

Underwater CAM photosynthesis elucidated by *Isoetes* genome

David Wickell^{1,2}, Li-Yaung Kuo³, Hsiao-Pei Yang², Amra Dhabalia Ashok⁴, Iker Irisarri^{4,5}, Armin Dadras⁴, Sophie de Vries⁴, Jan de Vries^{4,5,6}, Yao-Moan Huang⁷, Zheng Li⁸, Michael S. Barker⁹, Nolan T. Hartwick¹⁰, Todd P. Michael¹⁰✉ & Fay-Wei Li^{1,2}✉

To conserve water in arid environments, numerous plant lineages have independently evolved Crassulacean Acid Metabolism (CAM). Interestingly, *Isoetes*, an aquatic lycophyte, can also perform CAM as an adaptation to low CO₂ availability underwater. However, little is known about the evolution of CAM in aquatic plants and the lack of genomic data has hindered comparison between aquatic and terrestrial CAM. Here, we investigate underwater CAM in *Isoetes taiwanensis* by generating a high-quality genome assembly and RNA-seq time course. Despite broad similarities between CAM in *Isoetes* and terrestrial angiosperms, we identify several key differences. Notably, *Isoetes* may have recruited the lesser-known ‘bacterial-type’ PEPC, along with the ‘plant-type’ exclusively used in other CAM and C4 plants for carboxylation of PEP. Furthermore, we find that circadian control of key CAM pathway genes has diverged considerably in *Isoetes* relative to flowering plants. This suggests the existence of more evolutionary paths to CAM than previously recognized.

¹Plant Biology Section, School of Integrative Plant Science, Cornell University, Ithaca, NY, USA. ²Boyce Thompson Institute, Ithaca, NY, USA. ³Institute of Molecular & Cellular Biology, National Tsing Hua University, Hsinchu, Taiwan. ⁴Department of Applied Bioinformatics, Institute for Microbiology and Genetics, University of Goettingen, Goettingen, Germany. ⁵Campus Institute Data Science, University of Goettingen, Goettingen, Germany. ⁶Department of Applied Bioinformatics, Goettingen Center for Molecular Biosciences, University of Goettingen, Goettingen, Germany. ⁷Taiwan Forestry Research Institute, Taipei, Taiwan. ⁸Department of Integrative Biology, University of Texas at Austin, Austin, TX, USA. ⁹Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ, USA. ¹⁰The Molecular and Cellular Biology Laboratory, The Salk Institute for Biological Studies, La Jolla, CA, USA. ✉email: tmichael@salk.edu; fl329@cornell.edu

Isoetes, commonly known as quillworts, is the only genus in the lycophyte order Isoetales, containing roughly 250 described species¹. It is the last remaining member of an ancient lineage with a fossil record that dates back to at least the late Devonian. As such, quillworts are believed to represent the closest living relatives of the giant, tree-like lycopsids such as *Sigillaria* and *Lepidodendron* that dominated the terrestrial landscape during the Carboniferous². However, in contrast to its arborescent ancestors, modern *Isoetes* species are diminutive and mostly aquatic with the vast majority of species growing completely or partially submerged. Underwater, *Isoetes* can conduct CAM³, a carbon concentrating mechanism involving the separation of carbon uptake and fixation in a time of day (TOD) fashion, with carbon being sequestered as malate at night, to be fed into the Calvin cycle during the day. CAM is a common strategy to improve water-use efficiency among xeric-adapted plants, allowing them to keep their stomata closed during the day. However, its prevalence in aquatic species of *Isoetes*³, as well as several aquatic angiosperms^{4,5}, highlights its utility for reducing photorespiration where CO₂ availability may be limited. While CO₂ limitation in terrestrial plants is caused by increased stomatal resistance, in aquatics it is largely the result of the relatively high diffusional resistance of water combined with significant diel fluctuation of dissolved CO₂ in the oligotrophic lakes and seasonal pools^{4,6}.

Though it has been nearly four decades since Keeley first described “CAM-like diurnal acid metabolism” in *Isoetes howellii*⁷, relatively little is known about the genetic mechanisms controlling CAM in *Isoetes* or any other aquatic plant. Previous genomic and/or transcriptomic studies that focused on terrestrial CAM have found evidence for regulatory neofunctionalization, enrichment of *cis*-regulatory elements, and/or reprogramming of gene regulatory networks that underlie the convergent evolution of CAM in *Sedum album*⁸, *Ananas comosus*⁹, *Kalanchoe fedtschenkoi*¹⁰, several orchids^{11–13}, and Agavoideae species^{14,15}. Furthermore, a case of amino acid sequence convergence in phosphoenolpyruvate carboxylase (PEPC), which catalyzes the carboxylation of phosphoenolpyruvate (PEP) to yield oxaloacetate (OAA), has also been reported among some terrestrial CAM plants. However, the lack of a high-quality genome assembly has made a meaningful comparison of *Isoetes* or any other aquatic CAM plant to terrestrial CAM species impossible.

The only lycophyte genomes available to date are from the genus *Selaginella*^{16–18}, leaving a deep, >300-million-year gap in our knowledge of lycophyte genomics and limiting inferences of tracheophyte evolution. *Selaginella* is the only genus in the Selaginellales, the sister clade to Isoetales. Notably, *Selaginella* is known for being one of few lineages of vascular plants for which no ancient whole-genome duplications (WGDs) have been detected. Conversely, there is evidence from transcriptomic data for as many as two rounds of WGD in *Isoetes tegetiformans*¹⁹. As such, a thorough characterization of the history of WGD in *Isoetes* is vital to future research into the effects and significance of WGD across lycophyte diversity.

With this study, we seek to investigate genome evolution as well as the genetic underpinnings of CAM in *Isoetes*. To that end, we present a high-quality genome assembly for *Isoetes taiwanensis* DeVol. We find evidence for a single ancient WGD event that appears to be shared among multiple species of *Isoetes*. Additionally, while many CAM pathway genes display similar expression patterns in *Isoetes* and terrestrial angiosperms, notable differences in gene expression suggest that the evolution of CAM in *Isoetes* may have followed a markedly different path than it has in terrestrial angiosperms.

Results

Genome assembly, annotation, and organization. Using Illumina short-reads, Nanopore long-reads, and Bionano optical mapping, 90.13% of the diploid ($2n = 2X = 22$ chromosomes) *I. taiwanensis* genome was assembled into 204 scaffolds (N50 = 17.40 Mb), with the remaining 9.87% into 909 unplaced contigs (Table 1). The total assembled genome size (1.66 Gb) is congruent with what was estimated by k-mers (1.65 Gb) and flow cytometry (1.55 Gb) (Supplementary Fig. 1). A circular-mapping plastome was also assembled, from which we identified a high level of RNA-editing (Supplementary Note 1 and Supplementary Fig. 2).

A total of 39,461 high confidence genes were annotated based on ab initio prediction, protein homology, and transcript evidence. The genome and proteome BUSCO scores are 94.5% and 91.0%, respectively, which are comparable to many other seed-free plant genomes (Supplementary Fig. 3) and indicative of high completeness. Orthofinder²⁰ analysis of 25 genomes placed 647,535 genes into 40,144 orthogroups (Supplementary Note 2 and Supplementary Fig. 4). Subsequent examination of lignin biosynthesis genes in *I. taiwanensis* suggests that evolution of particular pathway steps to S-lignin likely predates the divergence of *Isoetes* and *Selaginella* (Supplementary Note 3 and Supplementary Figs. 5–17). In addition, analysis of key stomatal and root genes (Supplementary Notes 4 and 6) in *I. taiwanensis* genome supported their homology (at the molecular level) with similar structures in other vascular plants (Supplementary Table 1 and Supplementary Figs. 18–20).

Repetitive sequences accounted for 38% of the genome assembly with transposable elements (TEs) accounting for the majority of those at 37.08% of the assembly length. Long terminal repeat (LTR) retrotransposons were the most abundant (15.72% of total genome assembly) with the Gypsy superfamily accounting for around 68% of LTR coverage (10.7% of total genome assembly; Supplementary Data 1). When repeat density was plotted alongside gene density, the distribution of both was found to be homogeneous throughout the assembly (Fig. 1). This even distribution of genes and repeats is markedly different from what has been reported in most angiosperm genomes²¹ where gene density increases near the ends of individual chromosomes. However, it is consistent with several high-quality genomes published from seed-free plants, including *Physcomitrium patens*²², *Marchantia polymorpha*²³, and *Anthoceros agrestis*²⁴. The result from *I. taiwanensis* thus adds to the growing evidence that the genomic organization might be quite different between seed and seed-free plants²⁵.

Table 1 *Isoetes taiwanensis* genome assembly statistics.

Assembly size (Mb)	1658.30
Scaffolds (#)	204
Scaffold length (Mb)	1494.58
N50 of scaffold length (Mb)	17.40
Scaffolded contigs (#)	1879
Scaffolded contig length (Mb)	1211.25
N50 length of scaffolded contigs (Mb)	1.48
Un scaffolded contig (#)	909
Un scaffolded contig (Mb)	149.46
N50 length of un scaffolded contigs (Mb)	0.26
Genome BUSCO score (Eukaryota) (%)	94.5
Proteome BUSCO score (Eukaryota) (%)	91.0
Predicted protein-coding genes (#)	39,461
Predicted repetitive sequence (%)	38

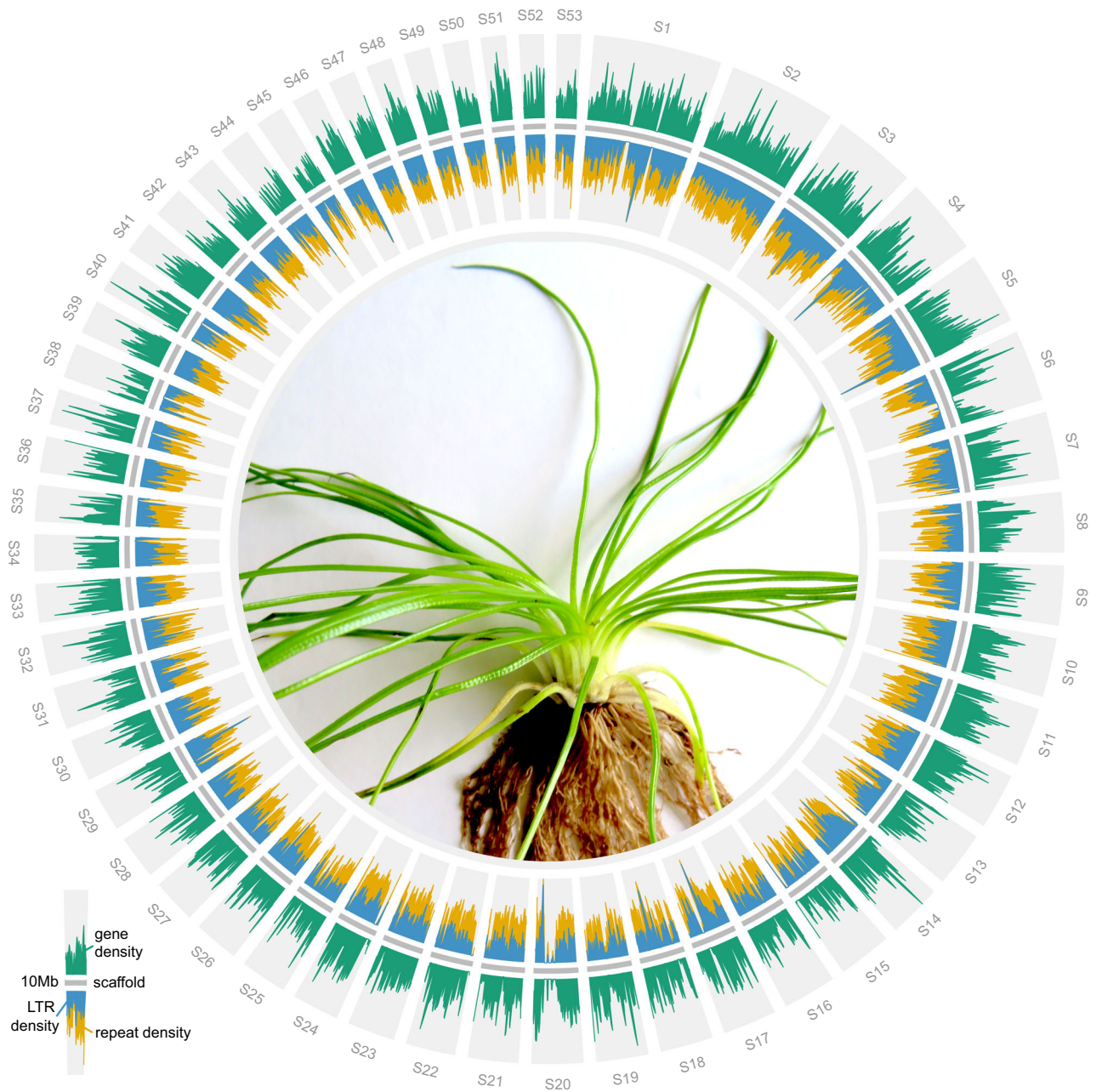


Fig. 1 Distribution of genes and repetitive elements in *I. taiwanensis*. The relatively even distributions differ from angiosperm genomes, but are similar to what have been reported in other seed-free plants. Only scaffolds longer than 10 Mb are plotted. Center: an image of *I. taiwanensis*. Source data are provided as a Source Data file.

Evidence for WGD in *Isoetes taiwanensis*. Using a combination of methods including synonymous substitutions per site (K_s), phylogenetic reconciliation, and synteny analyses, we identified a single ancient WGD in *I. taiwanensis*. This is in contrast to a previous K_s analysis using 1KP transcriptome data, which found evidence for two rounds of WGD, named ISTE α and ISTE β , in the North American species *I. tegetiformans* and *I. echinospora*²⁶. These two WGDs have median K_s values of ~ 0.5 and ~ 1.5 ²⁶ (Supplementary Fig. 21). Our K_s analysis of the whole paranome (i.e., all of the paralogous gene copies in the genome) in *I. taiwanensis* revealed a single peak at $K_s \sim 1.8$ (Fig. 2a), suggesting that the earlier of the two duplications (ISTE β) in *I. tegetiformans* and *I. echinospora* is shared by *I. taiwanensis* while the more recent event (ISTE α) is not. This result was corroborated using

four-fold degenerate site transversion rates (4d_{tv}; Supplementary Fig. 21). Further analysis of orthologous divergence between *I. taiwanensis* and *I. lacustris* indicated that ISTE β predates the divergence of these two species (Supplementary Fig. 22). The ISTE β event was subsequently confirmed by gene tree-species tree reconciliation using genomic data in the WhALE package²⁷. WhALE returned a posterior distribution of gene retention centered on $q = \sim 0.12$. This result compares favorably with a previously documented WGD event in *Azolla filiculoides*²⁸ ($q = \sim 0.08$) and is in stark contrast to our negative control, *Marchantia polymorpha*²³ ($q = \sim 0$) (Fig. 2b, c).

While self-self syntenic analysis revealed 6196 genes (15.7%) with a syntenic depth of $1\times$ in 107 clusters (Supplementary Fig. 23), we do not believe they resulted from WGD. Our K_s

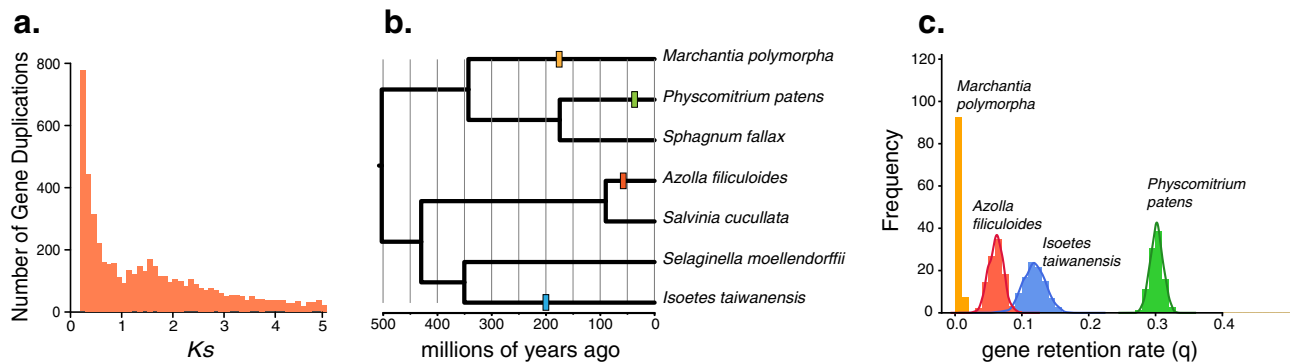


Fig. 2 Evidence for WGD in *I. taiwanensis*. **a** K_s plot showing a peak centered on 1.8 corresponding to the ISTE β event. **b** Hypothesized WGD events that were tested (colored rectangles) in our WhALE analysis are shown on a phylogeny. **c** *I. taiwanensis*' posterior distribution of gene retention rates (q) falls between that of *A. filiculoides* and *P. patens*, both are known to have at least one WGD. This provides additional support for the ISTE β event. Conversely, the gene retention rate is close to zero for *M. polymorpha*, consistent with its lack of WGD. Source data are provided as a Source Data file.

analysis restricted to syntenic gene pairs failed to recover the peak at $K_s \sim 1.8$ and instead consisted of an initial slope toward a much lower K_s value (Supplementary Fig. 24). Given their high degree of similarity and location on separate scaffolds, it is possible that these low K_s gene pairs are the result of relatively recent segmental duplications. The absence of conserved synteny from ISTE β is unsurprising. The high K_s value implies that ISTE β is ancient; long enough ago for extensive genomic restructuring and fractionation to have taken place. Altogether, of the two hypothesized WGDs in *Isoetes*, we confirmed the presence of ISTE β while the younger ISTE α might be either specific to *I. tegetiformans* and *I. echinospora* or an artifact stemming from the quality or completeness of the transcriptomes.

Similarities to terrestrial CAM plants. The CAM in *Isoetes* is unusual for at least two reasons. First, *Isoetes* diverged from other CAM plants more than 300 million years ago and second, *Isoetes* has an aquatic lifestyle. Here, we demonstrated that when submerged, titratable acidity in the leaves of *I. taiwanensis* increased throughout the night, reaching peak acidity in the morning and decreased throughout the daylight hours (Fig. 3b), consistent with the cycle of carbon sequestration and assimilation seen in dry-adapted CAM plants. To identify the underlying genetic elements, we generated TOD RNA-seq (Supplementary Data 2), sampling every 3 h over a 27 h period under 12 h light/12 h dark and continuous temperature (LDHH). A multidimensional scaling (MDS) plot of normalized expression data showed that the samples were generally clustered in a clockwise fashion as expected for TOD expression analysis (Supplementary Fig. 25).

We found that some of the CAM pathway genes in *I. taiwanensis* exhibited TOD expression patterns that largely resemble those found in terrestrial CAM plants (Fig. 3c–i and Supplementary Data 3). For example, the strong dark expression of PHOSPHOENOLPYRUVATE CARBOXYLASE KINASE (PEPCK) appears to be conserved in *I. taiwanensis* as well as in all three terrestrial taxa (Fig. 3i). Likewise, we found one copy of β -CARBONIC ANHYDRASE (β -CA) that cycled similarly with homologs in *A. comosus* and *K. fedtschenkoi* (Fig. 3g)—increasing during the night and peaking in the early morning—although this is different from *S. album* in which no β -CA genes showed a high dark expression. Similar to *A. comosus* where two copies of MALATE DEHYDROGENASE (MDH) were found to cycle in green leaf tissue⁹, we found multiple copies of MDH that appear to cycle in *I. taiwanensis* with one copy appearing to exhibit a similar peak expression to its orthologue in pineapple (Fig. 3e).

However, neither of the other two MDH genes that cycle in *I. taiwanensis* exhibit similar expression to their orthologues in terrestrial CAM species (Supplementary Fig. 26).

During the day, decarboxylation typically occurs by one of two separate pathways (Fig. 3a). The first utilizes NADP-MALIC ENZYME (NADP-ME) and PYRUVATE PHOSPHATE DIKINASE (PPDK), and appears to be favored by *K. fedtschenkoi* and *S. album*^{8,10}. The second utilizes MDH and PHOSPHOENOLPYRUVATE CARBOXYKINASE (PEPCK) and is favored by *A. comosus*⁹. Based on its TOD expression of multiple copies of MDH and associated expression dynamics, it is possible that *I. taiwanensis* utilizes the MDH/PEPCK pathway. While all four genes have elevated expression levels during the day, the expression of NADP-ME is inverted compared to *K. fedtschenkoi* and *S. album* (Fig. 3c), and PPDK exhibits relatively weak cycling overall ($R = 0.637$; Fig. 3d). Additionally, PEPCK and one copy of MDH have similar TOD expression in *I. taiwanensis* and *A. comosus* (Fig. 3f, e, respectively), which may indicate a shared affinity for MDH/PEPCK decarboxylation. Interestingly, the copy of PEPCK that cycles in *I. taiwanensis* is not orthologous to the copy that cycles in *A. comosus*, being placed in a different orthogroup by Orthofinder²⁰.

***I. taiwanensis* likely recruited bacterial-type PEPCK.** While TOD expression of many key CAM pathway genes was broadly similar to that seen in terrestrial CAM plants, one important difference can be found in the PEPCK enzyme, which is the entry point of carboxylation in CAM and C4 photosynthesis (Fig. 3a). PEPCK is present in all photosynthetic organisms as well as many non-photosynthetic bacteria and archaea. It is a vital component of plant metabolism, carboxylating PEP in the presence of HCO_3^- to yield OAA. In plants, the PEPCK gene family consists of two clades, the “plant-type” and the “bacterial-type”. The latter was named because of its higher sequence similarity with proteobacteria PEPCK than other plant-type PEPCK genes²⁹. All CAM and C4 plants characterized to date recruited only the plant-type PEPCK³⁰, with the bacterial-type often being expressed at relatively low levels and/or primarily in non-photosynthetic tissues³¹.

Interestingly, in *I. taiwanensis* we found that both types of PEPCK were cycling and that the bacterial-type was expressed at much higher levels than plant-type PEPCK (Fig. 3h). Copies from both types had similar expression profiles in *I. taiwanensis*, peaking at dusk and gradually tapering off during the night. While this may seem counterintuitive as PEPCK is an important component of the dark reactions, it is consistent with what has

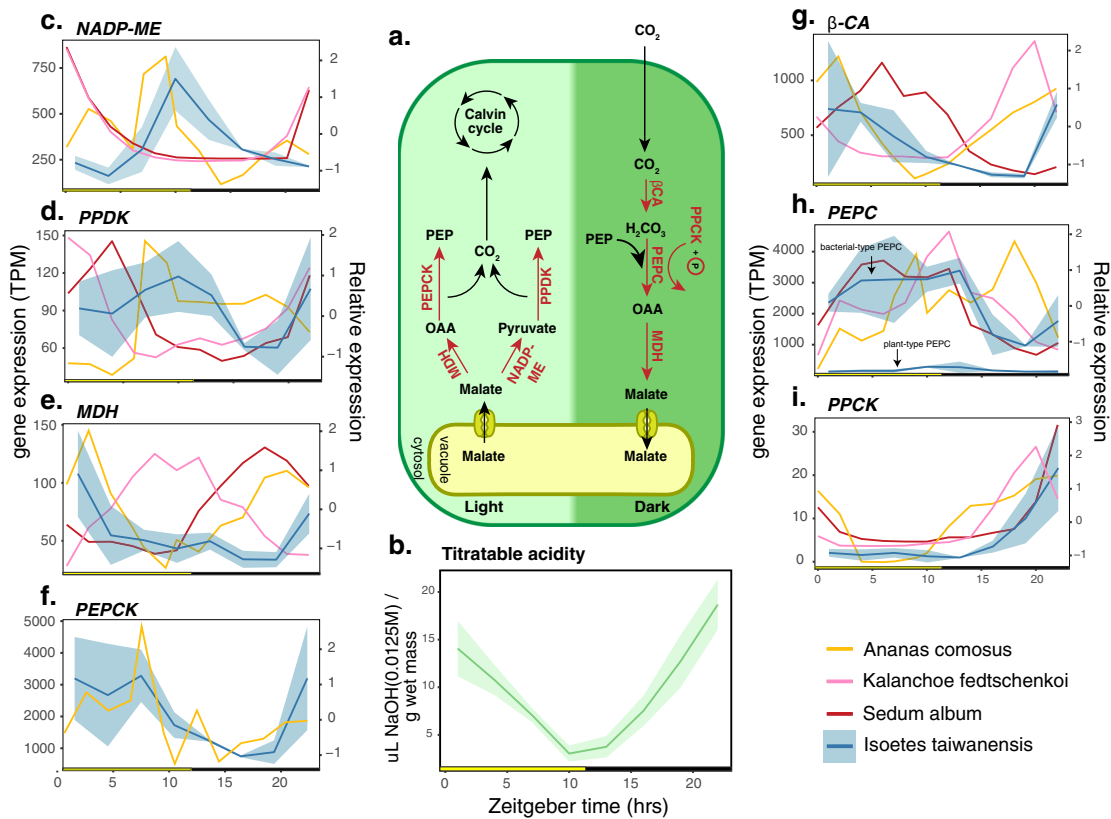


Fig. 3 Key CAM pathway genes and their expression patterns in *I. taiwanensis*. **a** The CAM pathway with important reactions and their enzymes shown in red. **b** Titratable acidity in *I. taiwanensis* exhibited a clear diel fluctuation. Diel expression patterns for highlighted genes are shown for the day (**c–f**) and night reactions (**g–i**). Average of TPM normalized expression data for *I. taiwanensis* is plotted in blue with a shaded ribbon representing the standard deviation. Relative expression profiles for homologous, cycling genes in other CAM species are plotted for comparison. All times are displayed in hours after lights-on (Zeitgeber time). Locus IDs for genes used in expression plots is provided in Supplementary Data 3. Source data for *I. taiwanensis* are provided in Supplementary Data 2. Source data for the other taxa are provided as a Source Data file.

previously been found in other terrestrial CAM plants, with the overall expression profile resembling that of *S. album*⁸. The advantage of recruiting bacterial-type PEPC is unclear. In vivo, both bacterial- and plant-type PEPC can interact with each other to form a hetero-octameric complex that is less sensitive to inhibition by malate³². Although the functional and physiological implications await future studies, the unusual involvement of bacterial-type PEPC is suggestive of a divergent evolutionary path to underwater CAM in *Isoetes*.

No evidence for convergent evolution of PEPC. Plant-type PEPC was recently shown to undergo convergent amino acid substitutions in concert with the evolution of CAM¹⁰. An aspartic acid (D) residue appears to have been repeatedly selected across multiple origins of CAM such as in *K. fedtschenkoi* and *P. equestris*¹⁰, although notably not in *A. comosus* nor *S. album*. This residue is situated near the active site, and based on in vitro assays, the substitution to aspartic acid significantly increased PEPC activity¹⁰. However, in *I. taiwanensis* we did not observe the same substitution in any copies of PEPC (Fig. 4); instead, they have arginine (R) or lysine (H) at this position like PEPC from many non-CAM plants. This lack of sequence convergence between *Isoetes* and few CAM angiosperms could be the result of their substantial phylogenetic distance and highly divergent life histories. Alternatively, it is also likely that the substitution is not as important as previously hypothesized, or relevant only in the context of plant-type PEPC. As *I. taiwanensis* may also utilize

bacterial-type PEPC, the aspartic acid residue might not serve the same purpose.

Circadian regulation in *Isoetes*. Previous analysis of the *A. comosus* genome found promoter regions of multiple key CAM pathway genes containing known circadian *cis*-regulatory elements (CREs) including Morning Element (ME: CCACAC), Evening Element (EE: AAATATCT), CCA1-binding site (CBS: AAAAATCT), G-box (CACGTG) and TCP15-binding motif (NGGNCCAC)⁹. This suggests that expression of CAM genes in pineapple is largely under the control of a handful of known circadian clock elements. The direct involvement of circadian CREs was corroborated by a later study of the facultative CAM plant *S. album* where shifts in diel expression patterns were tied to a shift in TOD-specific enrichment of CREs: EE and Telobox (TBX: AAACCCT)⁸.

In order to examine the role of the circadian clock and light/dark cycles in regulating *I. taiwanensis* CAM, we used the HAYSTACK pipeline³³ to identify all genes with TOD expression patterns. We predicted 3241 cycling genes, which is 10% of the expressed genes (Supplementary Fig. 27 and Supplementary Note 6). While 10% is low compared to land plants that have been tested under this condition (LDHH)—usually at 30–50% genes^{8,33,34}, a recent study found a reduced number of cycling genes in another aquatic plant *Wolffia australiana* (duckweed/watermeal)³⁵. Accordingly, decreased cycling may be a feature of aquatic plants. Further discussions and comparisons of *I. taiwanensis* TOD gene expression with other species can be

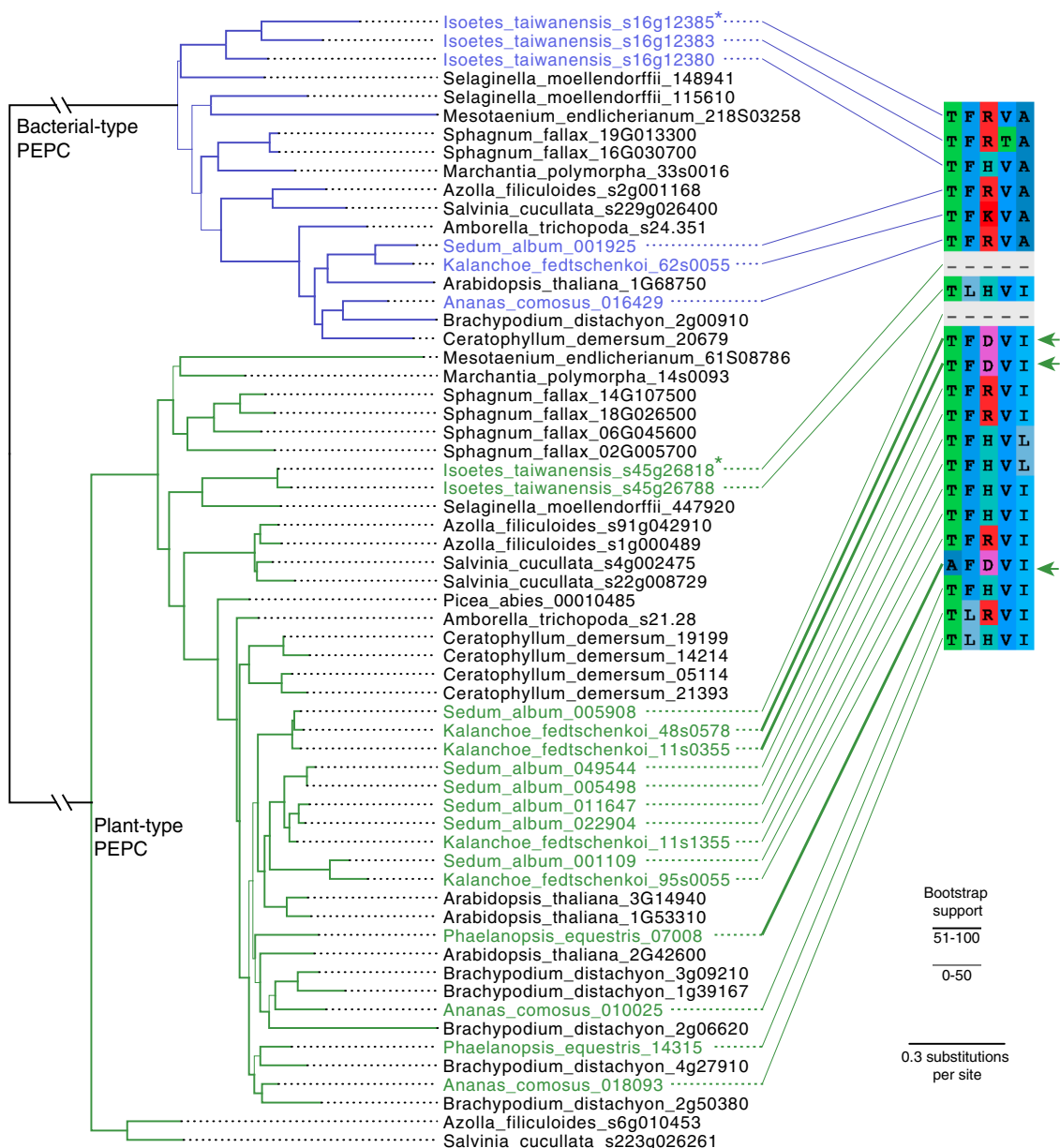


Fig. 4 A lack of PEPC sequence convergence in *I. taiwanensis*. Copies with putative convergent amino acid sequence (D at position 3 in alignment) are indicated by thickened connecting lines and green arrows. Copies of bacterial-type and plant-type PEPC shown to cycle in *I. taiwanensis* are marked with asterisks (*). Branch thickness indicates bootstrap support.

found in Supplementary Note 7, Supplementary Figs. 27, 28, and Supplementary Data 4.

Core circadian clock genes such as *LATE ELONGATED HYPOCOTYL (LHY)*, *PSEUDO-RESPONSE REGULATOR 7 (PRR7)*, *LUX ARRHYTHMO (LUX)*, and *EARLY FLOWERING 3 (ELF3)*, cycle with the expected TOD expression seen in their *Arabidopsis* orthologs (Fig. 5, Supplementary Note 8, and Supplementary Figs. 29, 30)³³. Furthermore, *TIMING OF CAB2 1/PSEUDO-RESPONSE REGULATOR 1 (TOC1/PRR1)* and *GIGANTEA (GI)*, which are typically single-copy genes in land plants, have, respectively, 3 and 5 predicted genes in distinct genomic locations. Similarly, an increased number of homologs was found in the facultative CAM plant *S. album*⁸. Closer inspection confirmed all 3 *TOC1/PRR1* paralogs are full length, while only 1 of the *GI* genes (*Gla*) is full length and 1 other (*Glb*) is a true partial/truncated (and expressed) paralog. Surprisingly, all 3 copies of

TOC1/PRR1 have dawn-specific expression compared to the dusk-specific expression found in all plants tested to date³⁶ (Fig. 5b) including terrestrial CAM species (Supplementary Fig. 30). In addition, *Gla* and *Glb* have antiphase expression, with the full-length *Gla* having dusk-specific expression, which is consistent with other plants, and *Glb* having dawn-specific expression (Fig. 5c and Supplementary Fig. 30).

The duplications and divergent expression patterns of *TOC1/PRR1* and *GI* in *I. taiwanensis* have important implications on circadian clock evolution. Despite the TOD expression of core circadian clock genes being highly conserved since the common ancestor of green algae and angiosperms, the mechanisms may be simpler in algae³⁷ and mosses³⁸. This idea is largely based on the lack of key components of the evening-phased loop including *PRR1*, *GI*, and *ZTL* in *P. patens* and the absence of the same along with morning-phased loop genes *ELF3* and *ELF4* in algae³⁹.

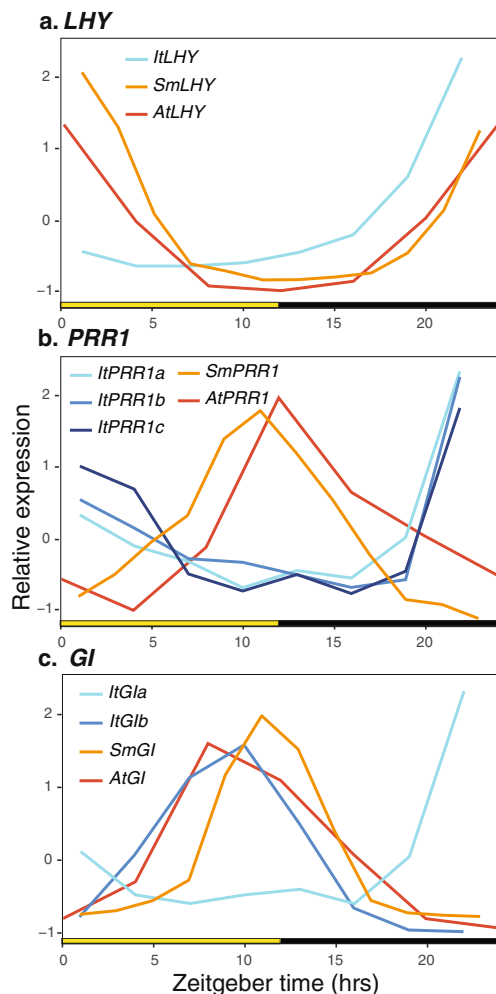


Fig. 5 Expression of key circadian associated genes is shifted in *I. taiwanensis*. **a** LATE ELONGATED HYPOCOTYL (*LHY*), **b** PSEUDO-RESPONSE REGULATOR 1 (*PRR1*), **c** GIGANTEA (*GI*) orthologs in *Isoetes* (blue lines), *Selaginella* (orange line), and *Arabidopsis* (red line) normalized expression over the day. Day (yellow box); night (black box); Zeitgeber time (ZT) is the number of hours (h) after lights on (0 h). Locus IDs for genes used in expression plots can be found in Supplementary Data 3. Source data for *I. taiwanensis* are provided in Supplementary Data 2. Source data for the other taxa are provided as a Source Data file.

While *I. taiwanensis* possesses all the major clock genes that are found in other vascular plants, lineage-specific expansion and phase-shifted gene expression in the evening-phased loop could indicate that circadian control was less conserved during the early evolution of land plants. However, *Selaginella* exhibits a very similar expression of various circadian modules relative to other vascular plants and likewise, possesses a single copy of both *GI* and *PRR1*³⁹. It is thus possible that the TOD architecture in *I. taiwanensis* represents a more recent adaptation to its aquatic CAM lifestyle. As a comparison, *S. album* similarly has multiple duplicated clock genes and its transition to CAM is associated with significant shifts in both phase and amplitude of gene expression⁸. To further investigate the relationship between clock and CAM in *I. taiwanensis*, we next focused on characterizing the circadian CREs.

Canonical circadian CREs are not enriched in *Isoetes* CAM cycling genes. We used ELEMENT³³ to exhaustively search the

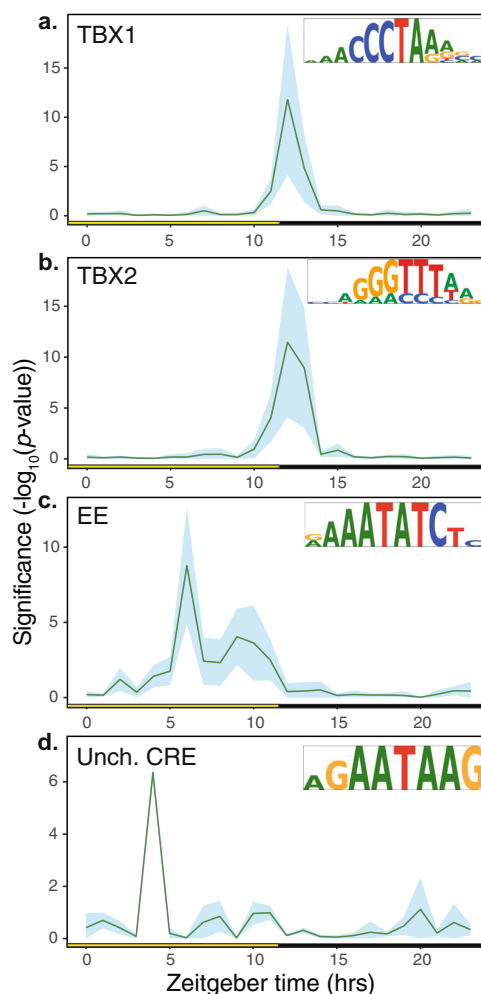


Fig. 6 Multiple CREs exhibit time-structured enrichment in *I. taiwanensis*. **a, b** Two telobox (TBX) containing motifs showed similar patterns to one another, both being enriched ingenes with peak expression at dusk. **c** Motif containing Evening Element (EE) was significantly enriched in genes with peak expression at mid-day. **d** A significantly enriched motif at mid-day. Day (yellow box); night (black box); Zeitgeber time (ZT) is the number of hours (h) after lights on (0 h). Shaded regions represent the standard deviation of log-transformed, FDR corrected, single-tailed *p*-values of component *k*-mers as calculated by ELEMENT³³. Source data are provided as a Source Data file.

promoter region of cycling genes for putative CRE motifs. Following de novo identification, putative CREs were compared to known transcription factor binding sites in *Arabidopsis* to determine to what degree their functions might be conserved between *Isoetes* and flowering plants. We identified 16 significantly enriched CREs motifs in the 500 bp 5' promoter region of cycling genes identified by HAYSTACK, and clustered them according to TOD expression (Supplementary Data 5). Half of the motifs shared some degree of sequence similarity to known circadian CREs previously identified in *Arabidopsis*, including the EE as well as two “ACGT”-containing elements (G-box-like) and two TBX-containing motifs³³. In the case of TBX, both motifs were associated with peak expression at dusk (at around 12 h after lights on; Zeitgeber Time [ZT]) in *I. taiwanensis* (Fig. 6a, b), similar to *Arabidopsis* under light/dark cycles alone³³. On the other hand, the EE appear to be associated with peak expression at different TOD. In *Arabidopsis*, the EE is enriched in genes with peak expression at dusk (ZT = 12), but in *I. taiwanensis*, this

pattern is shifted, with the EE associating with genes that peak in expression around mid-day (ZT = 6) (Fig. 6c). Additionally, while the two “ACGT”-containing elements were found upstream of genes that exhibited significant cycling behavior, neither was strongly associated with peak expression at a particular TOD. We also found an unidentified CRE (AGAATAAG) that is strongly associated with peak expression in the morning (ZT = 4) (Fig. 6d).

We next examined the connection between circadian CREs and CAM genes in *I. taiwanensis*. Interestingly, with the exception of the RVE1/2 motif, we did not find significant enrichment of any known circadian CREs in CAM cycling genes relative to non-cycling paralogues. While a targeted search of CAM cycling gene promoters did uncover circadian CREs including the CBS, TCP15, TBX, and EE (Supplementary Data 6), none were strongly associated with either light or dark phase CAM gene expression. In addition, both ME and G-box were conspicuously absent from the promoter regions of cycling CAM photosynthetic genes.

In sum, TOD-specific enrichment of CREs appears to differ in various aspects from *Arabidopsis*. While some CRE sequences themselves are conserved between lycophytes and angiosperms, their interaction with various transcription factors and subsequent regulatory function could be quite different in *Isoetes*. Importantly, our results stand in contrast to other CAM plants such as *S. album*⁸ and *A. comosus*⁹ where CAM genes appeared to be under the direct control of a handful of strictly conserved circadian CREs. These results either suggest that the circadian clock network that emerged in *Isoetes*, which included the addition of central components *GI* and *PRR1*, was quite different than that found to be highly conserved in seed plants, or there is significant TOD innovation associated with the evolution of underwater CAM. Additional *Isoetes* genomes and TOD analysis of underwater CAM plants will be required to test these hypotheses.

The assembly and analyses of the *I. taiwanensis* genome bridges a substantial gap in our knowledge of vascular plant evolution. We have combined genomic and transcriptomic data to corroborate one of the two hypothesized WGDs in *Isoetes* relative to its closest extant relative *Selaginella*, highlighting the contrasting history of WGD in these two lineages. Additionally, comparison of TOD gene expression with genomic sequence data has given us insights into the convergent evolution of CAM photosynthesis, not only in a lycophyte, but also in the aquatic environment. As such, our analysis stands as a necessary counterpoint to similar studies previously conducted in terrestrial angiosperms. Shifts in expression of CAM pathway genes and the possible recruitment of bacterial-type PEPC in *I. taiwanensis* demonstrate a remarkable degree of plasticity in the convergent evolution of this complex trait throughout vascular plants. Likewise, differences in the enrichment of CREs associated with circadian gene expression suggest that control of CAM, as well as other processes tied to the circadian clock, may have diverged markedly since the common ancestor of *Isoetes* and flowering plants. We propose that the emergence of underwater CAM may have followed a distinct route in *Isoetes*, shedding light on a classic example of convergent evolution of a complex plant trait.

Methods

Plant sample. *Isoetes taiwanensis* is endemic to a small pond in Northern Taiwan and has been ex situ propagated in Taiwan Forestry Research Institute. This species is expected to have a low genetic diversity due to a very restricted distribution and a small population size. The voucher specimen (Kuo4500) was deposited at TAIF herbarium.

Genome size estimate. The genome size of *I. taiwanensis* was first determined by flow cytometry following the protocols outlined in Kuo et al.⁴⁰ and Li et al.²⁸. The

flow cytometric experiments were performed on BD FACScan system (BD Biosciences, USA), and the Beckman buffer⁴¹ was used with 0.5% (v/v) 2-mercaptoethanol, 40 mg mL⁻¹ PVP-40, and 0.1 mg mL⁻¹ RNaseA added. We used *Zea mays* (1C = 5.57pg⁴²) as the internal standard. To confirm the flow cytometry-based measurement, a k-mer frequency distribution was generated from Illumina 2 × 150 bp paired reads (described below) using Jellyfish⁴³, which was then input into GenomeScope⁴⁴ to estimate genome size.

Genome sequencing. High molecular weight (HMW) DNA was extracted using a modified CTAB method on isolated nuclei. First, leaf tissues were ground in liquid nitrogen, and the powder was resuspended in the Beckman buffer (same as in our flow cytometric experiments). We then used 30 µm nylon circular filters (Partec, Germany) to remove tissue debris, and precipitated nuclei with 100×g centrifugation under 4 °C for 20 min. DNA was extracted following a modified CTAB protocol⁴⁵. HMW DNA was QC'd on an agarose gel for length and quantified on a bioanalyzer. Unsheared HMW DNA was used to make Oxford Nanopore Technologies (ONT) ligation-based libraries (Oxford, UK). Libraries were prepared starting with 1.5 µg of DNA and following all other steps in ONT's SQK-LSK109 protocol. Final libraries were loaded on an ONT flowcell (v9.4.1) and run on the GridION. Bases were called in real-time on the GridION using the flip-flop version of Guppy (v3.1). The resulting fastq files were concatenated and used for downstream genome assembly steps. The same batch of HMW genomic DNA was used to construct Illumina (Illumina, USA) libraries for estimating genome size (above) and correcting residual errors in the ONT assembly. Libraries were constructed using the KAPA HyperPrep Kit (Kapa Biosystems, Switzerland) followed by sequencing on an Illumina NovaSeq6000 with 2 × 150 bp paired-ends.

Genome assembly. ONT reads were assembled using minimap2 and miniasm⁴⁶, and the resulting draft assembly was then polished by racon⁴⁷ (with ONT reads) and pilon⁴⁸ (with Illumina reads). Because the plants were grown non-axenically under water, the assembly inevitably contained contaminations. We, therefore, used blobtools⁴⁹ to identify non-plant contigs based on a combination of contig read coverage, taxonomic assignment, and GC content.

To further scaffold the assembly, we generated a genome map using Bionano with the Direct Label and Stain chemistry and DLE-1 labeling. For this, high molecular weight DNA was extracted using the Bionano Plant DNA Isolation Kit. Hybrid scaffolding, combining the nanopore draft and Bionano map, was done on the Bionano Saphyr computing platform at the McDonnell Genome Institute at Washington University. We then gap-filled the scaffolded genome using two rounds of LR_Gapcloser⁵⁰ (3 iterations each and a pilon polishing in between). Finally, to remove redundancy the purge_haplotigs pipeline⁵¹ was used to obtain the v1 assembly. The circular chloroplast genome was assembled from Illumina data using the GetOrganelle⁵² toolkit.

Repeat annotation. We generated a custom *I. taiwanensis*-specific repeat library using LTR-retriever⁵³ and RepeatModeler⁵⁴. The *I. taiwanensis* genome contains a high number of genes encoding pentatricopeptide repeat proteins which are often misclassified as repetitive elements in the genome. Thus, in order to identify and remove repeats with homology to plant proteins, we used BLASTx to query each repeat against the uniprot plant protein database (*e*-value threshold at 1e−10). The resulting library was then input into RepeatMasker⁵⁵ to annotate and mask the repetitive elements in the *I. taiwanensis* genome.

Gene annotation. We trained two ab initio gene predictors, AUGUSTUS⁵⁶ and SNAP⁵⁷, on the repeat-masked genome using a combination of protein and transcript evidence. For the protein evidence, we relied on the annotated proteomes from *Selaginella moellendorffii*¹⁶ and *S. lepidophylla*¹⁸, and for the transcript evidence, we used the RNA-seq data from our time-course experiment and a separate corm sample. To train AUGUSTUS, BRAKER2⁵⁸ was used and the transcript evidence was input as an aligned bam file. SNAP was trained under MAKER with 3 iterations, and in this case, the transcript evidence was supplied as a de novo assembled transcriptome done by Trinity⁵⁹. After AUGUSTUS and SNAP were trained, they were fed into MAKER⁶⁰ along with all the evidence to provide a synthesized gene prediction. Gene functional annotation was done using the eggNOG-mapper v2⁶¹. To filter out spurious gene models, we removed genes that met none of the following criteria: (1) a transcript abundance greater than zero in any sample (as estimated by Stringtie⁶²), (2) has functional annotation from eggNOG, and (3) was assigned into orthogroups in an Orthofinder²⁰ run (see below). The resulting gene set was used in all subsequent analyses.

Homology assessment and gene family analysis. Homology was initially assessed with Orthofinder²⁰ using genomic data from a range of taxa from across the plant tree of life including all CAM plant genomes published to date: *Amborella trichopoda*⁶³, *Ananas comosus*⁹, *Anthoceros agrestis*²⁴, *Arabidopsis thaliana*⁶⁴, *Azolla filiculoides*²⁸, *Brachypodium distachyon*⁶⁵, *Ceratophyllum demersum*⁶⁶, *Isoetes taiwanensis* (this study), *Kalanchoe fedtschenkoi*¹⁰, *Marchantia polymorpha*²³, *Medicago truncatula*⁶⁷, *Nelumbo nucifera*⁶⁸, *Nymphaea colorata*⁶⁹, *Phalaenopsis equestris*¹¹, *Physcomitrium patens*²², *Picea abies*⁷⁰, *Salvinia cucullata*²⁸, *Sedum album*⁸, *Selaginella moellendorffii*¹⁶, *Sphagnum fallax* (*Sphagnum fallax* v0.5, DOE-

JGI, <http://phytozome.jgi.doe.gov/>), *Spirodela polyrhiza*⁷¹, *Utricularia gibba*⁷², *Vitis vinifera*⁷³, and *Zostera marina*⁷⁴, and one algal genome: *Mesotetanium endlicherianum*⁷⁵. Following homology assessment, the degree of overlap between gene families was assessed using the UpsetR⁷⁶ package in R.

RNA-editing analysis. RNA-seq data were first mapped to combined nuclear and chloroplast genome assemblies using HISAT2⁷⁷. The reads mapping to the chloroplast genome were extracted using samtools⁷⁸. SNPs were called using the mpileup function in bcftools⁷⁹. The resulting vcf files were filtered using bcftools to remove samples with a depth < 20, quality score < 20 and mapping quality bias < 0.05. After filtering, C-to-U and U-to-C edits were identified using an alternate allele frequency threshold of 10%. Finally, RNA-editing sites were related to specific genes using the intersect command in bedtools⁸⁰ and characterized using a custom python script (available at https://github.com/dawickell/Isoetes_CAM).

Ks analysis. Ks divergence was calculated by several different methods. Initially, a whole paranome Ks distribution was generated using the “wgd mcl” tool⁸¹. Self-synteny was then assessed in i-Adhore and Ks values were calculated and plotted for syntenic pairs only using the “wgd syn” tool⁸¹. To conduct Ks analysis of related species, RNA-seq data was downloaded from the SRA database for *Isoetes yunguisensis* (SRR6920723)⁸², *I. sinensis* (SRR1648119)⁸³, *I. drummondii* (SRR4762161), *I. echinospora* (SRR6853338)⁸⁴, *I. lacustris* (SRR9620527)⁸⁵, and *I. tegetiformans* (ERR2040873)¹⁹. Transcriptomes were assembled using SOAPdenovo-Trans⁸⁶ with a k-mer length of 31. Next, for each *Isoetes* genome and transcriptome, we used the DupPipe pipeline to construct gene families and estimate the age distribution of gene duplications^{87,88}. We translated DNA sequences and identified ORFs by comparing the Genewise⁸⁹ alignment to the best-hit protein from a collection of proteins from 25 plant genomes from Phytozome⁹⁰. For all DupPipe runs, we used protein-guided DNA alignments to align our nucleic acid sequences while maintaining the ORFs. We estimated Ks divergence using PAML⁹¹ with the F3X4 model for each node in the gene family phylogenies.

Four-fold transversion substitution rate analysis. For each *Isoetes* genome and transcriptome, we used the DupPipe pipeline as described above to generate gene alignments. We estimated a four-fold transversion substitution rate (4dvtv) using an existing perl script for each duplicate gene pair (https://github.com/chaimol/KK4D/blob/master/calculate_4DTV_correction.pl).

Estimation of orthologous divergence. To place putative WGDs in relation to lineage divergence, we estimated the synonymous divergence of orthologs among pairs of species that may share a WGD based on their phylogenetic position and evidence from the within-species Ks plots. We used the RBH Orthologous pipeline⁸⁸ to estimate the mean and median synonymous divergence of orthologs, and compared those with the synonymous divergence of inferred paleopolyploid peaks. We identified orthologs as reciprocal best BLAST hits in pairs of transcriptomes. Using protein-guided DNA alignments, we estimated the pairwise synonymous divergence for each pair of orthologs using PAML⁹¹ with the F3X4 model.

Phylogenetic assessment of ancient whole-genome duplication. WGD inference was conducted by phylogenomic reconciliation using the WhALE package implemented in Julia²⁷. First, prior to WhALE analysis, Orthofinder²⁰ was used to identify groups of orthologous genes among 7 species representing 3 taxonomic groups (bryophytes, lycophytes, and ferns): *Azolla filiculoides*²⁸, *Isoetes taiwanensis* (this study), *Marchantia polymorpha*²³, *Physcomitrium patens*²², *Salvinia cucullata*²⁸, *Selaginella moellendorffii*¹⁶, and *Sphagnum fallax* (*Sphagnum fallax* v0.5, DOE-JGI, <http://phytozome.jgi.doe.gov/>). These species were chosen based on phylogenetic relatedness, availability of a high-quality genome assembly, and previous assessment for the presence or absence of WGD. The resulting orthogroups were filtered using a custom python script to remove the 5% largest orthogroups and those with less than 3 taxa. Additionally, WhALE requires removal of gene families that do not contain at least one gene in both bryophytes and ferns to prevent the inclusion of gene families originating after divergence from the most recent common ancestor. Alignments were generated for the filtered orthogroups in PRANK⁹² using the default settings. A posterior distribution of trees was obtained for each gene family in MrBayes 3.2.6⁹³ using the LG model. Chains were sampled every 10 generations for 100,000 generations with a relative burn-in of 25%. Following the Bayesian analysis, conditional clade distributions (CCDs) were determined from posterior distribution samples using ALEobserve in the ALE software suite⁹⁴. CCD files were subsequently filtered using the ccddata.py and ccdfilter.py scripts provided with the WhALE program. A dated, ultrametric species tree was generated using the “ape” package in R⁹⁵, in which branch lengths were constrained according to 95% highest posterior density of ages, assuming that bryophytes are monophyletic, as reported by Morris et al.⁹⁶. Finally, the filtered CCD files were loaded in Julia along with the associated species phylogeny. A hypothetical WGD node was inferred at 200 million years ago (MYA) along the branch leading to *I. taiwanensis*, prior to the estimated crown age of extant *Isoetes*⁹⁷. Modifying the hypothetical age of this WGD node did not affect the outcome. Additional WGD nodes were placed as positive controls along branches leading to *Physcomitrium patens* and *Azolla filiculoides* at 40 MYA and 60 MYA,

respectively, based on previous studies^{22,28}. A false WGD event was also placed arbitrarily in *Marchantia polymorpha* at 160 MYA as a negative control. A WhALE “problem” was constructed using an independent rate prior and MCMC analysis was conducted using the DynamicHMC library in Julia (<https://github.com/tpapp/DynamicHMC.jl>) with a sample size of 1000.

Phylogenetic analysis of root, stomata, and CAM pathway genes. Following clustering of homologs in Orthofinder, we conducted a phylogenetic analysis of several gene families of interest, including those containing *SMF*, *FAMA*, *TMM*, *RSL*, and *PEPC* genes, which were identified based on homology using gene annotations from *Arabidopsis*. Gene trees from Orthofinder were initially used to identify paralogues and remove fragmented genes where appropriate. In the case of *PEPC*, orthogroups containing “bacterial-type” and “plant-type” *PEPC* were combined prior to alignment. Next, amino acid sequences were aligned using MUSCLE⁹⁸ under default settings and trimmed using TrimAL with the -strict flag. An amino acid substitution model was selected according to the Bayesian Information Criterion (BIC) in ModelFinder⁹⁹ prior to phylogenetic reconstruction by maximum likelihood in IQ-TREE v1.6.12¹⁰⁰ with 1000 ultrafast¹⁰¹ bootstrap replicates.

Phylogenetic and gene expression analysis of genes salient to the phenylpropanoid and lignin biosynthesis pathway. The datasets used for phylogenetic analyses were based on de Vries et al.¹⁰² with added *I. taiwanensis* sequences. In brief, we assembled a dataset of predicted proteins from (A) the genomes of seventeen land plants: *Anthoceros agrestis* as well as *Anthoceros punctatus*²⁴, *Amborella trichopoda*⁶³, *Arabidopsis thaliana*⁶⁴, *Azolla filiculoides*²⁸, *Brachypodium distachyon*⁶⁵, *Capsella grandiflora*¹⁰³, *Gnetum montanum*¹⁰⁴, *Isoetes taiwanensis* (this study), *Marchantia polymorpha*²³, *Nicotiana tabacum*¹⁰⁵, *Oryza sativa*¹⁰⁶, *Physcomitrium patens*²², *Picea abies*⁷⁰, *Salvinia cucullata*²⁸, *Selaginella moellendorffii*¹⁶, and *Theobroma cacao*¹⁰⁷; (B) the genomes of seven streptophyte algae: *Chlorokybus atmophyticus*¹⁰⁸, *Chara braunii*¹⁰⁹, *Klebsormidium nitens*¹¹⁰, *Mesotetanium endlicherianum*⁷⁵, *Mesostigma viride*¹⁰⁸, *Penium margaritaceum*¹¹¹, *Spiroglaea muscicola*⁷⁵—additionally, we included sequences found in the transcriptomes of *Spirogyra pratensis*¹¹², *Coleochaete scutata* as well as *Zygnema circumcarinatum*¹¹³, and *Coleochaete orbicularis*¹¹⁴; (C) the genomes of eight chlorophytes: *Bathycoccus prasinos*¹¹⁵, *Chlamydomonas reinhardtii*¹¹⁶, *Coccomyxa subellipsoidea*¹¹⁷, *Micromonas* sp. as well as *Micromonas pusilla*¹¹⁸, *Ostreococcus lucimarinus*¹¹⁹, *Ulva mutabilis*¹²⁰, *Volvox carteri*¹²¹. For phenylalanine ammonia-lyase, additional informative sequences were added based on de Vries et al.¹²².

Building on the alignments published in de Vries et al.¹⁰², homologs of each gene family (detected in the aforementioned species via BLASTp) were (re-)aligned using MAFFT v7.475¹²³ with a L-INS-I approach; both full and partial sequences from *I. taiwanensis* were retained. We constructed maximum likelihood phylogenies using IQ-TREE 2.0.6¹²⁴; 1000 ultrafast¹⁰¹ bootstrap replicates were computed. To determine the best model for protein evolution, we used ModelFinder⁹⁹ and picked the best models based on BIC (PAL: LG + F + R7; CSE: LG + F + R8; C4H: LG + R8; C3H: LG + R10; COMT: JTT + R7; HCT: WAG + R9; F5H: LG + F + R10; CCaOAMT: WAG + R5; 4CL: LG + R9; CAD: LG + R8; CCR: LG + R6). Residue information was mapped next to the tree based on structural analyses by Hu et al.¹²⁵, Pan et al.¹²⁶, Louie et al.¹²⁷, Youn et al.¹²⁸, and Ferrer et al.¹²⁹.

Raw read counts from expression data of different tissues (leaf time series and the corm) of *I. taiwanensis* were extracted for those genes highlighted. Data were filtered to retain genes with more than 1 count per million (CPM) in at least 2 samples and normalized by applying the trimmed mean of M values procedure with the edgeR package¹³⁰. The count data were transformed to log2-CPM via the limma package¹³¹. Finally, heatmaps of gene expression levels were produced via the R packages gplots and RColorBrewer.

Time-course titratable acidity and RNA-seq experiments. Leaves of *I. taiwanensis* were taken from five individuals (as five biological replicates) every 3 h over a 27-h period on a 12-hour light/dark cycle and constant temperature. To measure changes in acidity over time, a portion of the leaf tissues was weighed, mixed with 3.5–5 mL of ddH₂O, and titrated with 0.0125 M NaOH solution until pH = 7.0. At the same time, we froze the leaf tissues in liquid nitrogen, and extracted RNA using a modified CTAB protocol¹³². RNA quality was examined on a 1% agarose gel and RNA concentration was quantified using the Qubit RNA HS assay kit (Invitrogen, USA). Based on the RNA quality and concentration, three samples per time point were picked for sequencing. 2 µg of total RNA was used to construct stranded RNA-seq libraries using the Illumina TruSeq stranded total RNA LT sample prep kit (RS-122-2401 and RS-122-2402). Multiplexed libraries were pooled and sequenced on an Illumina NovaSeq6000 with 2 × 150 bp paired-ends.

Differential expression analysis. RNA-seq reads were mapped to the combined nuclear and chloroplast genome using HISAT2⁷⁷. Stringtie⁶² was used to assemble transcripts and estimate transcript abundance. A gene count matrix was produced using the included prepDE.py script. We imported gene count data into the DESeq2 package in R¹³³ for read normalization using its median of ratios method as well as identification and removal of outlier samples using multidimensional

scaling. A single outlier sample from each of six time points (1, 4, 7, 10, 13, and 19 h) was removed from the final dataset. The resulting dataset was used to analyze temporal gene expression patterns in the R package maSigPro¹³⁴. Using maSigPro, genes with significantly differential expression profiles were identified by computing a regression fit for each gene and filtered based on the associated *p*-value ($p < 0.001$).

Identification of CAM-associated genes. CAM-associated gene identification was accomplished by a combination of functional annotation, homology assessment, and differential expression analysis. Initially, genes previously identified as being involved in the CAM pathway in terrestrial plants were identified using their functional annotations assigned by eggNOG-mapper v2⁶¹ according to sequence similarity. Next, additional putative CAM photosynthesis genes were identified from Orthofinder results if they belonged to any group containing genes identified in the previous step or a group having known CAM-associated genes from *Ananas comosus*⁹, *Kalanchoe fedtschenkoi*¹⁰, and/or *Sedum album*⁸. Finally, genes identified in the previous steps were submitted to differential expression analysis to determine whether or not they showed TOD expression. Thus, genes were considered to be “CAM associated” if they exhibited homology to known CAM photosynthetic genes in terrestrial CAM species and cycled in *I. taiwanensis* in a TOD manner.

HAYSTACK global cycling prediction. Genes with mean expression across all the time points below 1 TPM were considered “not expressed” and filtered prior to cycling prediction with HAYSTACK (https://gitlab.com/NolanHartwick/super_cycling)³³. HAYSTACK operates by correlating the observed expression levels of each gene with a variety of user specified models that represent archetypal cycling behavior. We used a model file containing sinusoid, spiking traces, and various rough linear interpolations of sinusoids with periods ranging from 20 to 28 h in one-hour increments and phases ranging from 0 to 23 h in one-hour increments. Genes that correlated with their best fit model at a threshold of $R > 0.8$ were classified as cyclers with phase and period defined by the best fit model. This threshold for calling cycling genes is based on previous validated observations^{8,33,34,135}. We also validated this threshold by looking at the cycling of known circadian clock genes (Fig. 5).

ELEMENT cis-regulatory elements analysis. Once cycling genes in *I. taiwanensis* were identified, we were able to find putative *cis*-acting elements associated with TOD expression. Promoters, defined as 500 bp upstream of genes, were extracted for each gene and processed by ELEMENT (https://gitlab.com/salk-tm/snake_pip_element)^{33,136,137}. Briefly, ELEMENT generates an exhaustive background model of all 3–7 k-mer using all of the promoters in the genome, and then compares the k-mers (3–7 bp) from the promoters for a specified gene list. Promoters for cycling genes were split according to their TOD expression into “phase” gene lists and k-mers that were overrepresented in any of these 24 promoter sets were identified by ELEMENT. By splitting up cycling genes according to their associated phase, we gained the power to identify k-mers associated with TOD-specific cycling behavior at every hour over the day. Our threshold for identifying a k-mer as being associated with cycling was an FDR < 0.05 in at least one of the comparisons. The significant k-mers were clustered according to sequence similarity (Fig. 6).

Promoter motif identification. Core CAM genes with significantly differential diel expression profiles (as identified in maSigPro) including β -CA, PEPC, PEPCK, ME, MDH, and PPDK were selected for motif enrichment analysis. Enriched motifs were identified relative to a background consisting of non-cycling paralogues of photosynthetic genes using the AME utility¹³⁸. Promoters were searched for known motifs from the *Arabidopsis* promoter binding motif database¹³⁹ with FIMO¹⁴⁰.

Reporting summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

A reporting summary for this Article is available as a Supplementary Information file. Additional data supporting the findings of this work are available within the paper and its Supplementary Information. All the raw sequences generated by this study have been deposited in the NCBI Sequence Read Archive under the BioProject PRJNA735564. Genome assembly and annotation are available at <https://genomevolution.org/coge/GenomeInfo.pl?gid=61511>. Source data are provided with this paper.

Code availability

Sequence alignments, tree files, and custom scripts can be found at GitHub [https://github.com/dawickell/Isoetes_CAM].

Received: 22 June 2021; Accepted: 12 October 2021;

Published online: 03 November 2021

References

1. PPG1. A community-derived classification for extant lycophytes and ferns. *J. Syst. Evol.* **54**, 563–603 (2016).
2. Pigg, K. B. Isoetalean lycopsid evolution: from the Devonian to the present. *Am. Fern J.* **91**, 99–114 (2001).
3. Keeley, J. E. Distribution of diurnal acid metabolism in the genus *Isoetes* *Am. J. Bot.* **69**, 254–257 (1982).
4. Keeley, J. E. CAM photosynthesis in submerged aquatic plants. *Bot. Rev.* **64**, 121–175 (1998).
5. Aulio, K. Differential expression of diel acid metabolism in two life forms of *Littorella uniflora* (L.) Aschers. *N. Phytol.* **100**, 533–536 (1985).
6. Suissa, J. S. & Green, W. A. CO₂ starvation experiments provide support for the carbon-limited hypothesis on the evolution of CAM-like behaviour in *Isoetes*. *Ann. Bot.* **127**, 135–141 (2021).
7. Keeley, J. E. *Isoetes howellii*: a submerged aquatic CAM plant? *Am. J. Bot.* **68**, 420–424 (1981).
8. Wai, C. M. et al. Time of day and network reprogramming during drought induced CAM photosynthesis in *Sedum album*. *PLoS Genet.* **15**, e1008209 (2019).
9. Ming, R. et al. The pineapple genome and the evolution of CAM photosynthesis. *Nat. Genet.* **47**, 1435–1442 (2015).
10. Yang, X. et al. The *Kalanchoë* genome provides insights into convergent evolution and building blocks of crassulacean acid metabolism. *Nat. Commun.* **8**, 1899 (2017).
11. Cai, J. et al. The genome sequence of the orchid *Phalaenopsis equestris*. *Nat. Genet.* **47**, 65–72 (2015).
12. Zhang, L. et al. Origin and mechanism of crassulacean acid metabolism in orchids as implied by comparative transcriptomics and genomics of the carbon fixation pathway. *Plant J.* **86**, 175–185 (2016).
13. Heyduk, K. et al. Altered gene regulatory networks are associated with the transition from C3 to crassulacean acid metabolism in *Erycina* (Oncidiinae: Orchidaceae). *Front. Plant Sci.* **9**, 2000 (2018).
14. Heyduk, K. et al. Shared expression of crassulacean acid metabolism (CAM) genes pre-dates the origin of CAM in the genus *Yucca*. *J. Exp. Bot.* **70**, 6597–6609 (2019).
15. Abraham, P. E. et al. Transcript, protein and metabolite temporal dynamics in the CAM plant *Agave*. *Nat. Plants* **2**, 16178 (2016).
16. Banks, J. A. et al. The *Selaginella* genome identifies genetic changes associated with the evolution of vascular plants. *Science* **332**, 960–963 (2011).
17. Xu, Z. et al. Genome analysis of the ancient tracheophyte *Selaginella tamariscina* reveals evolutionary features relevant to the acquisition of desiccation tolerance. *Mol. Plant* **11**, 983–994 (2018).
18. VanBuren, R. et al. Extreme haplotype variation in the desiccation-tolerant clubmoss *Selaginella lepidophylla*. *Nat. Commun.* **9**, 13 (2018).
19. One Thousand Plant Transcriptomes Initiative. One thousand plant transcriptomes and the phylogenomics of green plants. *Nature* **574**, 679–685 (2019).
20. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biol.* **20**, 238 (2019).
21. Schmutz, J. et al. Genome sequence of the palaeopolyploid soybean. *Nature* **463**, 178–183 (2010).
22. Lang, D. et al. The *Physcomitrella patens* chromosome-scale assembly reveals moss genome structure and evolution. *Plant J.* **93**, 515–533 (2018).
23. Diop, S. I. et al. A pseudomolecule-scale genome assembly of the liverwort *Marchantia polymorpha*. *Plant J.* **101**, 1378–1396 (2020).
24. Li, F.-W. et al. *Anthoceros* genomes illuminate the origin of land plants and the unique biology of hornworts. *Nat. Plants* **6**, 259–272 (2020).
25. Szövényi, P., Gunadi, A. & Li, F.-W. Charting the genomic landscape of seed-free plants. *Nat. Plants* **7**, 554–565 (2021).
26. Li, Z. & Barker, M. S. Inferring putative ancient whole-genome duplications in the 1000 Plants (1KP) initiative: access to gene family phylogenies and age distributions. *Gigascience* **9**, gaa004 (2020).
27. Zwaenepoel, A. & Van de Peer, Y. Inference of ancient whole-genome duplications and the evolution of gene duplication and loss rates. *Mol. Biol. Evol.* **36**, 1384–1404 (2019).
28. Li, F.-W. et al. Fern genomes elucidate land plant evolution and cyanobacterial symbioses. *Nat. Plants* **4**, 460–472 (2018).
29. Sánchez, R. & Cejudo, F. J. Identification and expression analysis of a gene encoding a bacterial-type phosphoenolpyruvate carboxylase from *Arabidopsis* and rice. *Plant Physiol.* **132**, 949–957 (2003).
30. Deng, H. et al. Evolutionary history of PEPC genes in green plants: Implications for the evolution of CAM in orchids. *Mol. Phylogenet. Evol.* **94**, 559–564 (2016).
31. Ting, M. K. Y., She, Y.-M. & Plaxton, W. C. Transcript profiling indicates a widespread role for bacterial-type phosphoenolpyruvate carboxylase in malate-accumulating sink tissues. *J. Exp. Bot.* **68**, 5857–5869 (2017).
32. Blonde, J. D. & Plaxton, W. C. Structural and kinetic properties of high and low molecular mass phosphoenolpyruvate carboxylase isoforms from the endosperm of developing castor oilseeds. *J. Biol. Chem.* **278**, 11867–11873 (2003).

33. Michael, T. P. et al. Network discovery pipeline elucidates conserved time-of-day-specific cis-regulatory modules. *PLoS Genet.* **4**, e14 (2008).
34. Filichkin, S. A. et al. Global profiling of rice and poplar transcriptomes highlights key conserved circadian-controlled pathways and cis-regulatory modules. *PLoS ONE* **6**, e16907 (2011).
35. Michael, T. P. et al. Genome and time-of-day transcriptome of *Wolffia australiana* link morphological minimization with gene loss and less growth control. *Genome Res.* **31**, 225–238 (2020).
36. Steed, G., Ramirez, D. C., Hannah, M. A. & Webb, A. A. R. Chronoculture, harnessing the circadian clock to improve crop yield and sustainability. *Science* **372**, eabc9141 (2021).
37. Corellou, F. et al. Clocks in the green lineage: comparative functional analysis of the circadian architecture of the picoeukaryote *Ostreococcus*. *Plant Cell* **21**, 3436–3449 (2009).
38. Holm, K., Källman, T., Gyllenstrand, N., Hedman, H. & Lagercrantz, U. Does the core circadian clock in the moss *Physcomitrella patens* (Bryophyta) comprise a single loop? *BMC Plant Biol.* **10**, 109 (2010).
39. Ferrari, C. et al. Kingdom-wide comparison reveals the evolution of diurnal gene expression in Archaeplastida. *Nat. Commun.* **10**, 737 (2019).
40. Kuo, L.-Y., Huang, Y.-J., Chang, J., Chiou, W.-L. & Huang, Y.-M. Evaluating the spore genome sizes of ferns and lycophytes: a flow cytometry approach. *N. Phytol.* **213**, 1974–1983 (2017).
41. Ebihara, A. et al. Nuclear DNA, chloroplast DNA, and ploidy analysis clarified biological complexity of the *Vandenboschia radicans* complex (Hymenophyllaceae) in Japan and adjacent areas. *Am. J. Bot.* **92**, 1535–1547 (2005).
42. Praça-Fontes, M. M., Carvalho, C. R., Clarindo, W. R. & Cruz, C. D. Revisiting the DNA C-values of the genome size-standards used in plant flow cytometry to choose the “best primary standards”. *Plant Cell Rep.* **30**, 1183–1191 (2011).
43. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* **27**, 764–770 (2011).
44. Vurtture, G. W. et al. GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics* **33**, 2202–2204 (2017).
45. Kuo, L. Y. Polyploidy and biogeography in genus *Deparia* and phylogeography in *Deparia lancea*. PhD Thesis (2015).
46. Li, H. Minimap and minimiasm: fast mapping and de novo assembly for noisy long sequences. *Bioinformatics* **32**, 2103–2110 (2016).
47. Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Res.* **27**, 737–746 (2017).
48. Walker, B. J. et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS ONE* **9**, e112963 (2014).
49. Laetsch, D. R. & Blaxter, M. L. BlobTools: interrogation of genome assemblies. *F1000Res.* **6**, 1287 (2017).
50. Xu, G.-C. et al. LR_GapCloser: a tiling path-based gap closer that uses long reads to complete genome assembly. *Gigascience* **8**, gyl157 (2019).
51. Roach, M. J., Schmidt, S. A. & Borneman, A. R. Purge Haplotigs: allelic contig reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics* **19**, 460 (2018).
52. Jin, J.-J. et al. GetOrganelle: a fast and versatile toolkit for accurate de novo assembly of organelle genomes. *Genome Biol.* **21**, 241 (2020).
53. Ou, S. & Jiang, N. LTR_retriever: a highly accurate and sensitive program for identification of long terminal repeat retrotransposons. *Plant Physiol.* **176**, 1410–1422 (2018).
54. Smit, A. F. A. & Hubley, R. RepeatModeler Open-1.0. <http://www.repeatmasker.org> (2008).
55. Smit, A. F. A., Hubley, R. & Green, P. RepeatMasker Open-4.0. <http://www.repeatmasker.org> (2014).
56. Stanke, M. & Morgenstern, B. AUGUSTUS: a web server for gene prediction in eukaryotes that allows user-defined constraints. *Nucleic Acids Res.* **33**, W465–W467 (2005).
57. Korf, I. Gene finding in novel genomes. *BMC Bioinformatics* **5**, 59 (2004).
58. Brůna, T., Hoff, K. J., Lomsadze, A., Stanke, M. & Borodovsky, M. BRAKER2: automatic eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database. *NAR Genomics Bioinformatics* **3**, lqaa108 (2021).
59. Grabherr, M. G. et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat. Biotechnol.* **29**, 644–652 (2011).
60. Holt, C. & Yandell, M. MAKER2: an annotation pipeline and genome-database management tool for second-generation genome projects. *BMC Bioinformatics* **12**, 491 (2011).
61. Huerta-Cepas, J. et al. Fast genome-wide functional annotation through orthology assignment by eggNOG-Mapper. *Mol. Biol. Evol.* **34**, 2115–2122 (2017).
62. Pertea, M. et al. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat. Biotechnol.* **33**, 290–295 (2015).
63. Amborella Genome Project. The Amborella genome and the evolution of flowering plants. *Science* **342**, 1241089 (2013).
64. Lamesch, P. et al. The Arabidopsis Information Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Res.* **40**, D1202–D1210 (2012).
65. International Brachypodium Initiative. Genome sequencing and analysis of the model grass *Brachypodium distachyon*. *Nature* **463**, 763–768 (2010).
66. Yang, Y. et al. Prickly waterlily and rigid hornwort genomes shed light on early angiosperm evolution. *Nat. Plants* **6**, 215–222 (2020).
67. Tang, H. et al. An improved genome release (version Mt4.0) for the model legume *Medicago truncatula*. *BMC Genomics* **15**, 312 (2014).
68. Ming, R. et al. Genome of the long-living sacred lotus (*Nelumbo nucifera* Gaertn.). *Genome Biol.* **14**, R41 (2013).
69. Zhang, L. et al. The water lily genome and the early evolution of flowering plants. *Nature* **577**, 79–84 (2020).
70. Nystedt, B. et al. The Norway spruce genome sequence and conifer genome evolution. *Nature* **497**, 579–584 (2013).
71. Wang, W. et al. The *Spirodela polyrhiza* genome reveals insights into its neoteneous reduction fast growth and aquatic lifestyle. *Nat. Commun.* **5**, 3311 (2014).
72. Ibarra-Laclette, E. et al. Architecture and evolution of a minute plant genome. *Nature* **498**, 94–98 (2013).
73. Jaillon, O. et al. The grapevine genome sequence suggests ancestral hexaploidization in major angiosperm phyla. *Nature* **449**, 463–467 (2007).
74. Olsen, J. L. et al. The genome of the seagrass *Zostera marina* reveals angiosperm adaptation to the sea. *Nature* **530**, 331–335 (2016).
75. Cheng, S. et al. Genomes of subaerial Zygnematophyceae provide insights into land plant evolution. *Cell* **179**, 1057–1067.e14 (2019).
76. Conway, J. R., Lex, A. & Gehlenborg, N. UpSetR: an R package for the visualization of intersecting sets and their properties. *Bioinformatics* **33**, 2938–2940 (2017).
77. Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**, 907–915 (2019).
78. Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
79. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
80. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
81. Zwaenepoel, A. & Van de Peer, Y. wgd—simple command line tools for the analysis of ancient whole-genome duplications. *Bioinformatics* **35**, 2153–2155 (2019).
82. Qi, X. et al. A well-resolved fern nuclear phylogeny reveals the evolution history of numerous transcription factor families. *Mol. Phylogenet. Evol.* **127**, 961–977 (2018).
83. Yang, T. & Liu, X. Comparative transcriptome analysis of *Isoetes sinensis* under terrestrial and submerged conditions. *Plant Mol. Biol. Rep.* **34**, 136–145 (2016).
84. Hetherington, A. J., Emms, D. M., Kelly, S. & Dolan, L. Gene expression data support the hypothesis that *Isoetes* rootlets are true roots and not modified leaves. *Sci. Rep.* **10**, 21547 (2020).
85. Wood, D., Besnard, G., Beerling, D. J., Osborne, C. P. & Christin, P.-A. Phylogenomics indicates the “living fossil” *Isoetes* diversified in the Cenozoic. *PLoS ONE* **15**, e0227525 (2020).
86. Xie, Y. et al. SOAPdenovo-Trans: de novo transcriptome assembly with short RNA-Seq reads. *Bioinformatics* **30**, 1660–1666 (2014).
87. Barker, M. S. et al. Multiple paleopolyploidizations during the evolution of the Compositae reveal parallel patterns of duplicate gene retention after millions of years. *Mol. Biol. Evol.* **25**, 2445–2455 (2008).
88. Barker, M. S. et al. EvoPipes.net: bioinformatic tools for ecological and evolutionary genomics. *Evol. Bioinformatics* **6**, 143–149 (2010).
89. Birney, E., Clamp, M. & Durbin, R. GeneWise and Genomewise. *Genome Res.* **14**, 988–995 (2004).
90. Goodstein, D. M. et al. Phytozome: a comparative platform for green plant genomics. *Nucleic Acids Res.* **40**, D1178–D1186 (2012).
91. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586–1591 (2007).
92. Loytynoja, A. & Goldman, N. Phylogeny-aware gap placement prevents errors in sequence alignment and evolutionary analysis. *Science* **320**, 1632–1635 (2008).
93. Ronquist, F. et al. MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Syst. Biol.* **61**, 539–542 (2012).
94. Szollosi, G. J., Tannier, E., Lartillot, N. & Daubin, V. Lateral gene transfer from the dead. *Syst. Biol.* **62**, 386–397 (2013).
95. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R. *Bioinformatics* **35**, 526–528 (2019).

96. Morris, J. L. et al. The timescale of early land plant evolution. *Proc. Natl Acad. Sci. USA* **115**, E2274–E2283 (2018).
97. Larsén, E. & Rydin, C. Disentangling the phylogeny of *Isoetes* (Isoetales), using nuclear and plastid data. *Int. J. Plant Sci.* **177**, 157–174 (2016).
98. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
99. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermini, L. S. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* **14**, 587–589 (2017).
100. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**, 268–274 (2015).
101. Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2: improving the ultrafast bootstrap approximation. *Mol. Biol. Evol.* **35**, 518–522 (2018).
102. de Vries, S. et al. The evolution of the phenylpropanoid pathway entailed pronounced radiations and divergences of enzyme families. *Plant J.* <https://doi.org/10.1111/tpj.15387> (2021).
103. Slotte, T. et al. The *Capsella rubella* genome and the genomic consequences of rapid mating system evolution. *Nat. Genet.* **45**, 831–835 (2013).
104. Wan, T. et al. A genome for gnetophytes and early evolution of seed plants. *Nat. Plants* **4**, 82–89 (2018).
105. Sierro, N. et al. The tobacco genome sequence and its comparison with those of tomato and potato. *Nat. Commun.* **5**, 1–9 (2014).
106. Ouyang, S. et al. The TIGR rice genome annotation resource: improvements and new features. *Nucleic Acids Res.* **35**, D883–D887 (2007).
107. Argout, X. et al. The genome of *Theobroma cacao*. *Nat. Genet.* **43**, 101–108 (2011).
108. Wang, S. et al. Genomes of early-diverging streptophyte algae shed light on plant terrestrialization. *Nat. Plants* **6**, 95–106 (2020).
109. Nishiyama, T. et al. The Chara genome: secondary complexity and implications for plant terrestrialization. *Cell* **174**, 448–464.e24 (2018).
110. Hori, K. et al. *Klebsormidium flaccidum* genome reveals primary factors for plant terrestrial adaptation. *Nat. Commun.* **5**, 3978 (2014).
111. Jiao, C. et al. The *Penium margaritaceum* genome: hallmarks of the origins of land plants. *Cell* **181**, 1097–1111.e12 (2020).
112. de Vries, J. et al. Heat stress response in the closest algal relatives of land plants reveals conserved stress signaling circuits. *Plant J.* **103**, 1025–1048 (2020).
113. de Vries, J., Curtis, B. A., Gould, S. B. & Archibald, J. M. Embryophyte stress signaling evolved in the algal progenitors of land plants. *Proc. Natl Acad. Sci. USA* **115**, E3471–E3480 (2018).
114. Ju, C. et al. Conservation of ethylene as a plant hormone over 450 million years of evolution. *Nat. Plants* **1**, 14004 (2015).
115. Moreau, H. et al. Gene functionalities and genome structure in *Bathycoccus prasinos* reflect cellular specializations at the base of the green lineage. *Genome Biol.* **13**, R74 (2012).
116. Merchant, S. S. et al. The *Chlamydomonas* genome reveals the evolution of key animal and plant functions. *Science* **318**, 245–250 (2007).
117. Blanc, G. et al. The genome of the polar eukaryotic microalga *Coccomyxa subellipsoidea* reveals traits of cold adaptation. *Genome Biol.* **13**, R39 (2012).
118. Worden, A. Z. et al. Green evolution and dynamic adaptations revealed by genomes of the marine picoeukaryotes *Micromonas*. *Science* **324**, 268–272 (2009).
119. Palenik, B. et al. The tiny eukaryote *Ostreococcus* provides genomic insights into the paradox of plankton speciation. *Proc. Natl Acad. Sci. USA* **104**, 7705–7710 (2007).
120. De Clerck, O. et al. Insights into the evolution of multicellularity from the sea lettuce genome. *Curr. Biol.* **28**, 2921–2933.e5 (2018).
121. Prochnik, S. E. et al. Genomic analysis of organismal complexity in the multicellular green alga *Volvox carteri*. *Science* **329**, 223–226 (2010).
122. de Vries, J., de Vries, S., Slamovits, C. H., Rose, L. E. & Archibald, J. M. How embryophytic is the biosynthesis of phenylpropanoids and their derivatives in *Streptophyte* algae? *Plant Cell Physiol.* **58**, 934–945 (2017).
123. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
124. Minh, B. Q. et al. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
125. Hu, Y. et al. Crystal structures of a *Populus tomentosa* 4-coumarate:CoA ligase shed light on its enzymatic mechanisms. *Plant Cell* **22**, 3093–3104 (2010).
126. Pan, H. et al. Structural studies of cinnamoyl-CoA reductase and cinnamyl-alcohol dehydrogenase, key enzymes of monolignol biosynthesis. *Plant Cell* **26**, 3709–3727 (2014).
127. Louie, G. V. et al. Structure-function analyses of a caffeic acid O-methyltransferase from perennial ryegrass reveal the molecular basis for substrate preference. *Plant Cell* **22**, 4114–4127 (2010).
128. Youn, B. et al. Crystal structures and catalytic mechanism of the Arabidopsis cinnamyl alcohol dehydrogenases AtCAD5 and AtCAD4. *Org. Biomol. Chem.* **4**, 1687–1697 (2006).
129. Ferrer, J.-L., Zubieta, C., Dixon, R. A. & Noel, J. P. Crystal structures of alfalfa caffeoyl coenzyme A 3-O-methyltransferase. *Plant Physiol.* **137**, 1009–1017 (2005).
130. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140 (2010).
131. Ritchie, M. E. et al. limma powers differential expression analyses for RNA-seq and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
132. Barbier, F. F. et al. A phenol/chloroform-free method to extract nucleic acids from recalcitrant, woody tropical species for gene expression and sequencing. *Plant Methods* **15**, 62 (2019).
133. Love, M., Anders, S. & Huber, W. Differential analysis of count data—the DESeq2 package. *Genome Biol.* **15**, 10–1186 (2014).
134. Conesa, A., Nueda, M. J., Ferrer, A. & Talón, M. maSigPro: a method to identify significantly differential expression profiles in time-course microarray experiments. *Bioinformatics* **22**, 1096–1102 (2006).
135. MacKinnon, K. J. M., Cole, B. J., Yu, C. & Coomey, J. H. Changes in ambient temperature are the prevailing cue in determining *Brachypodium distachyon* diurnal gene regulation. *N. Phytol* **227**, 1709–1724 (2020).
136. Michael, T. P. et al. A morning-specific phytohormone gene expression program underlying rhythmic plant growth. *PLoS Biol.* **6**, e225 (2008).
137. Mockler, T. C. et al. The DIURNAL project: DIURNAL and circadian expression profiling, model-based pattern matching, and promoter analysis. *Cold Spring Harb. Symp. Quant. Biol.* **72**, 353–363 (2007).
138. McLeay, R. C. & Bailey, T. L. Motif enrichment analysis: a unified framework and an evaluation on ChIP data. *BMC Bioinformatics* **11**, 165 (2010).
139. Franco-Zorrilla, J. M. et al. DNA-binding specificities of plant transcription factors and their potential to define target genes. *Proc. Natl Acad. Sci. USA* **111**, 2367–2372 (2014).
140. Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif. *Bioinformatics* **27**, 1017–1018 (2011).

Acknowledgements

We would like to thank Karolina Heyduk, Peter Schafran, and Arthur Zwaenepoel for their advice and support regarding various aspects of this project; Tai-Chung Wu, Wen-Yuan Kao, Ta-Chun Lin, and Chun-Neng Wang for their help on time-course experiments. J.d.V. is supported through funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant no. 852725; ERC Starting Grant 'TerreStriAL'), and MaLand (DFG priority programme 2237; VR132/4-1). A.D.A. and A.D. are grateful for being supported through the International Max Planck Research School (IMPRS) for Genome Science.

Author contributions

D.W., L.-Y.K., T.P.M. and F.-W.L. coordinated the project. Y.-M.H. provided the plant materials. L.-Y.K. carried out the time-course experiment and nucleic acid extraction. T.P.M. and F.-W.L. sequenced and assembled the genome. H.-P.Y. and F.-W.L. annotated the genome. D.W. assembled the plastome and profiled RNA-editing. D.W. circumscribed gene families and examined genes related to stomata and root development. A.D.A., I.I., A.D., S.d.V. and J.d.V. characterized lignin biosynthesis genes. D.W., Z.L. and M.S.B. carried out WGD analysis. D.W. analyzed expressions of the CAM pathway genes. N.T.H. and T.P.M. carried out HAYSTACK and ELEMENT analyses. D.W., T.P.M. and F.-W.L. synthesized and wrote the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-021-26644-7>.

Correspondence and requests for materials should be addressed to Todd P. Michael or Fay-Wei Li.

Peer review information *Nature Communications* thanks Marek Mutwil, Xiaohan Yang and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available.

Reprints and permission information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021