/1652/2013-06-04/2t28z7 |
| Snakemake v4.3.1 | Köster and Rahmann, 2012 | https://snakemake.readthedocs.io/en/stable/ |
| Tandem Repeats Finder (TRF) | Benson, 1999 | http://tandem.bu.edu/trf/trf.html |
| tera-BLASTn 9.0.0 | Active Motif, Inc. | http://www.timelogic.com/catalog/757/tera-blast |
| topGO v2.22.0 | Adrian Alexa and Jörg Rahnenführer | https://bioconductor.org/packages/release/bioc/html/topGO.html |
| Tophat v2.1.0 | Kim et al., 2013 | https://ccb.jhu.edu/software/tophat/index.shtml |
| TransDecoder v2.0.1 | Haas et al, 2013 | https://github.com/TransDecoder/TransDecoder |
| Trimmomatic v0.36 | Bolger et al., 2014 | http://www.usadellab.org/cms/?page=trimmomatic |
| Vmatch v2.3.0 | Abouelhoda et al., 2004 | http://www.vmatch.de/ |
| Web Apollo | Lee et al., 2013 | http://genomearchitect.github.io |
| Other | | |

| *Chara braunii* genome interface for gene models open for human curation | This study | http://bioinformatics. psb.ugent.be/orcae/ |
|---|---|---|
| DeCypher 9.0.0.25 (Biocomputing Platform) | TimeLogic | http://www.timelogic. com/catalog/752/bio computing-platforms |

Figure 1                                                    Click here to download Figure Figure1.pdf ⬇



Figure 1

Figure 2

**Oogonium and antheridium**
- Oogonium
- Antheridium

**Branchlet cells**
- Nuclei
- Cell wall
- Central vacuole
- ♀
- Plastid
- ♂
- Nodal cell
- Branchlet cell

**Internodal cell**
- Plastid
- Nuclei
- Central vacuole
- Charasome
- Plasma membrane
- Cell wall

**Node**
- Internodal cell
- Branchlet cell
- Nodal cell
- Central cell
- Stipulode

**Thallus**
- Oogonium and antheridium
- Branchlet
- Stipulode
- Stem
- Rhizoid

**Rhizoid**
- Rhizoidal plate
- Rhizoidal cell

**Life cycle of *Chara braunii***
- Thallus (n)
- Antheridium (n)
- LS antheridium (n)
- Antheridial filament (n)
- Sperm cell (n)
- Protonema (n)
- Young rhizoid and protonema (n)
- Meiosis
- Diploid zygote (2n)
- Fertilization
- LS oogonium (n)
- Oogonium (n)

Figure 3                                                    Click here to download Figure Figure3.png

Figure 4

Click here to download Figure Figure4.png



Figure 4

Figure 6

Click here to download Figure Figure6.png ⬇

Nishiyama et al.

**The *Chara* genome: secondary complexity and implications for plant terrestrialization**

<u>Data S1: Phylogenetic trees and alignments</u>

Throughout this document, *C. braunii* sequences/clades are marked green where suitable. Similarly, *K. nitens* (previously: *K. flaccidum*) is marked blue, *A. thaliana* / blast query sequences in red.

Figures are referenced by Data S1 and an alphabetic index here and when referred to in the manuscript. All alignments, showing the accession numbers of all sequences used, are available as fasta files from Mendeley (doi:10.17632/9hzzf9m4kh.1).

Methodological details of phylogenetic tree inference are in STAR methods.

Five letter codes abbreviating the species names are comprised of the first three letters of the genus followed by the first two letter of the species, e.g. ARATH = *Arabidopsis thaliana* and can be looked up in the key resource table. Transcriptomic datasets are marked by "tr" after the five letter species code.

Data S1A:    Maximum Likelihood GT43 phylogeny.
Data S1B:    Maximum Likelihood GT47 phylogeny.
Data S1C:    PIN phylogeny.
Data S1D:    PIN alignment.
Data S1E:    Maximum likelihood phylogenetic reconstruction for the ARF auxin signaling gene family.
Data S1F:    Maximum likelihood phylogenetic reconstruction for the AUX/IAA auxin signaling gene family.
Data S1G:    Maximum likelihood phylogenetic reconstruction for the TIR1/AFB auxin signaling gene family.
Data S1H:    Maximum likelihood tree of response regulators.
Data S1I:     *Cb*GUN1 clusters with land plant GUN1 proteins.
Data S1J:    Phylogenetic tree of trihelix transcription factors.
Data S1K:    Phylogenetic tree of plant MADS-box genes.
Data S1L:    Condensed Maximum likelihood tree of the LysM-RLK family.
Data S1M:    Alignment of the Chara LysM-RLK with representative LYK and LYR from angiosperms.
Data S1N:    Midpoint-rooted Maximum Likelihood tree of the LysM-RLK family.
Data S1O:    Outgroup-rooted phylogenetic analysis of *C. braunii* class III peroxidases
Data S1P:    *In silico* modeling of *C. braunii* LRR FBPs.
Data S1Q:    Phylogenetic tree of cyclin proteins.
Data S1R:    Phylogenetic tree of EXOCYST 70A proteins.
Data S1S:    Phylogenetic tree of NPSN proteins.
Data S1T:    Phylogenetic overview tree of canonical actin proteins.
Data S1U:    Phylogenetic tree of canonical actin proteins.

**IRX9/9L**

**IRX14/14L**

**Data S1A: Maximum Likelihood GT43 phylogeny.**

GT43 contain activities involved in xylan biosynthesis, in a so far unknown manner. Both clades are not, at least in higher plants, suggested to be direct xylan synthases. *C. braunii* has a single member (green) with an apparent ancestral location in neither of the defined clades of GT43. *A. thaliana* sequences are marked red, *K. nitens* sequences blue.

**Data S1B: Maximum Likelihood GT47 phylogeny.**

Phylogenetic tree of GT47 showing the identified *C. braunii* sequences in Clade B and Clade E, putative arabinosyltransferases involved in the biosynthesis of arabinan and extensin, respectively. Note the absence of *C. braunii* sequences from Clade A (xyloglucan galactosylation) and Clade D (xylan synthase). Presence of *C. braunii* and *K. nitens* sequences in clades is shown by a green or blue box on the branch of that clade.

*#*



**Data S1C: PIN phylogeny.**

Maximum likelihood reconstruction of PIN protein family focused on algal species. 500 bootstrap replications, result (in %) indicated above branches. *C. braunii* sequences are marked by a green box, the *K. nitens* sequence by a blue box, and the clade harboring the *A. thaliana* sequences by a red box.

# N-terminal domain



# C-terminal domain



**Data S1D: PIN alignment.**

*C. braunii* sequences are shown on top.

PHYPAT_Pp3c4_12970V3.1.p_14-693
PHYPAT_Pp3c4_13010V3.1.p_14-693
PHYPAT_Pp3c6_26890V3.1.p_14-608
MARPOL_OAE23454.1_8-625
ARF10; ARF16
PICABI_MA_2237g0020.1_12-704
PICABI_MA_2421g0010.1_5-687
PICABI_MA_98506g0010.1_4-681
ARF17
AMBTRI_XP_011624048.1_4-535
AMBTRI_XP_011624395.1_4-535
COLORB_GBSL01017215.1_17-456
SELMOE_XP_002961221.1_42-735
SELMOE_XP_002969449.1_1-694
SELMOE_XP_002976921.1_9-664
SELMOE_XP_002980664.1_9-683
COLORB_GBSL01031614.1_119-532
COLORB_GBSL01031615.1_196-609
COLORB_GBSL01031616.1_163-576
MARPOL_OAE31269.1_23-796
PHYPAT_Pp3c16_6100V3.1.p_4-716
PHYPAT_Pp3c27_60V3.1.p_4-704
PHYPAT_Pp3c5_9420V3.1.p_5-714
PHYPAT_Pp3c6_21370V3.1.p_4-695
ARF3; ARF4
PICABI_MA_10432528g0010.1_7-338
PICABI_MA_10431460g0020.1_16-352
PICABI_MA_10432349g0010.1_13-413
ARF2
PICABI_MA_10432580g0010.1_8-803
ARF1; ARF9; ARF11; ARF12; ARF13; ARF14;
ARF15; ARF18; ARF20; ARF21; ARF22
PICABI_MA_40335g0010.1_20-809
PICABI_MA_92651g0010.1_8-812
SELMOE_XP_002961629.1_4-740
SELMOE_XP_002964700.1_4-369
SELMOE_XP_002964106.1_8-699
SELMOE_XP_002992064.1_8-699
ARF6; ARF8
PICABI_MA_10121946g0020.1_1-124
PICABI_MA_10022392g0010.1_1-204
PICABI_MA_164845g0010.1_1-227
PICABI_MA_167372g0010.1_1-227
PHYPAT_Pp3c13_4720V3.1.p_3-713
PHYPAT_Pp3c26_11550V3.1.p_3-708
PHYPAT_Pp3c1_14440V3.1.p_1-269
PHYPAT_Pp3c1_14480V3.1.p_3-745
PHYPAT_Pp3c2_25890V3.1.p_3-685
PHYPAT_Pp3c14_16990V3.1.p_3-494
PHYPAT_Pp3c17_19900V3.1.p_3-769
PHYPAT_Pp3c1_40270V3.1.p_3-775
MARPOL_OAE20672.1_3-770
SELMOE_XP_002982854.1_3-711
SELMOE_XP_002990954.1_3-712
SELMOE_XP_002984982.1_3-722
SELMOE_XP_002986171.1_3-711
SELMOE_XP_002984990.1_4-759
SELMOE_XP_002986179.1_4-759
PICABI_MA_122218g0010.1_3-291
ARF7; ARF19
PICABI_MA_85955g0020.1_1-266
ARF5
CHABRA_g17664.t1_97-743
NITMIR_JV767279.1_85-719

1.0

b

**Data S1E: Maximum likelihood phylogenetic reconstruction for the ARF auxin signaling gene family.**

The *C. braunii* ARF candidate is highlighted in bold/green. Angiosperm ARF proteins are clustered and the corresponding *A. thaliana* members are given for the collapsed clades. The ML tree shown is rooted on the *C. braunii* sequence clade to visualize the relationship of the examined genes without suggesting any ancestry. The bootstrap values were assessed with the ultrafast bootstrap approximation using 1,000 replicates and the corresponding values are depicted next to the tree nodes on an amino acid alignment of a set of species as indicated in Table S11.

**Data S1F: Maximum likelihood phylogenetic reconstruction for the AUX/IAA auxin signaling gene family.**

The *C. braunii* AUX/IAA candidates most similar to *A. thaliana* IAA33 (red) are highlighted in bold/green; the *K. nitens* sequence is shown in blue. Angiosperm AUX/IAA proteins are clustered and the corresponding *A. thaliana* members are given for the collapsed clades. The ML tree shown is rooted on the Charales sequences to visualize the relationship of the examined genes without suggesting any ancestry. The bootstrap values were assessed with the ultrafast bootstrap approximation using 1,000 replicates and the corresponding values are depicted next to the tree nodes on an amino acid alignment of a set of species as indicated in Table S11.

COLORB_GBSL01019006.1
COLORB_GBSL01019007.1
COLORB_GBSL01022275.1
SELMOE_XP_002976462.1
SELMOE_XP_002994052.1
PHYPAT_Pp3c17_7680V3.1.p
PHYPAT_Pp3c15_1200V3.1.p
PHYPAT_Pp3c9_10120V3.1.p
MARPOL_OAE28114.1
SELMOE_XP_002960189.1
SELMOE_XP_002967516.1
SELMOE_XP_002965932.1
SELMOE_XP_002968358.1
ARATHA_AT2G39940.1 COI1
AMBTRI_XP_011624824.1
AMBTRI_XP_006848115.1
AMBTRI_XP_011624728.1
PICABI_MA_10428988g0010.1
PICABI_MA_108477g0010.1
PHYPAT_Pp3c3_8580V3.1.p
PHYPAT_Pp3c20_4050V3.1.p
PHYPAT_Pp3c23_18960V3.1.p
SELMOE_XP_002966544.1
SELMOE_XP_002982385.1
ARATHA_AT3G62980.1 TIR1
ARATHA_AT4G03190.1 AFB1
AMBTRI_XP_006844497.1
ARATHA_AT1G12820.1 AFB3
ARATHA_AT3G26810.1 AFB2
PICABI_MA_14836g0010.1
PICABI_MA_121693g0010.1
ARATHA_AT4G24390.2 AFB4
AMBTRI_XP_006841829.1
PICABI_MA_92309g0010.1
MARPOL_OAE27142.1
PHYPAT_Pp3c24_9814V3.1.p
PHYPAT_Pp3c24_9800V3.1.p
PHYPAT_Pp3c23_11300V3.1.p
PHYPAT_Pp3c20_17150V3.1.p
PHYPAT_Pp3c23_5870V3.1.p
ARATHA_AT1G47056.1 VFB1
ARATHA_AT4G07400.1 VFB3
ARATHA_AT3G50080.1 VFB2
ARATHA_AT5G67250.1 VFB4
PICABI_MA_179294g0010.1
SELMOE_XP_002970190.1
SELMOE_XP_002978355.1
SELMOE_XP_002987863.1
MARPOL_OAE20892.1
SELMOE_XP_002969264.1
SELMOE_XP_002986847.1
PHYPAT_Pp3c17_21920V3.1.p
PHYPAT_Pp3c14_26462V3.1.p
PHYPAT_Pp3c14_26460V3.1.p
AMBTRI_XP_006858612.1
KLEFLA_kfl00434_0030_v1.1
CHABRA_g48450.t1
CHABRA_g23719.t1
AMBTRI_XP_006849175.2
PICABI_MA_64990g0010.1
SELMOE_XP_002981899.1
SELMOE_XP_002986021.1
PHYPAT_Pp3c4_4760V3.1.p
PHYPAT_Pp3c12_11790V3.1.p
CHABRA_g8969.t1
NITMIR_JV762907.1
COLORB_GBSL01024071.1
COLORB_GBSL01024072.1
COLORB_GBSL01024073.1
COLORB_GBSL01024074.1
SPIPRA_JO183506.1
COLORB_GBSL01025383.1
COLORB_GBSL01025384.1
KLEFLA_kfl00020_0460_v1.1
COLORB_GBSL01056204.1
AMBTRI_XP_006850831.1
ARATHA_AT4G15475.1 F-box/RNI-like superfamily
SELMOE_XP_002983796.1
SELMOE_XP_002989023.1
PICABI_MA_10426597g0010.1
CHABRA_g11992.t1
NITMIR_JV739635.1
PICABI_MA_189740g0010.1
PICABI_MA_62909g0010.1
AMBTRI_XP_006858775.1
ARATHA_AT3G07550.1 RNI-like superfamily
CHLATM_JO192414.1
ARATHA_AT1G80570.2 RNI-like superfamily
PICABI_MA_183768g0010.1
PENMAR_JO209618.1
MESVIR_GBSK01013456.1
CHABRA_g10791.t1
CHLATM_JO196002.1
CHABRA_g23992.t1
KLEFLA_kfl00169_0150_v1.1
KLEFLA_JO262959.1

0.4

10

**Data S1G: Maximum likelihood phylogenetic reconstruction for the TIR1/AFB auxin signaling gene family.**

*Chara* proteins containing F-box or F-box-like domains with considerable similarity to the TIR1/AFB proteins from *A. thaliana* are highlighted in bold/green. For proteins from *A. thaliana* the gene alias is given. The midpoint rooted ML tree shown should visualize the relationship of the examined genes without suggesting any ancestry. The bootstrap values were assessed with the ultrafast bootstrap approximation using 1,000 replicates and the corresponding values are depicted next to the tree nodes on an amino acid alignment of a set of species as indicated in Table S11. The *A. thaliana* TIR1 sequences is marked in red, the *K. nitens* sequences in blue.

**Data S1H: Maximum likelihood tree of response regulators.**

*C. braunii* transcripts (derived from the RNA-seq based transcriptome representation) clustering with *A. thaliana* (marked in red) response regulators ARR1 and ARR14 (RRB), ARR4 and ARR9 (RRA) and ARR 22 (RRC). No response regulators could be identified in the genome; in order to make sure that this result is not due to a missing or fragmentary gene model we screened the transcriptome. While two response regulator transcripts could be detected clustering with the A-type (ARR4 and ARR9, RRA), no RRB (ARR1 and ARR14) could be detected (*cf*. STAR Methods).

**Data S1I: *Cb*GUN1 clusters with land plant GUN1 proteins**.

Neighbour-joining phylogeny of 381 tetratricopeptide (TPR) and pentatricopeptide repeat (PPR) proteins from a broad range of photosynthetic eukaryotes (see color code at the top). Note the well-supported clades of (i) GUN1 proteins, including *Cb*GUN1 (g9159, green frame), and (ii) the streptophyte algae- and bryophyte- specific PPRs, including the *K. nitens* protein kfl00096_0090 (blue frame) and the *C. braunii* protein g31394 (green frame). Values from 1,000 bootstrap replicates are shown at the nodes.

**Data S1J: Phylogenetic tree of trihelix transcription factors.**

Phylogenetic analysis based on the trihelix domain shows that most of the trihelix paralogs were probably secondarily gained in the *Chara* lineage; Bayesian inference based on the trihelix domain, midpoint-rooted, numbers at the nodes represent posterior probabilities. Land plant sequences (*A. thaliana* and *P. patens*) are shown in red, *C. braunii* sequences in green, *K. nitens* sequences in blue. Sequences of other charophytes are shown in black. Colored branches correspond to trihelix sub families

according to Kaplan-Levy et al., 2012: GT-1 blue, GT-2 green, SIP1 lavender and SH4 beige. Three *Chara* trihelix proteins group within previously defined clades, namely g38396 in SH4, g34370 in GT-1 and g29899 in GT-2. All three are expressed in vegetative tissue (supplemental file 3).

**Data S1K: Phylogenetic tree of plant MADS-box genes.**

The three previously described groups of Type I MADS-box genes in angiosperms Mα, Mβ and Mγ as well as lycophyte Type I genes are highlighted in purple, dark blue, light blue and cyan, respectively. The two groups of MIKC$^{C}$- and MIKC*-type genes found in land plants and the charophyte MIKC-type genes are highlighted in green, orange and red, respectively. The posterior probabilities of main branches are depicted next to the tree. The unrooted tree shown here only visualizes the relationship of the examined genes without suggesting any ancestry. *C. braunii* branches are marked by a green box.

**Data S1L: Condensed Maximum likelihood tree of the LysM-RLK family.**

The charophyte sequences form a single cluster (green box), presumably a clade (blue branches) encompassing seven *C. braunii* sequences. Assuming these charophyte sequence serve as outgroup to land plant LysM-RLK genes, the duplication (red circle) leading to the LYK (orange) and LYR (green) subclades occurred before the divergence of extant embryophytes. The mosses, liverworts, hornworts and *Selaginella* genes form a clade which is a sister group to a clade of angiosperm genes in each of LYK and LYR subclass.

```
                        10         20         30         40         50         60
                 ....|....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1          ---------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1  ----------------------------------------------------
g49675           MWHGEEWCNVVSAAVCAHTIDLSMDLPLWFAGTNIEDRPEDDNMAVYQESTVICIAHAFR
g50598           ---------------------------------------------------------------
g44510           MAGACYIHVLLMLLRFLATIVLFAVRPFSAAVMSFAVRPFAA----------------
g91196           ---------------------------------------------------------------
g30047           ---------------------------------------------------------------
g8984            MSQRHHDTVLWSAVRRNASLPAPVFAADFVHLLWAEGWIAVHLRDGLWFGVVVAFKTTKL
g8619            ---------------------------------------------------------------
MtNFP_-_Medtr5g019040.1  -----------------------------------------------------

                        70         80         90        100        110        120
                 ....|....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1          ---------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1  ----------------------------------------------------
g49675           RAVQMGAHIDGDFISYDRLCRVADCFRLLFAACMWIMRMAGDDPRSHYKAFYLANLLAKP
g50598           ---------------------------------------------------------------
g44510           ---------------------------------------------------------------
g91196           ---------------------------------------------------------------
g30047           ---------------------------------------------------------------
g8984            TCEQVCQQHRAEQSRAEQNRADPD--------------------------------------
g8619            ---------------------------------------------------------------
MtNFP_-_Medtr5g019040.1  -----------------------------------------------------

                       130        140        150        160        170        180
                 ....|....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1          ---------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1  ----------------------------------------------------
g49675           TLVASMHRPFDHRRSVVRAAKVVTERLGKVNATFGEYPDYIPEWEPYGIGFRHDMSITGP
g50598           ---------------------------------------------------------------
g44510           ---------------------------------------------------------------
g91196           ---------------------------------------------------------------
g30047           ---------------------------------------------------------------
g8984            ---------------------------------------------------------------
g8619            ---------------------------------------------------------------
MtNFP_-_Medtr5g019040.1  -----------------------------------------------------

                       190        200        210        220        230        240
                 ....|....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1          ---------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1  ----------------------------------------------------
g49675           EYAKKLDWQEEKAGELYVAAMGEERADKSFIEKDSKKTTVSRGGIVERILLCPLSSSCPL
g50598           ---------------------------------------------------------------
g44510           ---------------------------------------------------------------
g91196           ---------------------------------------------------------------
g30047           ---------------------------------------------------------------
g8984            ----------------------------------------------------------RSSSM
g8619            ---------------------------------------------------------------
MtNFP_-_Medtr5g019040.1  -----------------------------------------------------

                       250        260        270        280        290        300
                 ....|....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1          ---------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1  ----------------------------------------------------
g49675           PPLSSYPRSTSEHCAGMASPPQQQQPNRCAGMASPPPQQQQPKRCAGMASPPQQQPNRCA
g50598           ---------------------------------------------------------------
g44510           ---------------------------------------------------------------
g91196           ---------------------------------------------------------------
g30047           ---------------------------------------------------------------
g8984            MESDVRRAATAVECVDCHHRNNGVAAKPSLSQASGTSTSSSACVPPAYWGLFLAFVQALL
g8619            ---------------------------------------------------------------
MtNFP_-_Medtr5g019040.1  -----------------------------------------------------

                       310        320        330        340        350        360
                 ....|....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1          ------------------------------MFSLPALLIGACAFAAAAVAASGDGCRAG
MtLYK3_-_Medtr5g086130.1  ----------------------MNLKNGLLLFILFLDCVFFKVESKCVKGCDVA
g49675           GKGSLRRQPDRCAGKASLQPHTPSPLHVPLLSFSQIVLLCSLFLTPPYIPIASGLRCQVD
g50598           -------MAFYHGYPRVTPFYCPRVDSPLVNAICTHVLTLLIILARVSCALAVPCSAV
g44510           -------VVRSFAVRMLPFAVAVRSFVVRPFVAVVLSTAVTLMCATLPSALAISCSAA
g91196           --------------------MVFAATALMSMSAILSSAVGIPCSVA
g30047           ---------------------------MIPLACGVTCRSS
g8984            AGRTGKIGGGVAGSRGPTTPVLMMRCRPTAALGWLVLLWLCSGHQSLCIPAACAFACSTE
g8619            ----------MGCGTRRLHHHCRGWNGPAVLHWLLMSLLCGSHFSVLFLPCSGYPCRPA
MtNFP_-_Medtr5g019040.1  --------MSAFFLPSSSHALFLVLMLFFLTNISAQPLYISETNFTCPVD

                       370        380        390        400        410        420
                 ....|....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1          -----CSLAIAAYYFSEGS--NLTFIATIFAIGGGGYQA---LLPYNPA-ITNPDYVVTG
MtLYK3_-_Medtr5g086130.1  ---------LASYYIIPSI--QLRNISNFMQSKIVLTNSFDVIMSYNRDVVFDKSGLISY
g49675           --DQQSCTAFVVYRIQVGD--SMRDICDRFRQYEGVVLV---INNM-----TEATRLTVN
g50598           --DQQQCTAYVAYRLQVGD--MEADIASRFGLGR------------------DPDLMSSIGS
g44510           MGDDMNCTAFVVYRPQRGD--TFDDVVRRFAQDPVVNRY---PFRNDSGPDNDL-----
g91196           MGDDKNCTAFVAYRPRRAD--TWHDLAIRFAQDPVINPY---PILIHGGSDNDLLTL---
g30047           --DQATCTSYVVYTWQEGE--RVADVAKRFGLNLANDSA---RITI----DPDHI-----
g8984            ---DSSCQAYAVYVVSVGD--RLEDICSRFQVTFADISA---VNPV------CNFTLHPG
```

22

```
g8619                    - - H D T S C Q A Y A V Y D V R P G D - - R V Q D I S K L F D D D Y R D V A S - - - V N K L - - - - D L G N L T L L P D
MtNFP_-_Medtr5g019040.1  - - S P P S C E T Y V A Y R A Q S P N F L S L S N I S D I F N L S P L R I A K - - - A S N I - - - - E A E D K K L I P D


                                   430         440         450         460         470         480
                          . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . |
OsCERK1                  D R V L V P F P C S C L G L P A A P A S T - - F L A G A I P Y P L P L P R G G G D - - - - - - - - - - - - - - - - - - - -
MtLYK3_-_Medtr5g086130.1 T R I N V P F P C E C I G G E - - - - - - - F L G H V F E Y T T - - K - E G D - - - - - - - - - - - - - - - - - - - -
g49675                   R E L I I P V S C S C - - - - N Q S R S A Q V Y W A R P T Y P V - - - Q - V D D
g50598                   G L Q I I P V F C S C T N R S S S Q S A S G - F E T A R L T Y T V - - - Q - S G D
g44510                   - L L L F P V A C S C - - - - N E S R R - Q Y E T D L F Y T V - - - Q - S G D
g91196                   - L S L F P V P C S C - - - - N E S R R - Q F E T D I F Y T V - - - E - T G D
g30047                   - E L I F P V T C S C - - - - N R T R Q A - F R T D R I S Y V V - - E A E G G A A G G G G G G E G G E G G G G G G G G R
g8984                    Q N L L I P L S C L C - - - - N D E G - - F F Q A S I Q H V V - - - R - D D Q
g8619                    Q Q L Y I P L T C H C - - - - D E R N A - V F Y A T V N H V I - - - A F D G E
MtNFP_-_Medtr5g019040.1  Q L L L V P V T C G C - - - - - T K N - - - - H S F A N I T Y S I - - - K - Q G D


                                   490         500         510         520         530         540
                          . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . |
OsCERK1                  - - T Y D A V A A N - Y A D L T T A A W L E A T N A Y P - - P G R I P G G D G R V N V T I N C S C G - D E R V S P R Y G
MtLYK3_-_Medtr5g086130.1 - - D Y D L I A N T Y Y A S L T T V E L L K K F N S Y - - D P N H I P V K A K I N V T V I C S C G - N S Q I S K D Y G
g49675                   - - T L E K I Q N D L Y Q N L T A I S D I V E V N R L G - G Q D Y L M A G W T L R I P I P C A C N - - T S D G G S F S
g50598                   - - T V D K I A T S K Y Q S L T S T T E I G A A T G L Q - N V D V I N P G D V L R I P I P C A C N R T S A K A S Q S R
g44510                   - - T L S R I A D R E Y D S L I S I S D M V K V N G I S - V E D F I D V A W R L K I A I R C A C N - S S S A G R E F E
g91196                   - - S L S G I A Y L K Y E S L L S V A D M V Q A N G V S - D P D Y V R V G W R L K I A I R C A C N S S S S A G R D F D
g30047                   S V S V L K I A S S R Y Q S L T S A E D I M D A N G I E D V N E N S V A A G Q T L K I P I R C A C N - - V S E G E Q F Q
g8984                    - - T L H E V A S T S F H N L T T A E D I A D A S R I N - T L Q S V H A G Q S L R I P I S C G C N A S S L S W L R A P
g8619                    - - T I E S I A S V L F Q N L S H A Q E I V N A S G L T - H G V A V Q P G Q T V S V P I R C G C N S T A L L S V R T P
MtNFP_-_Medtr5g019040.1  - - N F F I L S I T S Y Q N L T N Y L E F K N F N P N L - S P T L L P L D T K V S V P L F C K C P S K N Q L N K G I K


                                   550         560         570         580         590         600
                          . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . |
OsCERK1                  L F L T Y - P L W D G E T L E S V A A Q Y G F S S P A E M E L I R R Y N P G M G G V S G K - - - - - - - - - - G I V F I
MtLYK3_-_Medtr5g086130.1 L F V T Y - P L R S D D T L A K I A T K A G L - - - - D E G L I Q N F N Q D A N F S I G S - - - - - - - - - - G I V F I
g49675                   S F V T Y V P G S T G E T A D A I A A K F Q T - - - - T V G N I R R A S M I A E E E G V V D V Y - - - - - - R P L I V
g50598                   S Y L T Y G P A S G S E K Y A T I A Q T F N S - - - T A R A I L D L N G V V G D P V T V S S T V N P S Q R P Q P L I I
g44510                   S Y L T I S S K S T S D T P S G V A Q R F N T - - - T V D A I Q R V N G L F G S D F V D F E - - - - - - G P L I V
g91196                   S Y I T Y C P N S T F E T P A D V A Q R F N T - - - T V A A I R T V N G I V G P R F D L F D R D N - - - - - R P L V I
g30047                   S Y V T Y R P R D T S D T P A S V A K K F S S - - - T V E A M Q R V N G Y V G A V V F R A N Q - - - - - - E P I I V
g8984                    S Y L T Y - V L Q Q G D T I D S V A N T F G S - - - - D N Q T I Q K S N G G A G G V L Q L R P G - - - - - - I K L I I
g8619                    S F L T Y - V V G N D T D V R S I A Q T F R S - - - - S E R L I L E A N A G A Q A D L P L D P G - - - - - - T R L V I
MtNFP_-_Medtr5g019040.1  Y L I T Y - V W Q D N D N V T L V S S K F G A S - - - Q V E M L A E N N H N F T A S T N - - - - - - - - - - R S V L I


                                   610         620         630         640         650         660
                          . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . |
OsCERK1                  P V K D P N G S Y H P L K S G G M G N S L S - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
MtLYK3_-_Medtr5g086130.1 P G R D Q N G H F F P L Y S R T G I A K G S - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g49675                   P T N S - V Y S L F S G S G G G D W G R W S - - - - - - - - - S G D G A A G G R P S S S S D Y M T R H G S E S A A Y Q S
g50598                   P T N N - S Y P R L P G G D T S S P P P L N - - - - - - - - - G G G A T N G G S P P V S G G G G T N G G S P P - - - - V
g44510                   P T H R - I F H R E P S A P D S L V S S S S I P I P M T P C V P T T V L P S L P T P V S P V W P S A V T S P P P S L T I
g91196                   P T H R - I F Q R P S A P P - - - - - - - - - - - - - - - - - - P V D A T P V M P S V A T P P A P S - - V
g30047                   P T N T W R T L R F S S S S S Q D Q D Q I F - - - - - - - - R G G G G G G P R E D -
g8984                    P S N W - T L P N V R S Q D R P P M V P R N - - - - - - - - - K A E G S E S P P L E P A S K P G S G N P S E T A H T S L
g8619                    P A Y L - S F P T R S D E G S R A R A R S S - - - - - - G G D P S L S I A T E P P A Y P R N G V R G G A S G A E D T T
MtNFP_-_Medtr5g019040.1  P V T S L P K L D Q P S S N G R K S S S Q N - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -


                                   670         680         690         700         710         720
                          . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . |
OsCERK1                  - - - - - - - - - - - - - - - - - - - - - - - - - - G G A I A G I V I A C I A I F I V A I W L I I M F Y R W Q K
MtLYK3_-_Medtr5g086130.1 - - - - - - - - - - - - - - - - - - - - - - - - A V G I A M A G I F G L L L F V I Y I Y A K Y F Q K K
g49675                   A T S S S S S S S S S P S S T R K N S R G E V G - - - - - T L G L W I W V I V G V G I L A V L A G A S A L T V R I L R E
g50598                   S G G G T N N N G S I P T D G G G - - - - - - - - - - - N Q A L S I W I G I A I A I V V G G M A M L I V F F V I W
g44510                   S M S R S Q S R P S A S T S R - - - - - - - - - - - I V P A W L W G L I A V P V V M F V A I C A L I V Y I V R Q F
g91196                   S G P R S R S S F A V S S S Q - - - - - - - - - - - - - - I A L S V L S A L V V V I F G V I C A L I V C F V R Q L
g30047                   - - - - - - - - - - - - - - - - - - - - - - - - T H Q V H E G V P P W Q V I L L A S T L F V L I P I F I A L F
g8984                    M G G S G V N A G S D P N P K S S H S D G N Q F - - - - - P I A M V V G V L A G A A V F L I A A V A A V I C L A Q R A
g8619                    S T Q A N G T T N T A S V P M D T P G K G K T G H D P A I P F E L I I G I L A G T T V V L M L A I I C A M I Y I A R R V
MtNFP_-_Medtr5g019040.1  - - - - - - - - - - - - - - - - - - - - - - - - L A L I I G I S L G S A F F I L V L T L S L V Y Y C L K M


                                   730         740         750         760         770         780
                          . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . |
OsCERK1                  F R K A T S R P S P E E T S H L D D A S Q A - - - - - - - - - - - - - -
MtLYK3_-_Medtr5g086130.1 E E E K T K L P Q T S R A F S T Q D A S G S A E Y E T S G S S G - - - - - -
g49675                   Q R P P S P P L P P S W C K D V S G S G P W A F - - - - - - - G R T L D L S V I S R S R S L S F L G
g50598                   R R R R R P - - - - - T P G P G Q V V E T S R F W H W L S K A G L L R H - - - - - S S S L G A M R Y S K S G R L - - - -
g44510                   P R E L S - - - - - - C R T C G P A F G L S Q Q C Q H Q G P A Y G I D D R - - - - D K V K G L - K A A K R P R Q H Y - -
g91196                   R L A R S - - - - - - - T G A R F G L C L R L R P S A P R - G I G D R - - - E V D - - - K A T R R P R Q C Y P L
g30047                   R R Y G W P T C S R S D D G G H S L S F S C R F L H N S S E S C G C A L S - - - - H K P L Y S T R G K Q E G D S S N - -
g8984                    G K P G G - - - - - D N W T D V D K W G W A D W R E W S R G S G R W G R F - - - S R G S G N S R G S F G G R G W R - -
g8619                    A K G T S - - - - - - V R T A A A D W E R R E W A E W S T E S A A R W Q L S V I G S R A L P V T R A S A T L P V P Q -
MtNFP_-_Medtr5g019040.1  K R L N R - - - - - - - -


                                   790         800         810         820         830         840
                          . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . |
OsCERK1                  - - - - - - - - - -
MtLYK3_-_Medtr5g086130.1 - - - - - - - - - -
g49675                   C H V A G S D Y S S G V E D S K Y I S Y G F D S Q L D R L R G G G V G G G G P G A E L S G S G V D V S S G G G G A G G G
g50598                   - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - P G G E A T T T A A G G G G G G G G G G G G G G G
g44510                   - - - T P D K R W R G V K N M D D Q A C R R K E - - - - - - - - - - - A S G C A F V N G C G G - - - - G G G K K L G K
```

23

```
g91196                  P P P P P P P P P Q R V K R G R D E E V I L H K F D C G E A A A G R A V A E V G C C E A N D V G R A G A E G G G R R G G M
g30047                  - - - - - - - - - - - - - - - - - - - D S R L W R S S G H H N V H D L N T D S D V D W G G G E G G G G G G G G
g8984                   - - - - - - - - - - - - - - - - - - G T S I P G T A V R Q Q A S S G V M R S S G S T L G G
g8619                   - - - - - - - - - - - - - - - - - - L M N T F R R R D I I H T S I S N T S S L E G F S S H G G N C A K F I
MtNFP_-_Medtr5g019040.1 - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

                                   850         860         870         880         890         900
                        . . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . |
OsCERK1                 - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
MtLYK3_-_Medtr5g086130.1 - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g49675                  G T A D R V A I V S L D S A R G K A V D N R H - - - - - - - - - - - - - - - - - - A F S V D D N R R G K V I P I E L Y
g50598                  G G G D N - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g44510                  G G G G K R G R Q G - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - V R R
g91196                  G G G G G G G G V R - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - R D H I K R H K C E T D
g30047                  G G R G G H A K G G G D A R G E V R E G G G E E E E R R G K E E K A G G V G C W P I P I T V E D T T K Q K Q K Q K E
g8984                   P A S L R T S G T A P D L L G L R G M D - - - - - - - - - - - - - - - - - - - - - - - - - - - - D V N R L G G S Q L D
g8619                   P G V D N G S G G N V - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - G R T
MtNFP_-_Medtr5g019040.1 - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - S T S

                                   910         920         930         940         950         960
                        . . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . |
OsCERK1                 - - - - - - - - - - E G I K V E R S I E F S Y E E I F N A T Q G F S M E H K I G Q G G F G S V Y Y A E L R G E K T A I
MtLYK3_-_Medtr5g086130.1 - - H A T G S A A G L T G I M V A K S T E F T Y Q E L A K A T N N F S L D N K I G Q G G F G A V Y Y A E L R G E K T A I
g49675                  S E K I P P T G G V G G V G G G G G I G T F S Y H E L A H A T N G F S P A N Q I G K G L S G P V F Y A K L R G R E V A I
g50598                  - - - - - - - W V S L D D E E E D R R V V F S L R E L E E S T D G F N E A N L I G R G G S C E V F Y A N L R G R E V A I
g44510                  E I K W Q S G K H H K T L E E D G R K F V I F S L K E L E E A T D G F S S D S K I G R G S S C E V F Y A Q L R G R E V A I
g91196                  N K V I N N L D L S E L E K E G R K F I V F S R E E L E Q A S N G F S S D N Q I G R G S S C R V F Y A Q L R D R E V A I
g30047                  K E R E K E R E V V Q L K E E A T G V S V F S L P E L E D A T D G F S N A N Q I G I G A S C Q V F Y G V L R G R E V A I
g8984                   R Q P G R G S L S A F G V E D P A R L I V F P Y K E L L K A T G R F N E V N K I G E G A Y G S V Y Y A R V R G R D A A V
g8619                   V E A S A S P H K S A A A M V S P S L S V F T Y K D L M K A T D H F S E L K K I A E G S Y G S T Y R G L L E G R D V T V
MtNFP_-_Medtr5g019040.1 S S E T A D K L L S G V S G Y V S K P T M Y E I D A I M E G T T N L S D N C K I G E - - - - S V Y K A N I D G R V L A V

                                   970         980         990         1000        1010        1020
                        . . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . |
OsCERK1                 K K M G M Q A - T Q E F L A E L K V L T H V H H L N L V R L I G Y C V E N - - C L F L V Y E F I D N G N L S Q H L Q R T
MtLYK3_-_Medtr5g086130.1 K K M D V Q A - S S E F L C E L K V L T H V H H L N L V R L I G Y C V E G - - S L F L V Y E H I D N G N L G Q Y L H G I
g49675                  K R L K V S W G P K E L R K E M K V L R K V H H K H L V E L I G F C A D S - - S L F L V F E Y C H G S L S G S L H S P
g50598                  K K M R F D Y - S R E F K K E L K I L S V H H R H L V Q L I G F C T D R - - N L F L V Y E Y C H Q G T L S S H L H E P
g44510                  K R M R L A W - R R E F Q K E L K I L S Q V H H R H L V E L I G F C T H G - - Y L L L V Y E Y C H Q G T L S S H L H S P
g91196                  K R M R L V W - R R Q F Q K E L R I L S Q V H H R H L V E L I G F C T E Q - - Y L F L V Y E Y C H Q G T L S S H L H S P
g30047                  K R L K F N Y - T R E F Q K E L K V L S R V H H R C L V G L I G Y C I E R - - S M F L V Y E Y C I E R - - S M F L V Y E Y C V E R
g8984                   K R L K S V K - E K E F Q S E L E M M C R V H H R H L L E L L G Y C T E H - - C L I L V H E Y A E G G A L K K H L H S P
g8619                   K K L K K S K - A E Q F Q S E M E V M S R V H H C H V L A L V G Y C V D R - - C L M L V H E F A E K G S L N R H L H S P
MtNFP_-_Medtr5g019040.1 K K I K K D A - - - - - S E E L K I L Q K V N H G N L V K L M G V S S D N D G N C F L V Y E Y A E N G S L E E W L F S E

                                   1030        1040        1050        1060        1070        1080
                        . . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . |
OsCERK1                 G Y A - - - - - - - - - - - - - P L S W A T R V Q I A L D S A R G L E Y L H E H V V P V Y V H R D I K S A N I L L D K
MtLYK3_-_Medtr5g086130.1 G T E - - - - - - - - - - - - - P L P W S S R V Q I A L D S A R G L E Y I H E H T V P V Y I H R D V K S A N I L I D K
g49675                  A V H G - - - - - - - - - - - - V L D W P T R A Q V A L D A A Q A L L Y I H D R V H P S Y V H G D V R S A N I L L D R
g50598                  A T H G - - - - - - - - - - - - A L E W P A R V Q I A L D A A E G L Q Y I H E H A K P S Y V H C D I K S V N I L L D Q
g44510                  D A C G G - - - - - - - - - - - P L D W R T R V Q I A L D A A Q A L L Y I H D Q V C P A Y I H C D I K S V N I L L D R
g91196                  D T C G G - - - - - - - - - - - P L D W R T R V Q I A L D A A Q A L L Y I H D R V H P A Y V H C D I K S L N I L L D R
g30047                  S S S S S S S S S S S S I T P P L P L D W P T R V Q I A I D A A K A L L Y I H D Q V Q P S Y V H C D I K S V N I L L D L
g8984                   T T N G F S - - - - - - - - - - P L P W R S R V Q I A V D V A S A L E Y L H E H T R P S Y V H R D I K S S N I L L D K
g8619                   S R N G F L - - - - - - - - - - P L A W S V R V Q I A V D I S S A L E Y L H E R S W P G F I H H S I S S R N I F L D K
MtNFP_-_Medtr5g019040.1 S S K T S N S V V - - - - - - - S L T W S Q R I T I A M D V A I G L Q Y M H E H T Y P R I I H R D I T T S N I L L G S

                                   1090        1100        1110        1120        1130        1140
                        . . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . |
OsCERK1                 D F R A K I A D F G L A K L T E V - - - - - - - - - - - - - - - - - - - - - - G S M S Q S L S T R V A G T F G Y M
MtLYK3_-_Medtr5g086130.1 N L R G K V A D F G L T K L I E V - - - - - - - - - - - - - - - - - - - - - G N S T L H T R L V G T F G Y M
g49675                  N L R G K V G D F G L N K L T T K D D L Q R G R S - - - - - - - - - - - - - S H T A G Q L R T L S G R N L A Y M
g50598                  E L R G K V A D F G V I K L T P A A N S G G D G G S L A M T G S V A I T T T A M - - V V Q D H T V S T K F A G T M G Y M
g44510                  D L R G K V A D F G V T K L I R P E G E S G G S A I S G D P R K P - - - - - - - R K K K D T V S T K V A G T W G Y M
g91196                  D L R G K V A D F G V I K L I M P E G E S L A A A G L R H S - - - - - - - - - - R K K N G T L S T K F A G T W G Y M
g30047                  Q L R G K V A D F G V M K L M R K E G E T V T P E M A T R T T G V Q S W H A T W T S M W R S T S A S T R F A G T M G Y M
g8984                   D M R A K V A D F G L I R L M G C E - - - - - - - - - - - - - - - - - - - - - E G K G S M V S T G L V G T V G Y M
g8619                   N M R A K T A N F G L S K R I D W K D D S - - - - - - - - - - - - - - - - - T L S L A V N T A E L P D V V G Y M
MtNFP_-_Medtr5g019040.1 N F K A K I A N F G M A R T S T N - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

                                   1150        1160        1170        1180        1190        1200
                        . . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . |
OsCERK1                 P P E - A R Y G E V S P K V D V Y A F G V V L Y E L L S A K Q A I V R S S E S V S E S - - - - - - - - - - - - - - - - -
MtLYK3_-_Medtr5g086130.1 P P E Y A Q Y G D V S P K I D V Y A F G V V L Y E L I T A K N A V L K T G E S V A E S - - - - - - - - - - - - - - - - -
g49675                  A P E Y L Q E G H V S T S V D V Y S L G V V L L E L L T G R N A A A F T V E Q G S S G M T L - - - - - - - - - - -
g50598                  A P E Y M I Y G R V S P M C D V Y S F G V V L L E L L T G Q K A I L P S E P R V - - - - - - - - - - - -
g44510                  A P E Y M I Y G E V T P S S D V Y S F G V V L L E L I A G Q R A I L - - - - - - P R E R V A - - - - - - - - - -
g91196                  A P E Y M I Y G E V T P S C D V Y S F G V V L L E L I T G Q R A I L - - - - - - P R E R L A - - - - - - - - - -
g30047                  A P E Y M I Y G H V S A S C D V Y S F G V V L L E L I T G R K A V L S P A E E R E R E K V V - - - - - - - - - - - -
g8984                   A P E Y L K W G H V S T K A D V Y S F G V V L L E L L S G Q E A L S S T I D S H S S T V Q I I G G G G W R H H R Q S D H
g8619                   A P E C L K L G Q V S T K A D V Y S F G V V L E L L S G Q E A V S C A - - - - A G S H I V Q L S G G L H C E S W A L D
MtNFP_-_Medtr5g019040.1 - - - - - - - S M M P K I D V F A F G V V L I E L L T G K K A M T T K - - - - E N G E V V

                                   1210        1220        1230        1240        1250        1260
                        . . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . | . . . |
OsCERK1                 - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
MtLYK3_-_Medtr5g086130.1 - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
```

24

```
g49675                     ------------------------------------------------------------
g50598                     ------------------------------------------------------------
g44510                     ------------------------------------------------------------
g91196                     ------------------------------------------------------------
g30047                     ------------------------------------------------------------
g8984                      IGGSGSMVVSETSGSFLESSGGGSESQSVSESFANVLVKRKLEKGKPEREREPPSASGAV
g8619                      YNGSSSTVVAETSGSFAESDFGSTPPDSSSELRAKRPV---IQKRVLDVEIRPYKTAKTR
MtNFP_-_Medtr5g019040.1    ------------------------------------------------------------

                                    1270       1280       1290       1300       1310       1320
                           ....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1                    ------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1   ------------------------------------------------------------
g49675                     ------------------------------------------------------------
g50598                     ------------------------------------------------------------
g44510                     ------------------------------------------------------------
g91196                     ------------------------------------------------------------
g30047                     ------------------------------------------------------------
g8984                      VAAQPTTPFENDPRNVAQHESPSGAASRWGNSDGSLGSPSDGTGAVANAVPSCRSDSSPN
g8619                      LNRMGKRQGEGTKQHPPGNEGCENSNRYLENTTNSGGKKEDA------------------
MtNFP_-_Medtr5g019040.1    ------------------------------------------------------------

                                    1330       1340       1350       1360       1370       1380
                           ....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1                    ------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1   ------------------------------------------------------------
g49675                     --------------------------IESSCTSSDCSYSTFSYSQGADFASL---------
g50598                     --------------------------QSPNLDLGPDAVPRPAYRERAPFR----------
g44510                     --------------------------AAAASELAPQPPSIP------------------
g91196                     --------------------------AAENSELATHVTQPPTLP---------------
g30047                     --------------------------SSPSSAPPPPPPPPPPPPPTTPSPPPPPPTPPAP
g8984                      PTGKILSVNGKMLRLHGLMVEIVPEGDTFACKPVRVQQQQPVTSVTGCGEPDPVDDDPST
g8619                      -CGKILSVNGSKLRLRGLTVEVLADLSGPSSCKPVLSSSVENSASAIPTQEEPANSTTGS
MtNFP_-_Medtr5g019040.1    ------------------------------------------------------------

                                    1390       1400       1410       1420       1430       1440
                           ....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1                    ------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1   ------------------------------------------------------------
g49675                     ------------------------------------------------------------
g50598                     ------------------------------------------------------------
g44510                     ------------------------------------------------------------
g91196                     ------------------------------------------------------------
g30047                     TATPTPPAAPFNAGLPLPAR----------------------------------------
g8984                      TTPLMPEGSPSSDSKSGSTSRRTRKGGFLMGDKGVGKGSSSQGSGKGREKAGNVEVQGEN
g8619                      MIAVSVPSSDP-----------------------LPPCPHSVQASANDGCAVPTKGKAKSVA
MtNFP_-_Medtr5g019040.1    ------------------------------------------------------------

                                    1450       1460       1470       1480       1490       1500
                           ....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1                    ------------------------------------------------------------
MtLYK3_-_Medtr5g086130.1   ------------------------------------------------------------
g49675                     ------------------------------------------------------------
g50598                     ------------------------------------------------------------
g44510                     ------------------------------------------------------------
g91196                     ------------------------------------------------------------
g30047                     ------------------------------------------------------------
g8984                      SDAAAGGEARGEIQEEKSSKLSKSNSSRSKRTARGSTSQSAHSKSKGGLRKTRGRHGKQE
g8619                      AFSKHGKEVKIGIGKPRKADIAAGSNKKPVHALAMTSSHRKGPRSRKVAGKKYENKRRGD
MtNFP_-_Medtr5g019040.1    ------------------------------------------------------------

                                    1510       1520       1530       1540       1550       1560
                           ....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1                    -----------------------KGLVFLFEEALSAPNPTEALDELIDPSL-QGDYP
MtLYK3_-_Medtr5g086130.1   -----------------------KGLVQLFEEALHRMDPLEGLRKLVDPRL-KENYP
g49675                     -----------------------PSLPAALVLLIKQMKSPKDVKAIADPRLSESGYP
g50598                     -----------------------PFLTEVMVPRIKDMKAPAQVCPIVDPQL-GSDYP
g44510                     -----------------------KTLASSVAPLIKQMTSAEEVSVIVDPQL-GIAYP
g91196                     -----------------------ATLASSVAPLIKQMTAAEDVSAIVDPRL-DIAYP
g30047                     -----------------------PFLTEVMCPLINEMTCPGHVAAIIDRQL-GNGYP
g8984                      KAVDFSSSLDVGASVGRSWPRGVVQRKSLAEWIVPAFKCLGSPIDIVNLVDPDL-RDQYP
g8619                      GVKLLGAATEGTSPILRNSAVQRMLRCSLVSWIVPAIRSLGSPADVVCLVDPDM-RGQYP
MtNFP_-_Medtr5g019040.1    -----------------------ILWKDFWKIFDLEGNREERLRKWMDPKL-ESFYP

                                    1570       1580       1590       1600       1610       1620
                           ....|....|....|....|....|....|....|....|....|....|....|....|
OsCERK1                    VDSALKIASLAKSCTHEEPGMRPTMRSVVVALMALTANTDLRDMDYHPF-----------
MtLYK3_-_Medtr5g086130.1   IDSVLKMAQLGRACTRDNPLLRPSMRSIVVALMTLSSPTEDCDDDSSYENQSLINLLSTR
g49675                     KEAALQLARLAVECVSDKPQKRPQMKRVVYVMEDILAMSRPRHQPSADT--------AD
g50598                     QDLALQIGRLAMQCVQHEPELRPQMRRVVYLLDEILADTRRRTHQPSSA-------PSSD
g44510                     RDLTFELARVAAQCVELEPEDRPNMRRILVYVLDETLETYLMRSLHADYE-------ISID
g91196                     TDLTLQLARLAAQCVELEPDDRPEMKRVVYVLDETLSLAISRVCSSNTDDKMLAVRIAQD
g30047                     QESALQVAKLAVQCVATDPDQRPKMKQVVYVLELVRNRSAASGR---------------
g8984                      IEGVIKMAEVAVRCVQENPEARPDMKRVAYELDELMLLTQRWESAGRTT-----------
g8619                      PEVVAKMADLAVRCVQENPQARPNMSRVAYELDEILLLTRRWEALQCKAQPENHSRVVVR
MtNFP_-_Medtr5g019040.1    IDNALSLASLAVNCTADKSLSRPTIAEIVLCLSLLNQPSSEPMLERSLT-------SGLD

                                    1630       1640       1650       1660
```

```
                                    . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . . | . . . .
OsCERK1                             - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
MtLYK3_-_Medtr5g086130.1            * - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g49675                              P D P Y H V A F P R P * - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g50598                              G S S H H Q H T G S S T G G F L * - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g44510                              S V P V - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g91196                              S I P V - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g30047                              - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
g8984                               - - - P G G A D E E V L G A Q - - - - - - - - - - - - - - - - - - - - - - - - - - - L L R *
g8619                               L Q G S R E S E K S E V S G M G Q Q N V V L W A S T V G P D V S T D G E G S I S R P Q L Q R
MtNFP_-_Medtr5g019040.1             A E A T H V V T S V I A R * - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
```

**Data S1M: Alignment of the Chara LysM-RLK with representative LYK and LYR from angiosperms.**

CXnC motif representative of the LysM domain are present in the predicted extracellular domain of the Chara LysM-RLK (Positions: 357/367; 429/431 and 527/529). The activation loop missing in LYR proteins such as Medicago NFP is present in the *C. braunii* sequences (Position 1123 – 1149).

**Data S1N: Midpoint-rooted Maximum Likelihood tree of the LysM-RLK family.**

The Charales sequences form a single cluster, encompassing seven *C. braunii* sequences (green). For putative relationship to land plant genes see text and Fig. M13.

**Data S1O: Outgroup-rooted phylogenetic analysis of *C. braunii* class III peroxidases**

The CIII Prx protein sequences from *K. nitens* (3 sequences), *C. braunii* (14 sequences), *P. patens* (57 sequences) and *A. thaliana* (73 sequences) as well as 2 APx sequences from *C. braunii* and *A. thaliana* (outgroup) were aligned using MAFFT and the tree constructed using Maximum Likelihood. Evolutionary analyses were conducted in MEGA7. All protein sequences are available using the CHBRA# or the CbraPrx# (http://peroxibase.toulouse.inra.fr). The collapsed triangle contains all *P. patens* and *A. thaliana* sequences, as well as one of the three *K. nitens* sequences. The other two *K. nitens* sequences are marked by a blue box, the *C. braunii* sequences by a green box.

| LRR FBP | Modeled (%) | Confidence (%) | Multi-template and ab initio modeling | | |
| --- | --- | --- | --- | --- | --- |
| | | | ab initio (# a.a.) | templates | i.d. (%) |
| CHABRA g10791_t1 | 87 | >90 | 88 | c3ogmB | 17 |
| | | | | c2p1nE | 18 |
| | | | | c3oglD | 17 |
| | | | | c5hywA | 18 |
| CHABRA g11992_t1 | 79 | >90 | 140 | c3ogmB | 21 |
| | | | | c2p1nE | 21 |
| | | | | c3oglD | 21 |
| CHABRA g23992_t1 | 97 | >90 | 253 | c3ogmB | 22 |
| | | | | c2p1nE | 20 |
| | | | | c3oglD | 22 |
| | | | | c5hywA | 19 |
| | | | | c4lxrA | 11 |
| CHABRA g23719_t1 | 57 | >90 | 528 | c2p1nE | 25 |
| | | | | c3ogmB | 19 |
| | | | | c3oglD | 19 |
| | | | | c5hywA | 21 |

The overall confidence in the final model (57% at >90% confidence) was considered too low (<70%) for submission to 3DLigandSite

**Data S1P: *In silico* modeling of *C. braunii* LRR FBPs.**

Leucine-Rich-Repeat (LRR)-containing F-Box Proteins (FBPs) of *C. braunii* with sequence similarity to land plant LRR FBPs were *in silico* modeled using "intensive" modeling mode in Protein Homology/analogY Recognition Engine V 2.0 (Phyre2). The final models (color-coded by the confidence of the match to the templates overall) were submitted to 3DLigandSite server to predict potential binding sites (gray structures cartoon depiction). % i.d. percentage identity. See STAR Methods for details.

**Data S1Q: Phylogenetic tree of cyclin proteins.**

Midpoint-rooted Bayesian inference phylogenetic tree, numbers at the nodes are posterior probabilities. The *A. thaliana* CYCB sequences are shown in red, the *C. braunii* sequences belonging to that cluster in green, and th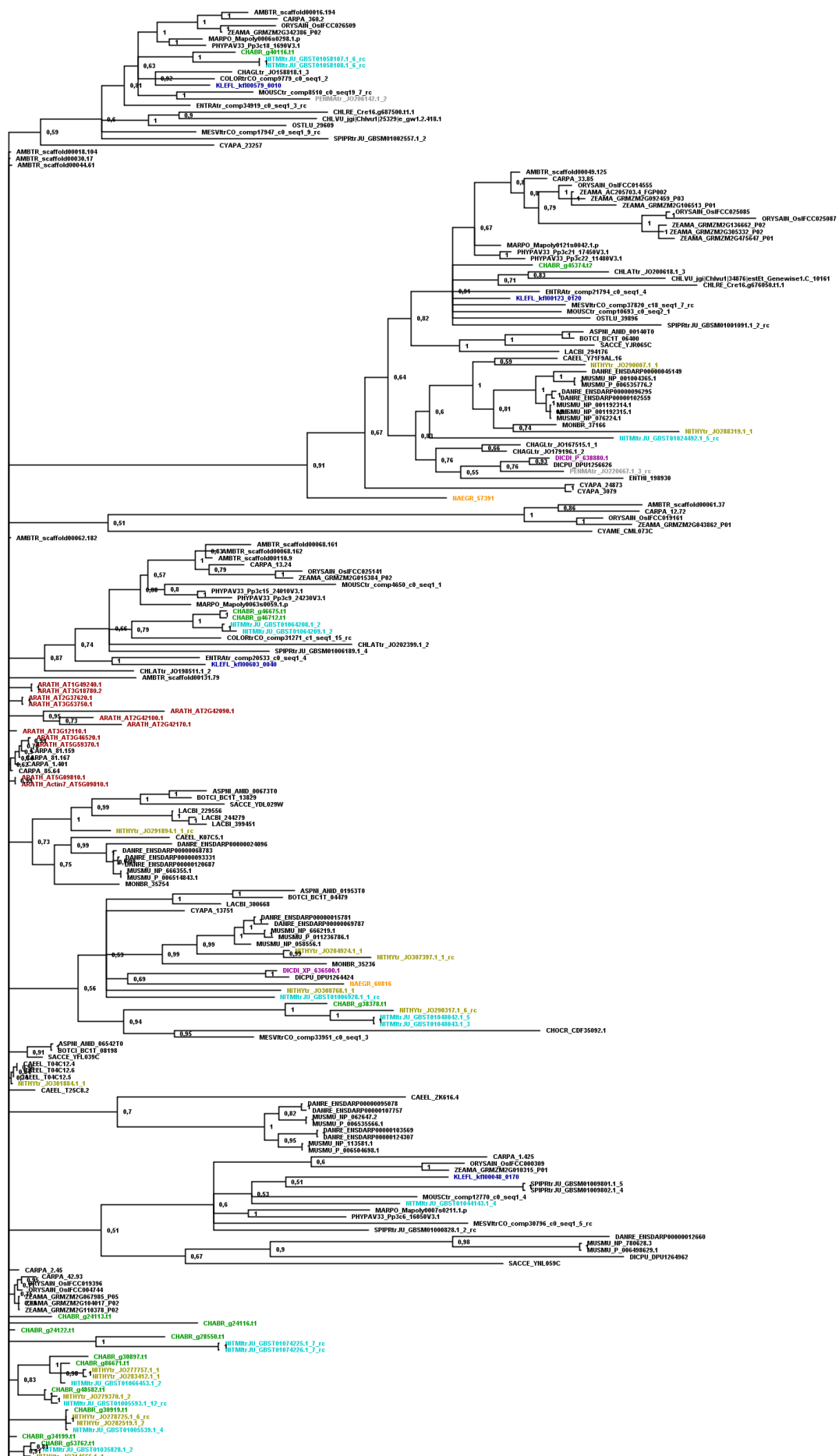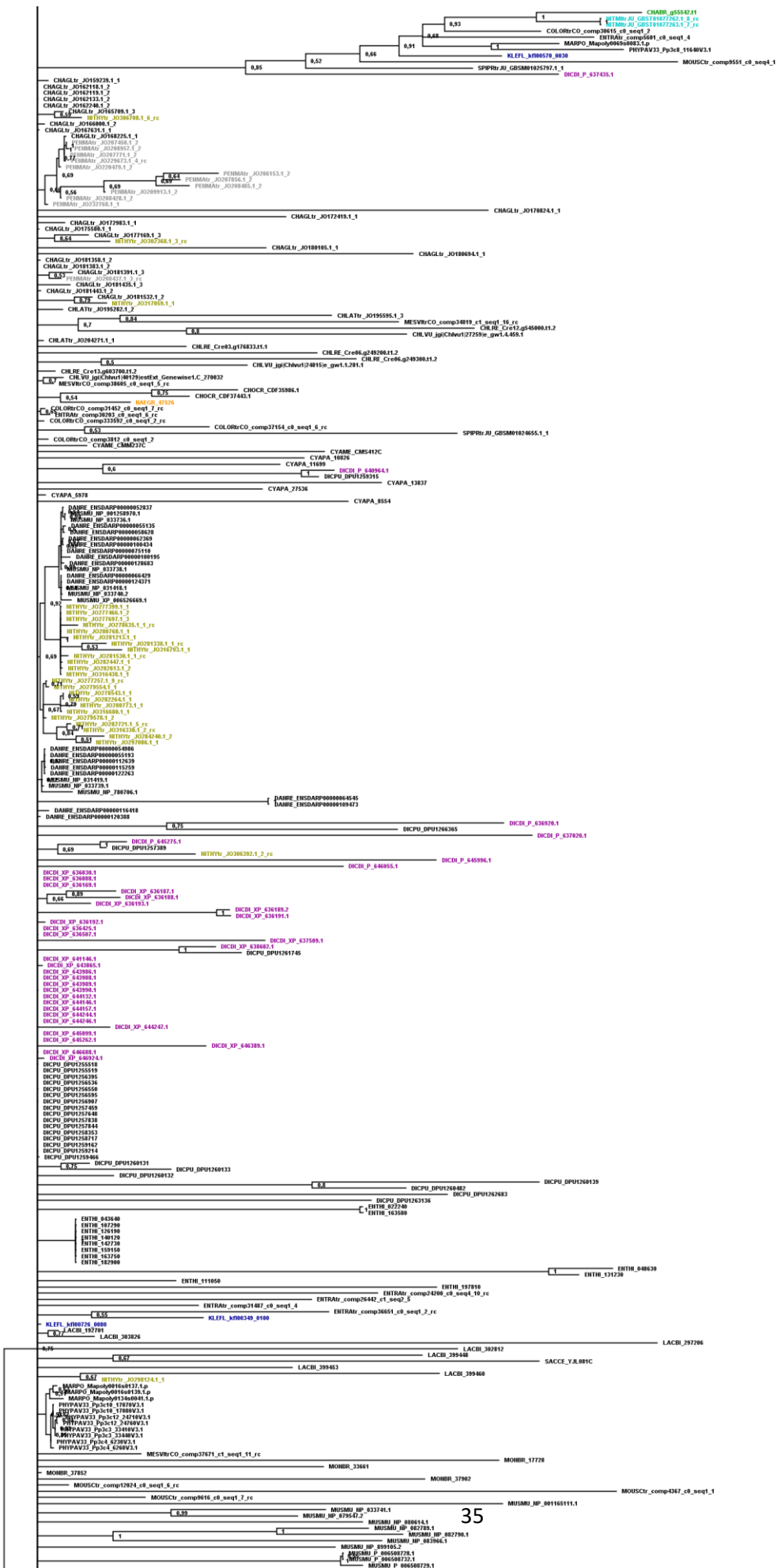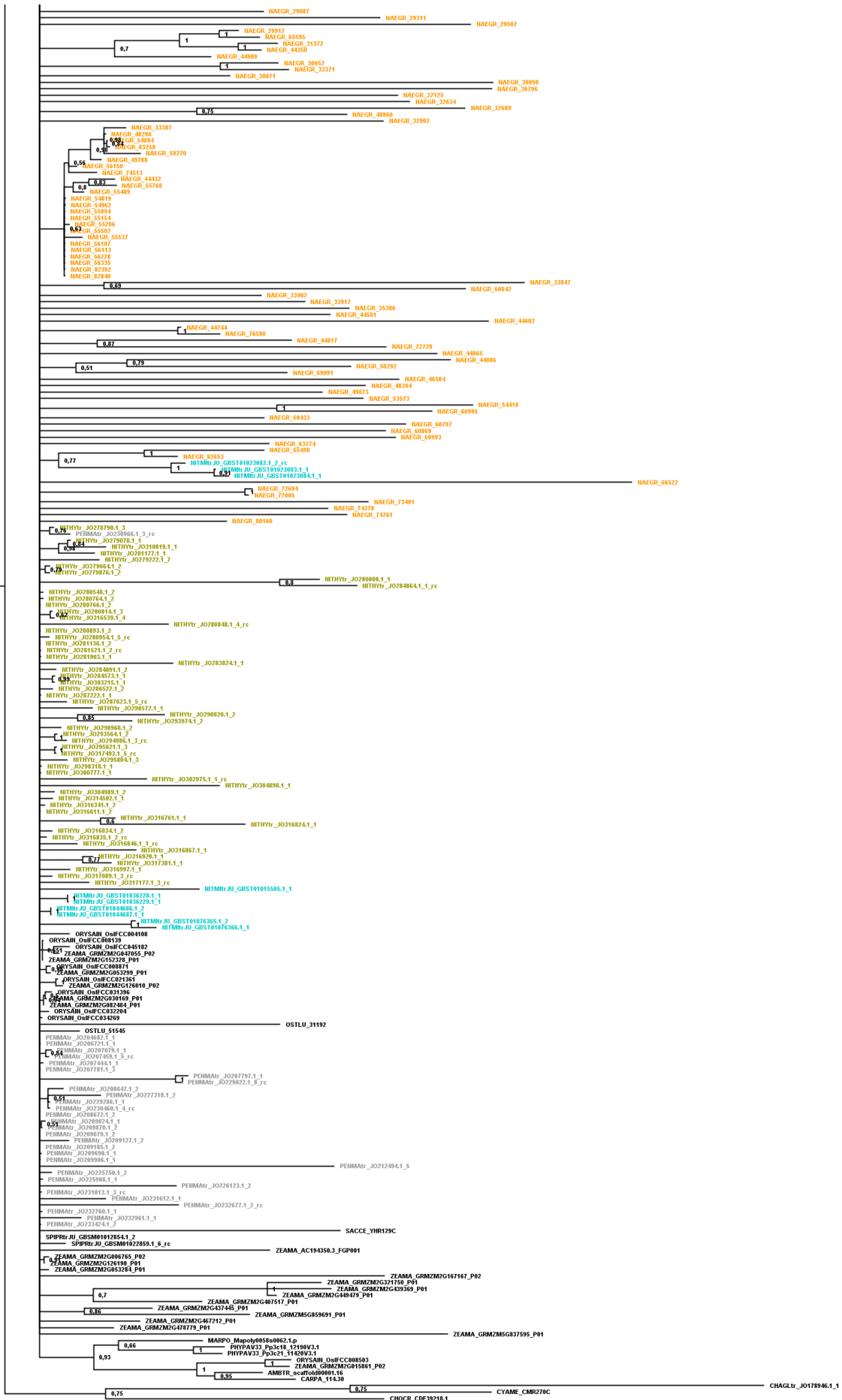e *K. nitens* sequence in blue. The clade containing the Archaeplastida CYCA sequences (one each of *C. braunii* and *K. nitens*) is collapsed.

30

**Data S1R: Phylogenetic tree of EXOCYST 70A proteins.**

Midpoint-rooted Bayesian inference phylogenetic tree, numbers at the nodes are posterior probabilities. Two clades containing exclusively seed plant species are collapsed. The *A. thaliana* EXO70A sequence is shown in red, the *C. braunii* sequences belonging to that cluster in green, and the *K. nitens* sequence in blue.

**Data S1S: Phylogenetic tree of NPSN proteins.**

Midpoint-rooted Bayesian inference phylogenetic tree, numbers at the nodes are posterior probabilities. The *A. thaliana* NPSN12/12/13 sequences are shown in red, the *C. braunii* sequences belonging to that cluster in green, and the *K. nitens* sequence in blue.

**Data S1T: Phylogenetic overview tree of canonical actin proteins.**
Midpoint-rooted Bayesian inference phylogenetic tree, numbers at the nodes are posterior probabilities. The *A. thaliana* actin sequences are shown in red, the *C. braunii* sequences in green, and the *K. nitens* sequences in blue. The sequences of *Naegleria gruberi* are shown in orange, those of *Dictyostelium discoideum* in purple. The transcriptomic sequences of *Nitella mirabilis* (cyan), *Nitella hyalina* (yellow) and *Penium margaritaceum* (grey) are also color-coded. See text for numbers of genes. See M21 for expanded tree.

34

35

36

**Data S1U: Phylogenetic tree of canonical actin proteins.**
Midpoint-rooted Bayesian inference phylogenetic tree, numbers at the nodes are posterior probabilities. The *A. thaliana* actin sequences are shown in red, the *C. braunii* sequences in green, and the *K. nitens* sequences in blue. The sequences of *Naegleria gruberi* are shown in orange, those of *Dictyostelium discoideum* in purple. The transcriptomic sequences of *Nitella mirabilis* (cyan), *Nitella hyalina* (yellow) and *Penium margaritaceum* (grey) are also color-coded. See text for numbers of genes. In the case of *A. thaliana*, the usage of the ACTIN7 query for the blast approach (*cf.* STAR Methods) resulted in recovery of all canonical actins according to TAIR, namely ACTIN 1, 2, 3, 4, 7, 8, 9, 11 and 12, and of only two additional, actin-like, proteins.

Supplementary Figures S1-S7 for Nishiyama et al.

## The *Chara* genome: secondary complexity and implications for plant terrestrialization
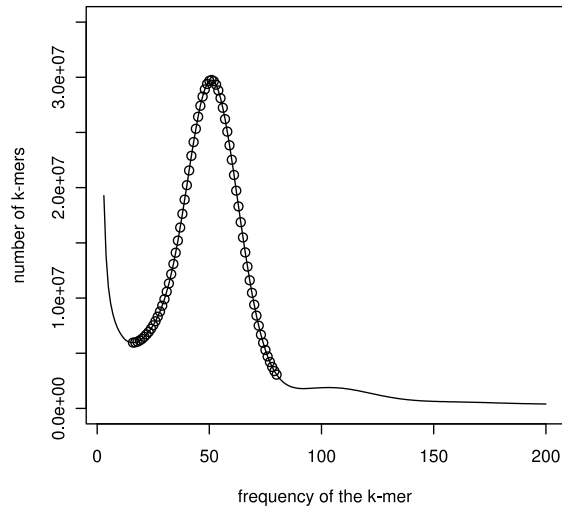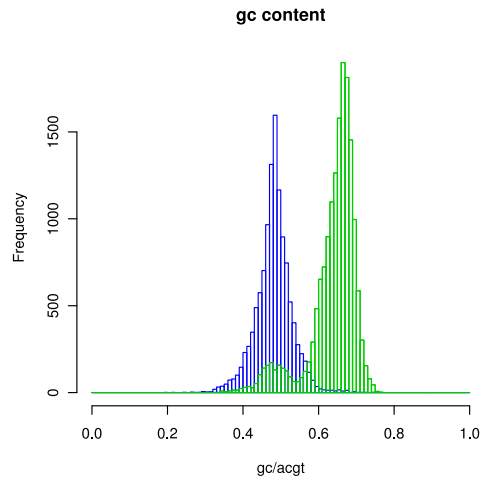


**Figure S1, related to STAR methods: Chromosomes in an antheridial filament of *C. braunii* (n=14, strain S276).**

The chromosomes during cell division in young antheridial filaments of strain S276 were observed after Feulgen staining. The chromosome number n=14 was confirmed by counts made on chromosomes during metaphase or anaphase. Most *Chara* species have either n=14 or n=28 chromosomes, Nitella and the other genera have different base numbers. There are numerous examples of monoicous/dioicous species pairs in the family, with the dioicous species always displaying half the number of chromosomes than their monoicous counterpart. For Chara typically dioicous=14, monoicous=28 (or other multiples of 14). *C. braunii* is monoicous, but is unique in having the dioicous chromosome number of 14. There are no known dioicous sister taxa to *C. braunii*, perhaps due to the already reduced genome. Scale bar = 2 μm.
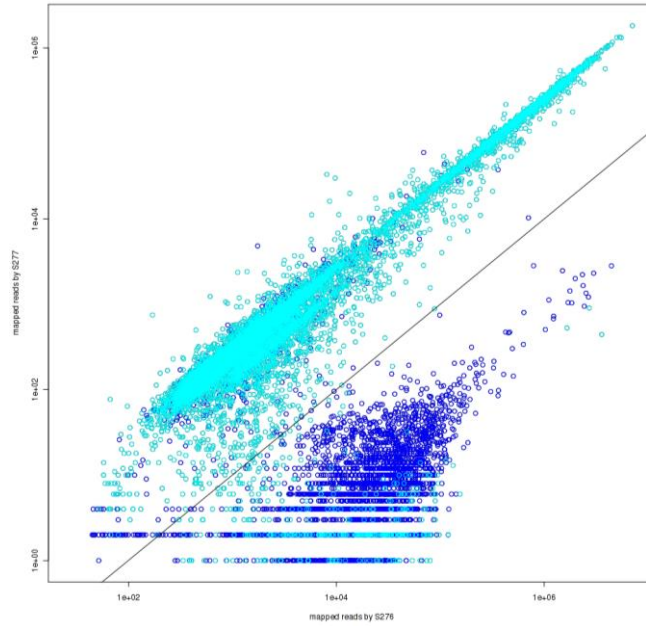
A



C



B



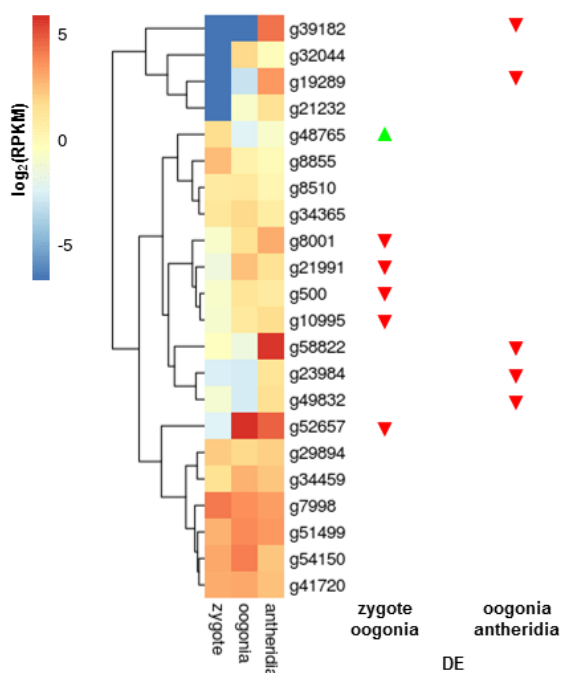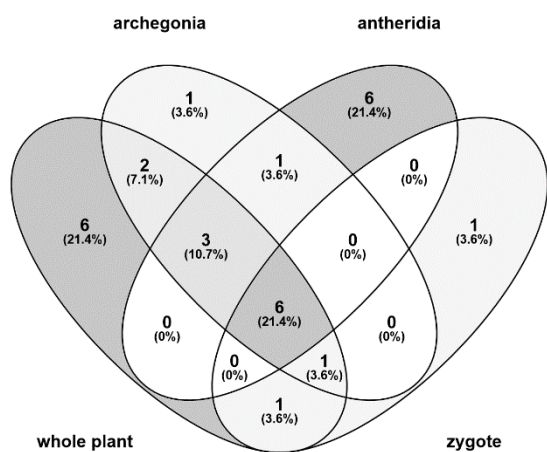**Figure S2, related to STAR methods: Assembly characteristics and decontamination**

A) k-mer frequency analysis of the S276 paired end read data with k=25. Number of 25-mers at frequency 3 to 200 are shown with the solid line. Circles shows the points from 16 to 80 as what was recognized the major peak, presumably representing the single copy region in *C. braunii*.

B) Scatter plot of mapped reads of two *C. braunii* strains on each scaffold. Blue and light blue points are scaffolds with GC content of at least 55% and less than 55%, respectively.

C) Frequency distribution of scaffold wise GC content compared between putative *C. braunii* derived scaffolds (blue) and other scaffolds (green).

**Figure S3, related to STAR methods: Ks-based analysis of *C. braunii* paralogs**

Paranome-based WGD signature prediction. (A) Ks frequency plot highlighting mixture model components mean and standard-deviation (top: #component, bottom: mean Ks) based on raw Ks value classification. (B) Ks frequency plot highlighting mixture model components mean and standard-deviation (top: #component, bottom: mean Ks) based on log-transformed Ks value classification. (C) Ks group assignment for raw Ks classification. (D) Ks group assignment for log-transformed Ks classification. (E) Significant zero crossing (SiZer) plot. (F) Significant convexity (SiCon) plot. (G-J) Significant features of kernel density estimates using indicated bandwidths, highlighting significant gradient regions in blue and significant curvature regions in green using a significance level of 0.05. Red vertical lines represent Ks value of 0.1 and 2.0, dotted red vertical line represents Ks value of 0.235 corresponding to 12.5 Ma ago (these events might be no WGDs but only more or less recent local duplication events). For *C. braunii* no single predicted WGD signature was supported by three different bandwidth kernel densities (cf. STAR Methods).
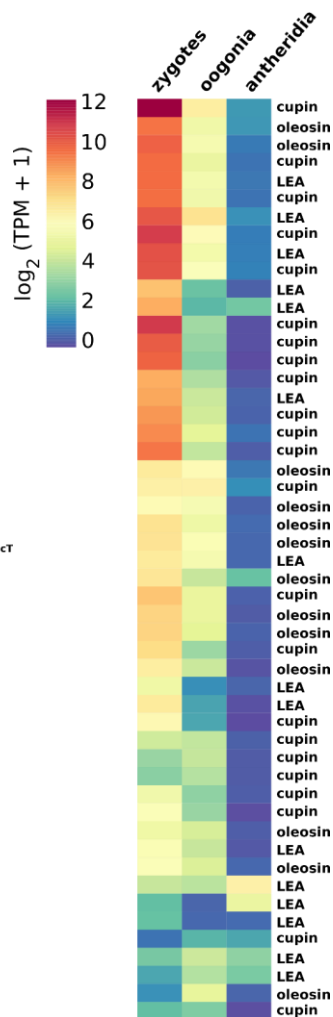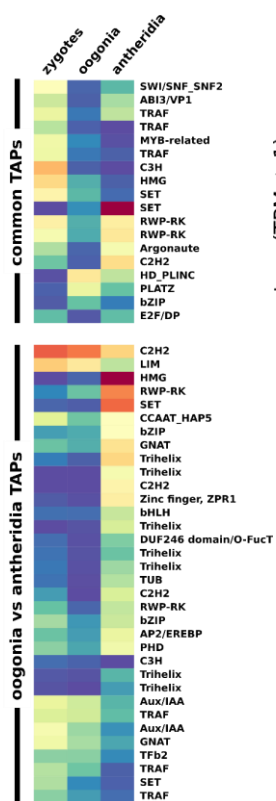
**Figure S4, related to Figure 6: Expression profiles during sexual reproduction.**

Expression profile of trihelix TF genes based on RNA-seq evidence (supplemental file 3) was visualized as A) a Venn diagram using venny (http://bioinfogp.cnb.csic.es/tools/venny/) and B) as a heatmap showing gene expression and DEGs from reproductive organs with RPKM > 1 in minimum two samples. C) Shows expression of differentially expressed TFs/TRs during sexual reproduction. D) Expression of DEGs associated with seeds during sexual reproduction. Transcripts per million (TPM) were transformed to log2 scale and clustered using the euclidean distance method and the complete clustering method (B, C, D).

**Figure S5, related to Figure 5: Exon-intron structure comparison of MIKC^C-type, MIKC*-type and charophyte MIKC-type genes.**

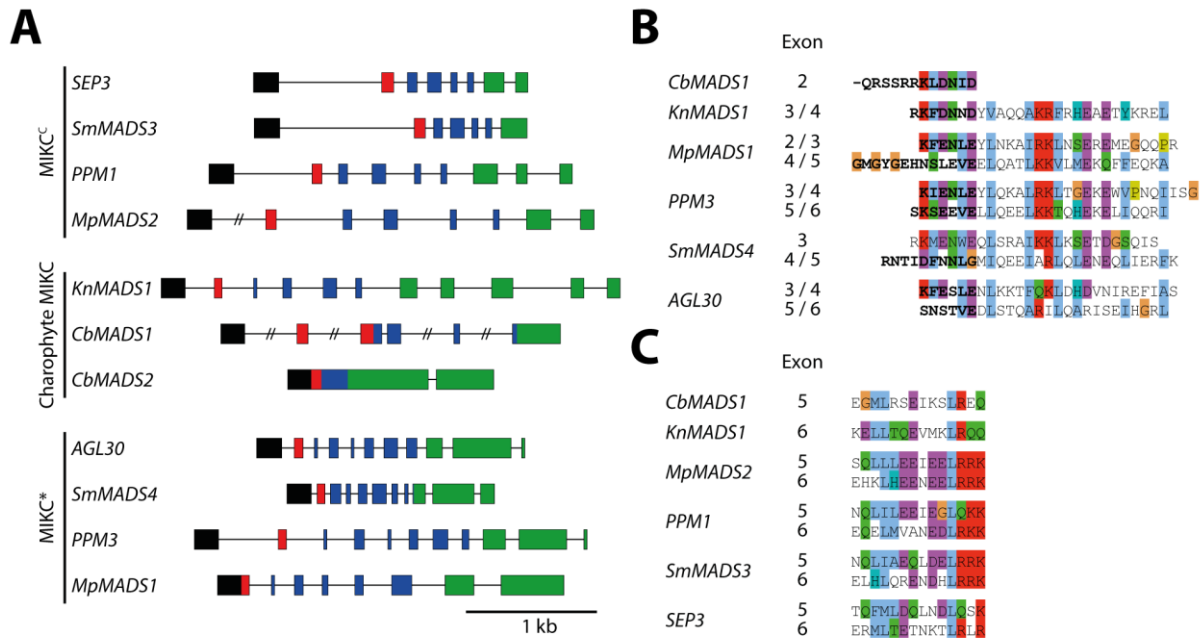(A) Exon-intron structures of representatives of MIKC^C-type and MIKC*-type genes together with the charophyte MIKC-type genes *CbMADS1*, *CbMADS2* and *KnMADS1*. The exons encoding MADS-, I- , K- and C-domains are color coded in black, red, blue and green, respectively. Among the three Type II genes that were identified in the *C. braunii* genome only *CbMADS1* shows a canonical MIKC-type gene sequence. In contrast *CbMADS2* lacks most (but not all) introns and thus probably evolved via a retrotransposition and recombination event. *CbMADS3* lacks the conserved K-box that encodes for the protein-protein interacting K-domain (data not shown). (B and C) Analysis of exon-intron structures suggest that *CbMADS1* directly descends from an ancestral MIKC-type gene that was a common ancestor of MIKC^C- and MIKC*-type genes. (B) It was previously suggested that the N-terminal part of the K-domain of MIKC*-type proteins evolved through a duplication of two K-domain exons of an ancestral MIKC-type gene (Kwantes et al., 2011). The aligned amino acid sequences encoded by exon 2 of *CbMADS1*, and by the first K-domain exons of *KnMADS1*, *MpMADS1*, *PPM3*, *SmMADS4* and *AGL30* indeed strongly support this hypothesis. (C) In addition, striking similarities between the aligned amino acid sequences encoded by exon 5 of *CbMADS1*, exon 6 of *KnMADS1* and exons 5 and 6 of *MpMADS2*, *PPM1*, *SmMADS3* and *SEP3*, respectively, suggest that also the K-domain of MIKC^C-type proteins evolved through an exon duplication of an ancestral MIKC-type gene. This is especially intriguing considering the fact that the last two K-domain exons of most if not all MIKC^C-type genes encode for a protein-protein interaction interface that facilitates tetramer formation of MIKC^C-type proteins (Theißen et al., 2016). It has already been suggested that the ability of MIKC^C-type proteins to tetramerize was an important precondition to evolve and diversify efficient developmental switches that

facilitated the transition to land and the evolution of complex body plans of land plants (Theißen et al., 2016). Thus it is tempting to speculate that an exon duplication of an ancestral MIKC$^C$-type gene in the MRCA of extant land plants created the molecular prerequisites for this evolutionary novelty.

A

**DNA integration**

microtubule-based movement

protein glycosylation

B

malate metabolic process
biosynthetic process
**oxidation-reduction process**
carbohydrate metabolic process
cellular iron ion homeostasis

C

transmembrane transport cell wall modification
photosynthetic electron transport chain
carbohydrate metabolic process
biosynthetic process
response to oxidative stress
protein phosphorylation
photosynthesis, light reaction
cellulose biosynthetic process
**oxidation-reduction process**
**photosynthesis**
lipid metabolic process
terpenoid biosynthetic process

D

base-excision repair
tricarboxylic acid cycle
mismatch repair embryo development
viral RNA genome replication cellular iron ion homeostasis
single organism reproductive process
regulation of transcription, DNA-templat...
bacterial-type flagellum-dependent cell ...
**transcription, DNA-templated**
carboxylic acid metabolic process
**protein phosphorylation** oxidation-reduction process superoxide metabolic process
glutamate biosynthetic process
chitin catabolic process
**microtubule-based movement** endocytosis
cofactor biosynthetic process
glycolytic process phosphorelay signal transduction system
cellular aromatic compound metabolic pro...
transmembrane transport post-embryonic development

**Figure S6, related to Figures 5/6: Transcriptome analyses of reproduction and early development.**
GO enrichment word clouds (category biological process); genes down-regulated (A) or up-regulated
(B) in oogonia as compared to antheridia, genes down-regulated (C) or up-regulated (D) in zygotes as
compared to oogonia. Antheridia are strongly enriched with the GO category GO:0015074 "DNA
integration" (A). 349 gene models expressed in antheridia were classified in this category; of these, 324
genes were found to be overlapping with a TE to at least 50 % (supplemental file 3). Most of these genes
were annotated as "integrase", "ribonuclease H-like", "reverse transcriptase", and "aspartyl protease"
by homology-based approach, terms typical of Ty3/Gypsy pol gene composition (Havecker et al., 2004).
Ty3/Gypsy elements represent 20 % of the *C. braunii* genome. These results might indicate mobilization

of retrotransposons and other mobile elements during male gametogenesis. This could be a consequence of genome rearrangement during male gamete formation. One could also imagine that mobilization and integration of retrotransposons might enhance genomic diversity during sexual reproduction.
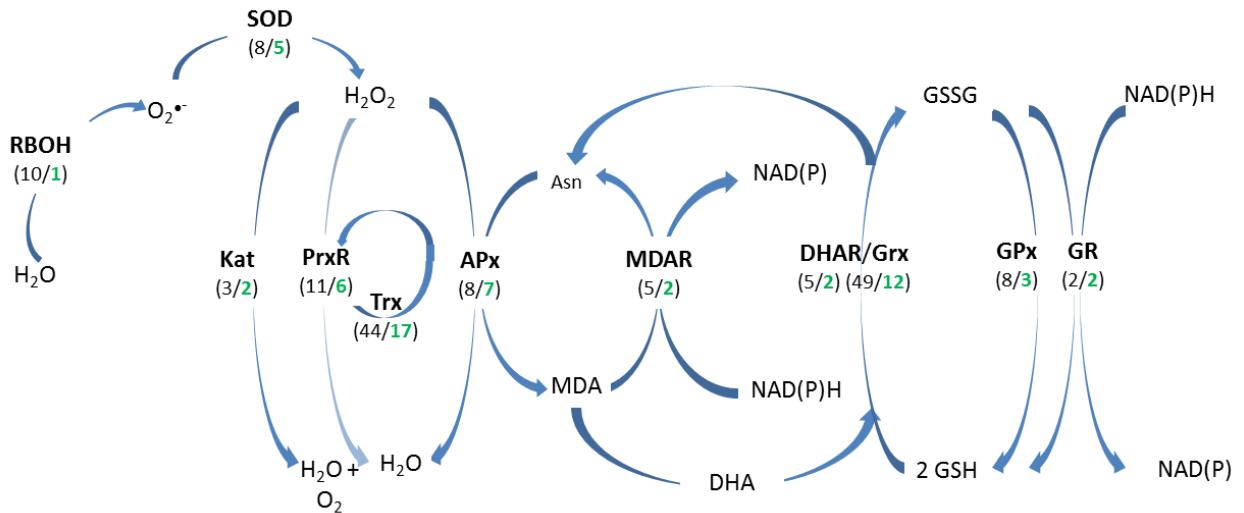
**Figure S7, related to Figure 6: Major reactive oxygen species scavenging pathway in plants.**
Proteins associated with ROS scavenging are in bold. Number of genes found for *A. thaliana* and *C. braunii* (in green) are indicated in brackets. APx: Ascorbate peroxidase, Asn: ascorbate, DHA: Dehydroascorbate, DHAR: Dehydroascorbate reductase, GPx: Plant glutathione peroxidase, GR: Glutathione reductase, Grx: Glutaredoxins superfamily, GSH: reduced glutathione, GSSH: oxidized glutathione. Kat: Catalase, MDAR: Monodehydroascorbate reductase, PrxR: Peroxiredoxins family, RBOH: Respiratory burst oxidase homolog also called NADPH oxidase, SOD: Superoxide dismutase, Trx: Thioredoxins, MDA: Monodehydroascorbate, adapted from (Inupakutika et al., 2016).