

RNA-seq highlights parallel and contrasting patterns in the evolution of the nuclear genome of holo-mycoheterotrophic plants

M.I. Schelkunov^{1,2,*}, A.A. Penin^{1,2,3}, M.D. Logacheva^{1,2,*}

1 - Institute for Information Transmission Problems, Russian Academy of Sciences, Moscow, Russia

2 - Lomonosov Moscow State University, A.N Belozersky Institute of Physico-Chemical Biology, Moscow, Russia

3 - Lomonosov Moscow State University, Faculty of Biology, Moscow, Russia

* - corresponding authors, shelkmike@gmail.com, maria.log@gmail.com

Summary

- While photosynthesis is the most notable trait of plants, several lineages of plants (so-called holo-heterotrophs) have adapted to obtain organic compounds from other sources. The switch to heterotrophy leads to profound changes at the morphological, physiological and genomic levels.
- Here, we characterize the transcriptomes of three species representing two lineages of mycoheterotrophic plants: orchids (*Epipogium aphyllum* and *Epipogium roseum*) and Ericaceae (*Hypopitys monotropa*). Comparative analysis is used to highlight the parallelism between distantly related holo-heterotrophic plants.
- In both lineages, we observed genome-wide elimination of nuclear genes that encode proteins related to photosynthesis, while systems associated with protein import to plastids as well as plastid transcription and translation remain active. Genes encoding components of plastid ribosomes that have been lost from the plastid genomes have not been transferred to the nuclear genomes; instead, some of the encoded proteins have been substituted by homologs. The nuclear genes of both *Epipogium* species accumulated mutations twice as rapidly as their photosynthetic relatives; in contrast, no increase in the substitution rate was observed in *H.monotropa*.
- Holo-heterotrophy leads to profound changes in nuclear gene content. The observed increase in the rate of nucleotide substitutions is lineage specific, rather than a universal phenomenon among non-photosynthetic plants.

Introduction

The capability for photosynthesis is the iconic trait of plants and is of the highest importance to the biosphere. However, some plants, including several thousands of flowering plant species, obtain organic substances from sources other than photosynthesis (Merckx *et al.*, 2009; Westwood *et al.*, 2010). These plants acquire organic compounds either from associated fungi (myco-heterotrophy) or by parasitizing other plants. Most of these species combine photosynthesis and heterotrophy, but several hundred species have totally lost photosynthetic ability and become completely heterotrophic (holo-heterotrophs). The acquisition of heterotrophic ability has occurred in the evolutionary history of plants more than 50 times (Merckx *et al.*, 2009; Westwood *et al.*, 2010). The switch to heterotrophy leads to profound changes at the phenotypic level (reduction of leaves, loss of green colour, reduction of the vegetation period) that are highly parallel in different lineages. The genotypic alterations that underlie these changes are for the most part unclear. The difficulty of cultivating heterotrophic plants under experimental conditions hampers classic genetic and physiological studies. Advances in DNA sequencing permit the application of a genomic approach to elucidate the genetic changes associated with heterotrophy.

Genetic and genomic studies of heterotrophic plants are currently focused on two aspects. The first is the interaction of parasitic plants with their hosts and their adaptations to parasitism (e.g., Yang *et al.*, 2015). Extensive exchange of transcripts occurs between hosts and parasites (Kim *et al.*, 2014). On an evolutionary scale, a large number of horizontal gene transfer (HGT) events from hosts to parasites have been found in the organellar and nuclear genomes of parasitic plants (Bellot *et al.*, 2016; Li *et al.*, 2013). A recent large-scale survey of HGT in Orobanchaceae showed that the number of these events correlates positively with the degree of heterotrophy (Yang *et al.*, 2016). The second aspect is the evolution of organellar (mainly plastid) genomes. As expected, the

plastomes of heterotrophic plants are reduced in size due to the loss of genes related to photosynthesis, and may retain only ~ 7.5% of the length of a typical plastome (Bellot & Renner, 2015). Despite the high degree of reduction, the plastomes of non-photosynthetic plants retain genes whose products are involved in translation, specifically transfer RNAs, components of the plastid ribosome and two other genes, *accD* and *clpP* (e.g., Wicke *et al.*, 2013; Barrett *et al.*, 2014; Schelkunov *et al.*, 2015; Bellot & Renner, 2015; Lam *et al.*, 2016). Although there are several considerations regarding the retention of the plastome, complete loss of the plastome is also apparently possible. At present, there are two known cases of such loss, one in *Polytomella*, a genus of unicellular algae (Smith & Lee, 2014), and one in the parasitic angiosperm *Rafflesia lagascae* (Molina *et al.*, 2014), although alternative explanations are possible in the latter case. Much less information on the mitochondrial genomes of non-photosynthetic plants is available, although there are indications that these genomes are not as extensively reduced in size (Fan *et al.*, 2016).

The changes in the nuclear genomes of non-photosynthetic plants have not been well studied. To date, the only work to deeply analyse the nuclear genomes of holo-heterotrophic plants was performed in Orobanchaceae, where a holo-heterotrophic species, *Orobanche aegyptiaca*, was compared with two of its relatives, one hemi-autotroph with obligatory parasitism and one hemi-autotroph with facultative parasitism (Wickett *et al.*, 2011). Surprisingly, the authors found evidence for conservation of the pathways responsible for chlorophyll synthesis. Additionally, in one study of the *Hypopitys monotropa* plastome, transcriptome analysis showed that many genes related to photosynthesis have been lost (Ravin *et al.*, 2016).

To obtain a more detailed understanding of the evolution of holo-heterotrophic plants, in this work, we analyse the nuclear genomes of *Epipogium aphyllum* and *E.roseum* (Orchidaceae) and *H.monotropa* (Ericaceae)

using transcriptome sequencing. These plants are good models for studying the characteristics of holo-heterotrophic plants since their plastomes are among the most reduced in size (Schelkunov *et al.*, 2015; Logacheva *et al.*, 2016) (19-35 Kb versus approximately 150 Kb in typical photosynthetic plants). Therefore, we expect that the nuclear genomes of these species may also differ profoundly from the genomes of photosynthetic plants and will therefore allow us to highlight characteristics that are specific to the genomes of holo-heterotrophic species.

Materials and methods

Sample collection, library preparation and sequencing

Information on the specimens used for transcriptome sequencing was reported earlier by Schelkunov *et al.* 2015 for *E.aphyllum* ("White Sea" sample) and *E.roseum* ("Vietnam 2" sample) and by Logacheva *et al.* 2016 for *H.monotropa*. RNA was extracted using the Qiagen RNeasy Plant Mini kit (Qiagen, Netherlands). To allow the characterization of non-polyadenylated transcripts (e.g., plastid and mitochondrial), we prepared libraries from ribosomal RNA-depleted samples using the Ribo-Zero plant leaf kit (Illumina, USA). Detailed information on the libraries and sequencing settings, as well as links to on-line databases where the reads and assembled sequences are deposited, is provided in Table S1.

Choice of datasets of photosynthetic plants

To compare the holo-heterotrophic species with typical photosynthetic plants, we obtained data from several RNA-seq experiments and from the genome assemblies of the plants that were available as of 2015. Although many transcriptomes have been sequenced, the corresponding data exhibit certain shortcomings, such as a low sequencing depth, the use of technologies that introduce errors in homopolymer regions, or the sequencing of organs (e.g., roots) in which the set of expressed genes may differ from that in above-ground organs. In view of

these considerations, we chose six photosynthetic species, including three for comparison with *Epipogium* and three for comparison with *H.monotropa*. The names and sources of the datasets on photosynthetic plants are listed in Tables S1 and S2. In the case of RNA-seq data, we assembled the transcriptomes using exactly the same methods that were used for the transcriptomes of *Epipogium* and *Hypopitys*.

Transcriptome assembly and postprocessing

Reads were trimmed with Trimmomatic 0.32 (Bolger *et al.*, 2014). Trimming involved the removal of sequencing adapters, bases with a Phred quality of less than 3 from both the 5' and 3' ends of reads, and bases from the 5' ends of reads starting from a region with 5 consecutive bases with an average score of less than 10 (trimming with a sliding window). Additionally, reads that had average quality of less than 20 after trimming and reads that after trimming became shorter than 30 bp were removed. Assembly was performed using Trinity version r20140717 (Grabherr *et al.*, 2011). Digital normalization of coverage to 50× was switched on. After assembly, we performed filtration, removing minor isoforms, low-coverage transcripts and contamination. We defined major isoforms as the isoforms to which the highest number of reads (estimated using RSEM (Li & Dewey, 2011)) were mapped relative to other isoforms. After removing minor isoforms, we filtered low-coverage transcripts by mapping reads to contigs using CLC Assembly Cell 4.2 and retaining only transcripts with an average coverage of at least 3. Potential CDSs in the transcripts were then determined using TransDecoder version r20140704 (Haas *et al.*, 2013). The criteria for considering an ORF to be a potential CDS consisted not only of hexanucleotide frequencies, which are employed in TransDecoder by default, but also all ORFs that possessed domains from the Pfam-A or Pfam-B databases. The minimum CDS length was decreased from the default of 100 amino acids to 30 amino acids. When there were several potential CDSs

from a transcript, only the longest one was taken. Then, to remove contamination, we performed BLASTP 2.2.29+ (Camacho *et al.*, 2009) alignment (maximum allowed e-value 10^{-5} , word size 3, the low-complexity filter switched off) of the translated CDS sequences against the NCBI NR database and against proteins of related plants whose genomes had been sequenced. A CDS was considered to represent contamination if the BLASTP match with the highest bit score was to a species that was not from Streptophyta. Conversely, if the best match was to Streptophyta member or the sequence presented no BLASTP matches, the CDS was considered to belong to a plant. Statistical parameters of the assemblies, such as N50 values and the number of sequences longer than 1000 bp, were calculated using custom scripts after each filtering step described above. Estimation of the completeness of the assemblies was performed with CEGMA 2.5 (Parra *et al.*, 2007). To evaluate gene expression without contamination, we calculated FPKM values using RSEM for transcripts that contained CDSs that were determined not to have arisen from contamination. Minor isoforms of the transcripts were also used for the analysis; the FPKM values provided in Table S3 are the sums of all gene isoforms.

Transcriptome annotation and Gene Ontology analysis

Three types of annotations of the CDSs were performed. The first was the computation of 1-1 orthologs from the CDSs of each species by reciprocal alignment with *Arabidopsis thaliana* proteins from TAIR10 database (Berardini *et al.*, 2015) using BLASTP (parameters are the same as indicated above). Also we performed GO annotation using results of BLASTP alignment of all proteins against the NCBI NR (parameters as above) and annotated them with B2G4Pipe 2.5, a command-line wrapper for Blast2GO (Conesa *et al.*, 2005). Then we assigned KEGG Orthology identifiers to the CDSs with the GhostKoala server (Kanehisa *et al.*, 2016b). The orthogroups in the studied species were calculated

separately for Orchidaceae and Ericales using OrthoMCL 2.09 (Li *et al.*, 2003) (inflation parameter 1.5).

For GO term enrichment analysis, we used a set of custom scripts written in Perl and R. Utilizing the results of the GO annotation performed with Blast2GO and a graph of GO terms (<http://purl.obolibrary.org/obo/go/go-basic.obo>, last accessed 15 April 2015), the scripts provide all GO terms corresponding to a gene, including terms that are paternal (i.e., with higher hierarchical levels) to those provided by Blast2GO. Comparison of the numbers of genes with specific GO terms between a pair of species was then performed via a series of Fisher's exact tests (one for each GO term). GO terms that were not assigned to any genes in either species were excluded from the analysis. After the Fisher's tests were performed, the scripts performed group comparisons between holo-heterotrophic and photosynthetic plants for each GO term, by taking all p-values for each pairwise comparison between pairs of species within the group (Ericales or Orchidaceae). The scripts then conducted Bonferroni correction for multiple testing to evaluate the statistical significance of the differences between these groups separately for Orchidaceae and Ericales. Next, another round of correction for multiple hypothesis testing was performed, taking into account the fact that an individual Fisher's test was performed for each GO term. This correction was performed via the method of Benjamini-Yekutieli (Yekutieli & Benjamini, 2001).

Since our goal was the analysis of nuclear genes, we excluded genes encoded in the plastid and mitochondrial genomes from the enrichment analysis by excluding all genes that were 1-1 orthologs of plastid and mitochondrial genes of *Arabidopsis thaliana*. Some plastid and mitochondrial transcripts may not have been present in the assemblies. To compensate for this, in Table S3, we indicated a plastid gene as present in the plastome of a studied species not only when it was found through the 1-1 ortholog method, but also when it was present in the annotation of the plastome of that species provided in

GenBank (accessions are provided in Table S2). Because the plastome of *O. italica* is not available, we employed the plastome of *Habenaria pantlingiana*, the only species with a characterized plastome from subtribe Orchidiae, which includes *O. italica*.

Analysis of substitution rates and selective pressure

To estimate the average selective pressure and substitution rates in the genomes of the studied species, we concatenated the sequences of all genes from orthogroups containing exactly one gene from each species. The concatenated sequences were then aligned using TranslatorX 1.1 (Abascal *et al.*, 2010). As a tool for the alignment of the amino acid sequences by TranslatorX, we used Muscle (Edgar, 2004). The topologies of the phylogenetic trees for Orchidaceae and Ericales were obtained from the APG III classification (THE ANGIOSPERM PHYLOGENY GROUP, 2009). dN/dS and substitution rate estimation was performed based on the alignment and the tree using PAML 4.8 (Yang, 2007), employing a branch model with free ratios, the Gy94_3×4 codon model and removal of all columns with at least one gap in the alignment.

To compare the magnitude of the selective pressure acting on individual genes in heterotrophic and autotrophic plants, the sequences of the genes were aligned using TranslatorX and Muscle as described above, and dN/dS ratios were calculated using PAML in the branch model. For the Orchidaceae, two calculations were performed. In the first calculation, one dN/dS was allowed for the branches of autotrophic plants; another dN/dS was allowed for the terminal branches of *E.aphyllum* and *E.roseum*; and a third dN/dS was allowed for the branch of their common ancestor. The second calculation was performed allowing one dN/dS for the common ancestor of *E.aphyllum* and *E.roseum* and a second dN/dS for all other branches, including autotrophic species, *E.aphyllum* and *E.roseum*. The P-values for the difference in dN/dS between the terminal branches of *Epipogium* and the

branches of autotrophic species were calculated using the likelihood ratio test. Allowing an individual dN/dS for the branch of the common ancestor of *Epipogium* was beneficial, as photosynthetic ability was lost on that branch, and it is unclear whether to group it with autotrophic or heterotrophic species. Similar calculations were performed for Ericales, allowing different dN/dS ratios for the *H.monotropa* branch and other branches the first time and demanding a single dN/dS for all branches the second time. Since *H.monotropa* was the only heterotrophic species from Ericales employed in this study, its terminal branch partially included its autotrophic ancestor. Thus, the values of dN/dS provided for *H.monotropa* in Table S3 describe selective pressure partially before and partially after the loss of photosynthetic ability. dN/dS was calculated only for genes that were present in all 5 species of Orchidaceae or in all 4 species of Ericales.

Analysis of protein targeting to organelles

To estimate the number of proteins targeted to various organelles, we predicted transit peptides using TargetP 1.1 (Emanuelsson *et al.*, 2000) and DualPred (Saravanan & Velan Lakshmi, 2015). Unlike TargetP, DualPred is also capable of predicting proteins that are dually targeted to mitochondria and plastids. TargetP classifies predictions into five "reliability classes", where the most confident predictions belong to the first class, and the least confident predictions belong to the fifth class. We considered a protein to be potentially targeted to an organelle if its transit peptide exhibited a reliability class of four or less. For DualPred, we considered a protein to be dually targeted to plastids and mitochondria if it presented a dual-targeting score of at least 0.5, as suggested by the author of DualPred in a personal communication. Prediction of transit peptides was performed only for proteins whose genes exhibited a completely assembled 5'-end according to TransDecoder and were assigned at least one GO term.

To predict the targeting of ribosomal proteins, we

utilized a more elaborate technique. It has been demonstrated that some proteins are dually targeted to plastids and mitochondria not because they have one transit peptide that allows them to enter both organelles, but because their mRNAs exhibit alternative translation start sites, resulting in proteins with different transit peptides (Mitschke *et al.*, 2009). To search for transit peptides that may originate from alternative translation, we truncated all of the ribosomal proteins under analysis before the first methionine occurring after the first 25 amino acids (as in Mitschke *et al.*, 2009) and performed TargetP and DualPred analyses for these shortened versions of the proteins as well. Additionally, all alternative isoforms of the ribosomal proteins were analysed. When DualPred suggested that a protein was dually targeted and TargetP suggested either plastid or mitochondrial targeting, we considered this to represent a non-contradictory prediction suggesting dual targeting.

The Pfam families to which the ribosomal proteins belonged were determined by aligning *Arabidopsis*

thaliana proteins to all families from the Pfam database, version 31.0 (Finn *et al.*, 2016), using the HMMER server (Finn *et al.*, 2015). To identify proteins that belonged to these families in *E.aphyllum*, *E.roseum* and *H.monotropa*, we conducted a search with hidden Markov models of these families using the hmmscan tool from HMMER package, version 3.1b1 (Eddy, 1998). The search was performed among all of the CDSs from transcripts, not only the longest ones, to make detection in polycistronic mitochondrial mRNAs possible. Orthologous proteins in the photosynthetic species employed for comparison were determined as belonging to orthogroups that contained the proteins identified for *E.aphyllum*, *E.roseum* and *H.monotropa*.

Graphic representation of results

Phylogenetic trees were drawn with TreeGraph 2.9.2 (Stöver & Müller, 2010). Maps of metabolic and signalling pathways were built the KEGG site (Kanehisa *et al.*, 2016a) (accessed 22 October 2015).

Table 1. Brief statistics of the transcriptome assemblies.

| | <i>Epipogium aphyllum</i> | | <i>Epipogium roseum</i> | | <i>Hypopitys monotropa</i> | |
|--|---------------------------|------------|-------------------------|------------|----------------------------|------------|
| | All transcripts | CDSs* | All transcripts | CDSs* | All transcripts | CDSs* |
| Number of sequences | 992 338 | 20 958 | 1 336 170 | 19 026 | 217 166 | 13 276 |
| Number of sequences longer or equal to 1000 bp | 39 403 | 4 321 | 28 284 | 3 259 | 33 560 | 5 421 |
| Total length of sequences | 290 947 719 | 12 988 731 | 321 048 790 | 10 721 418 | 116 614 033 | 13 496 709 |
| N50 | 371 | 1 173 | 257 | 990 | 1 199 | 1470 |
| Median length of sequences | 178 | 306 | 164 | 315 | 237 | 807 |

*Here, the CDS is defined as the longest ORF in the major isoforms of the transcripts that were assigned at least one GO term (low-coverage transcripts and contaminating transcripts are not considered).

Results and Discussion

Characteristics of transcriptome assembly

The statistics of transcriptome assembly are provided in Table 1 (for details, see Table S4). Simple statistical measures, such as N50 values and the mean contig length, were smaller than those in many studies involving transcriptome assemblies, mainly because we analysed contigs with a minimum length of 100 bp instead of 200 bp (the default cut-off in the Trinity assembler). Several genes of interest (for example, the plastid genes *rpl32* and *rpl36*) are shorter than 200 bp, and use of the default cut-off would incur a risk of missing their transcripts. The assembly statistics were also biased by contamination. The symbiosis of mycoheterotrophic plants with fungi may be quite deep; for example, the rhizome of *E.aphyllum* has been described as "heavily and permanently infected" (Rasmussen, 1995). Despite the fact that we sequenced RNA from the above-ground parts of the plants, we observed a large number of bacterial and fungal transcripts in the assemblies (Fig. S1), especially in *E.roseum*, in which they represented ~80% of the total contigs. In *E.aphyllum* and *H.monotropa*, the corresponding values were approximately 10% and 5%, respectively. This difference may reflect the correlation between soil microbiome biomass and climate (*E.aphyllum* and *H.monotropa* were collected from colder regions than *E.roseum*). Among the transcriptomes of photosynthetic plants used for comparison, the transcriptome of *Orchis italica* was also highly contaminated. Sequences originating from fungi and bacteria were discarded prior to further analysis. Estimation of the completeness of the assemblies based on the set of genes that are expected to be present in all eukaryotic genomes (Parra *et al.*, 2007) showed that > 90% of these genes were at least partially assembled (Table S5). Notably, the genomes were assembled less completely than the transcriptomes on average, with a median of 95.6% of the genes being assembled at least partially in the transcriptome assemblies, compared with 90.8% of the genes in the

genome assemblies. These results show that, given sufficient coverage, RNA-seq is as good as complete genome sequencing in terms of the number of retrieved genes, while being less costly and using an assembly process that is computationally faster.

Gene retention and reduction

Plastid-targeted proteins carry specific amino acid sequences known as targeting signals or transit peptides that interact with the translocon system and enable the import of these proteins into plastids. Thus, we expected that in non-photosynthetic plants, the fraction of proteins with plastid-targeting signals relative to other proteins will be decreased. However, a comparison revealed a more complex situation. In *Epipogium*, the fraction of proteins targeted to plastids relative to the total number of proteins is approximately 2 times lower than in autotrophic orchids on average, whereas the fraction of proteins that are targeted to mitochondria or the endoplasmic reticulum is approximately the same. In contrast, in *H.monotropa*, the fraction of proteins targeted to plastids does not appear to be decreased. However, because the fraction of plastid-targeted proteins differs greatly between the two photosynthetic Ericales species that we used for comparison (Table S6), these results should be treated with caution, as they may be biased by lineage-specific genome duplications and/or by differences in the quality of the assemblies. To obtain a deeper understanding of the patterns of gene reduction, we performed Gene Ontology (GO) enrichment analyses. GO analysis of *Epipogium aphyllum* and *E.roseum* versus three photosynthetic orchids revealed that the genes associated with 60 GO terms in *Epipogium* were underrepresented, while the genes associated with 38 terms were overrepresented, with q -values ≤ 0.05 . All of the overrepresented GO terms are related to genes associated with mobile elements. This result is presumably caused by methodological differences, as the *Epipogium* transcriptomes were sequenced without selection of polyadenylated transcripts, whereas the

transcriptomes of *Cymbidium ensifolium* and *O. italica* represented polyA fractions. Thus, mobile elements, whose RNAs are not polyadenylated (Chang & Schulman, 2008) are expected to be overrepresented in *Epipogium*. In the genome of *Phalaenopsis equestris*, mobile elements are masked as repeats, producing a similar effect. Among

the underrepresented GO categories, almost all were related to photosynthesis and plastids. The least underrepresented GO terms were general and difficult to interpret (e.g., "Single-organism metabolic process" and "Membrane"). The most underrepresented GO terms are listed in Table 2; for a complete list, see Table S3.

Table 2. GO terms showing the greatest differences in the fraction of genes between mycoheterotrophic species and their photosynthetic relatives.

| Photosynthetic orchids compared with <i>Epipogium</i> | | | | |
|---|----------------------|---|--|--|
| GO term | Type of GO term | Median number of genes with this GO term in <i>Epipogium aphyllum</i> and <i>Epipogium roseum</i> | Median number of genes with this GO term in photosynthetic orchids | Median ratio of fractions between <i>Epipogium aphyllum</i> and <i>Epipogium roseum</i> versus photosynthetic orchids* |
| chlorophyll binding | molecular function | 0.5 | 34 | 0.011 |
| photosynthesis, light harvesting | biological process | 1 | 23 | 0.035 |
| photosystem I | cellular compartment | 2 | 38 | 0.041 |
| plastid thylakoid lumen | cellular compartment | 1.5 | 31 | 0.041 |
| plastid thylakoid membrane | cellular compartment | 55.5 | 201 | 0.20 |
| photosynthesis, light reaction | biological process | 50 | 174 | 0.22 |
| plastid thylakoid | cellular compartment | 89 | 275 | 0.25 |
| photosystem II assembly | biological process | 28 | 93 | 0.25 |
| photosynthesis | biological process | 125 | 294 | 0.35 |
| organelle subcompartment | cellular compartment | 203.5 | 403 | 0.38 |

| Photosynthetic Ericales compared with <i>Hypopitys monotropa</i> | | | | |
|---|----------------------|---|---|--|
| GO term | Type of GO term | Number of genes with this GO term in <i>H.monotropa</i> | Median number of genes with this GO term in photosynthetic Ericales | Median ratio of fractions between <i>H.monotropa</i> versus photosynthetic Ericales* |
| plastid thylakoid membrane | cellular compartment | 51 | 247 | 0.34 |
| photosynthesis, light reaction | biological process | 51 | 243 | 0.37 |
| plastid thylakoid | cellular compartment | 80 | 341 | 0.39 |
| photosynthesis | biological process | 97 | 347 | 0.47 |
| tetrapyrrole binding | molecular function | 136 | 395 | 0.57 |
| plastid | cellular compartment | 1 051 | 2 177 | 0.81 |
| oxidation-reduction process | biological process | 1 122 | 2 256 | 0.83 |
| single-organism process | biological process | 6 479 | 11 495 | 0.94 |

GO terms that are similar to each other, such as "plastid thylakoid lumen" and "chloroplast thylakoid lumen", are combined here. *- The fraction of genes with a specific GO term is the number of genes with that GO term in a species divided by the total number of genes with GO terms in that species. The median ratio of fractions is a measure of the difference in the numbers of genes with a specific GO term between two groups of species, calculated as a median value among all pairwise comparisons in which the first member in a pair comes from the first group, and the second member of the pair comes from the second group.

GO enrichment analyses between *H.monotropa* and three photosynthetic Ericales showed similar results; 17 GO terms were overrepresented, and 9 are underrepresented. The overrepresented terms were related to mobile elements, and the underrepresented terms were related to photosynthesis and plastids (Table 2; Table S3).

Notably, we did not observe changes in the list of GO

terms other than related to photosynthesis and plastids, despite dramatic differences in the morphology and physiology of the studied species. This finding indicates that these differences are controlled not at the level of the presence or absence of specific genes, but rather by the regulation of gene expression. To gain insight into this regulation, more detailed transcriptome data are required. Alternatively, these morphological and physiological

changes may have originated from the disruption or loss of a small number of genes which do not produce statistically significant results in GO enrichment analysis.

The statistical analysis of GO terms is quite rough method and may not reflect differences at the level of individual genes. Thus, we also searched for orthologs of genes that are known to participate in processes occurring in the plastids in the model plant *Arabidopsis thaliana*. As expected, genes related to photosynthesis have been lost from nuclear genomes of *Epipogium* and *H.monotropa* (Table 3). In particular, no nuclear-encoded components of the plastid electron transfer chain (*PSB* and *PSA* genes) were found, which is consistent with the absence of plastid-encoded *psa* and *psb* genes from the plastomes of all three species. Components of the light-harvesting antenna (*LHC* genes) were completely absent from *H.monotropa*; one such gene (*LHCA4*) was present in the *E.aphyllum* transcriptome, but its expression (measured in FPKM) was ~7-25-fold lower than in photosynthetic orchids. Nevertheless, this gene exhibited a full-length open reading frame that has evolved under negative selection (dN/dS 0.11). Since the product of this gene participates in chlorophyll binding, its retention may be in some way related to the retention of chlorophyll synthesis in *Epipogium*, as described below. Additionally, as shown in Table S7, a few transcripts of genes that encode proteins in the plastid electron transfer chain were found in the transcriptomes of *Epipogium aphyllum*, *E.roseum* and *H.monotropa*. All of the observed sequences were shorter than 50% of the length of their orthologs in photosynthetic plants and are likely to be pseudogenes. The apparent partial conservation of

enzymes in the Calvin cycle is due to that some of these enzymes are involved in metabolic processes that are not related to photosynthesis, e.g., glycolysis. Many genes that encode Calvin cycle enzymes belong to small gene families (e.g., *GAPDH*, *PGK*, *TKL*, *TPI*) in which different members encode proteins that exhibit the same enzymatic activity (isoenzymes) but show different cellular localizations and act in different pathways. For example, the *GAPDH* and *TPI* proteins function in both glycolysis (cytosolic isoenzymes) and the Calvin cycle (plastidic isoenzymes). We assume that most of the transcripts corresponded to cytosolic isoenzymes. Consistent with this assumption, 8 of 10 of the transcripts found in *Epipogium* and 7 of the 14 transcripts found in *H.monotropa* did not have plastid-targeting signals (Table S7), and those that have (*TKL*, *RPI*) exhibit additional functions in plastids that are not related to the Calvin cycle. Expression of genes encoding proteins that act exclusively in photosynthesis (e.g., *rbcS*, *SBPase*) was absent. The nuclear-encoded components of plastid ATP synthase have been completely lost, as well as the plastid-encoded components. The sigma subunits of plastid-encoded RNA polymerase (*PEP*) and *PEP*-associated proteins have also been lost. *PEP* is involved in the transcription of genes related to photosynthesis, unlike nuclear-encoded RNA polymerase (*NEP*), which mainly transcribes plastid genes that are unrelated to photosynthesis (Shiina *et al.*, 2005). The redox-sensitive components of the plastid translocon, which is the system responsible for the import of proteins from the cytoplasm to plastids, have also been lost, with the exception of *TIC32* of *Hypopitys*.

Table 3. Number of selected photosynthesis-related nuclear genes and nuclear genes encoding proteins involved in plastid functions other than photosynthesis in *Epipogium aphyllum* and photosynthetic orchids as well as *Hypopitys monotropa* and photosynthetic Ericales.

| Genes | Number in <i>E. aphyllum</i> | Median number in photosynthetic orchids | Number in <i>H. monotropa</i> | Median number in photosynthetic Ericales |
|--|------------------------------|---|-------------------------------|--|
| Photosynthesis-related | | | | |
| Components of photosystem I | 0 | 9 | 0 | 8 |
| Components of photosystem II | 0 | 9 | 0 | 8 |
| Components of electron transfer chain (others than PSI and PSII) | 0 | 5 | 0 | 6 |
| Light-harvesting complex | 1 | 10 | 0 | 8 |
| Calvin cycle | 10 | 20 | 14 | 21 |
| Sigma subunits of PEP and PEP-associated proteins | 0 | 14 | 1 | 14 |
| Plastid ATP synthase | 0 | 3 | 0 | 3 |
| Non-photosynthesis-related | | | | |
| Plastid ribosome | 33 | 34 | 31 | 31 |
| Clp subunits | 11 | 12 | 10 | 10 |
| ACC subunits | 5 | 6 | 4 | 7 |
| Plastid translocon | 15 | 20 | 11 | 18 |

In contrast to genes encoding proteins that are necessary for photosynthesis, genes that are responsible for other functions of plastids have been retained (Table 3). In particular, components of clp-protease and acetyl-CoA carboxylase whose counterparts (clpP and accD) are encoded in the plastome have been retained. Consistent with this finding, NEP, which transcribes genes not related to photosynthesis, has been retained in both species, as have most components of the plastid ribosome (but see the discussion below). We also found transcripts of most proteins responsible for the replication and repair of the plastome. However, this situation is more complex because many of these proteins are targeted not only to plastids but also to mitochondria and the nucleus, and information on these proteins is sometimes inconsistent between different experiments (Tanz *et al.*, 2013; Huang *et al.*, 2013).

As shown above, the genomes of *Epipogium* and *H.monotropa* encode a number of proteins that must be imported into plastids. Accordingly, the components of the plastid translocons for both the outer- and inner-envelope membranes must be retained. Recent studies in *A. thaliana* have shown that the plastid-encoded protein *ycf1* (TIC214) plays an essential role in plastid translocation and that it acts at the inner membrane (Kikuchi *et al.*, 2013). However, the *ycf1* gene is absent from both the *Epipogium* and *H.monotropa* plastomes. A transcript similar to *ycf1* was only found in the transcriptome assembly of *E.aphyllum*, in which it carried a signal for targeting of the protein to mitochondria. Homologs of TIC100 and TIC56, which encode proteins that interact with Ycf1 within the TIC complex, were also absent from all three species. It should be noted that Ycf1 and its interacting proteins are also absent from several photosynthetic species, including grasses and *Vaccinium macrocarpon*, which raises a question regarding the universality of the function of Ycf1 (de Vries *et al.*, 2015). The current model of TIC postulates the existence of two

complexes. The first, referred to as the “photosynthetic-type” complex, consists of Ycf1, TIC56, TIC100 and TIC20-I and is a major TIC complex that functions in most land plants to import proteins involved in photosynthesis. The second complex is a “non-photosynthetic-type” complex, which imports proteins that are not related to photosynthesis (Nakai, 2015b). It has been hypothesized that the switch from the major to the alternative system of protein import occurred in grasses (Nakai, 2015b; Nakai, 2015a) and that the major system then degraded. The pattern of gene loss and retention observed in *Epipogium* and *H.monotropa* suggests that this could also be the case in these species. The existence of two complexes, one of which mainly imports photosynthetic proteins, while the other is non-photosynthetic, has also been postulated for the outer translocon TOC (Nakai, 2015b; Nakai, 2015a). *Epipogium* and *H.monotropa* possess orthologs of TOC proteins, but not a complete set (Table 2), suggesting that only one of these complexes is retained in these plants.

Annotation of the transcriptomes using the KEGG database of biological pathways revealed reductions in the number of genes associated with several other pathways, some of which are common to *Epipogium* and *H.monotropa*, while others are lineage-specific (Figs. S2-S5). For example, a reduction in the number of proteins involved in light reception and circadian rhythms was observed; although this reduction was found in *Epipogium* as well as in *Hypopitys*, it was notably different in these holo-heterotrophic plants (Fig. S2). In particular, *Epipogium* appears to lack genes encoding the photoreceptors phytochrome B and cryptochrome, whereas *H.monotropa* lacked several proteins (LHY, ZTL and GI) that regulate the circadian clock. The absence of these proteins may be related to the distinctive lifestyles of these plants, which spend most of their life cycle completely underground and appear above-ground only for several weeks during blossoming (Bjorkman, 1960; Yagame *et al.*, 2007; Taylor & Roberts, 2011), and may therefore have different requirements for the regulation of

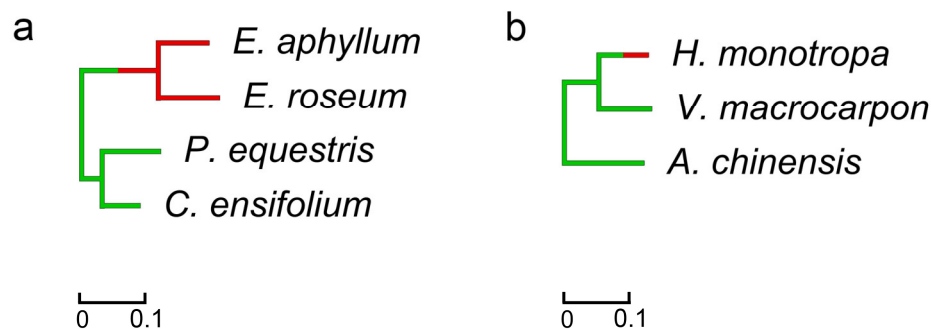
circadian rhythms than typical autotrophic plants. However, both *Epipogium* and *H.monotropa* have retained essential elements of circadian clock regulation (HY5, ELF3), photoreceptor phytochrome A and proteins interacting with them (PIF, COP1), indicating that the core proteins that regulate plant development under the influence of light have been conserved. In addition, we found a reduction in the number of genes associated with the carotenoid (Fig. S3) and phyloquinone (Fig. S4) biosynthetic pathways. As carotenoids are a component of light-harvesting antennae, the reduction in the number of genes related to carotenoid biosynthesis is presumably linked to the disappearance of photosynthesis. One interesting and unexplained observation was a reduction in the thiamine synthesis pathway that affected *Epipogium*, but not *H.monotropa* (Fig. S5). Another contrasting characteristic was the retention of the chlorophyll *a* synthesis pathway in *Epipogium*, whereas only the gene responsible for the first step in this pathway has been retained in *H.monotropa* (Fig. 1). In *H.monotropa*, the dN/dS for the only retained gene in this pathway was significantly higher than in autotrophic Ericales (p-value of 0.02 by the likelihood ratio test), showing a value of 0.27 in *H.monotropa* versus 0.11 in its autotrophic relatives. Among the 5 considered genes in *Epipogium*, only 1 showed increased dN/dS with a significant p-value

(< 0.05 by the likelihood ratio test with Bonferroni correction). The mean dN/dS for these 5 genes in *Epipogium* was only slightly higher than the mean in its photosynthetic relatives, at 0.11 versus 0.05, respectively (Table S7). However, all of the genes in this pathway were expressed at levels many times lower than in photosynthetic species (Table S7), suggesting that it could be not functional being the remnant of one that was active in a photosynthetic ancestor of *Epipogium*. Alternatively, chlorophyll *a* could indeed be synthesized in these species, where it may act in cellular processes that are unrelated to photosynthesis. Chlorophyll *a* has been found in many heterotrophic plants via chromatography (Cummings & Welschmeyer, 1998) (notably, among these species is *Monotropa uniflora*, a close relative of *H.monotropa*); transcriptome sequencing in parasitic *Orobanche aegyptiaca* has also demonstrated the presence of genes responsible for chlorophyll *a* synthesis, with no signs of relaxed selection (Wickett *et al.*, 2011). It has been shown that in *A. thaliana*, pheophorbide *a*, a product of chlorophyll *a* catabolism, causes cell death in a light-independent manner (Tanaka & Tanaka, 2006; Hirashima *et al.*, 2009). Conservation of this mechanism in *Arabidopsis* and non-photosynthetic plants could explain the conservation of the chlorophyll *a* synthesis pathway.

levels (approximately 20 and 2.5 times respectively) (Schelkunov *et al.*, 2015; Logacheva *et al.*, 2016). Characterization of transcriptome sequences allowed us to test whether this increase was confined to plastid genes. To calculate the mutation accumulation rate in the nuclear genomes, we used concatenated sequences of genes from orthogroups containing exactly one gene in each species. There were 4479 and 2547 of these orthogroups in

Orchidaceae and Ericales, respectively. The mutation accumulation rates in both *Epipogium* species were approximately two times higher than the mutation rates in their photosynthetic relatives. The rates of non-synonymous and synonymous mutation accumulation in *Epipogium* were increased proportionally (Figs. S6-S8). In contrast, the mutation accumulation rate in *H.monotropa* was not increased (Fig. 2).

Figure 2. Mutation accumulation rates in nuclear genes of *Epipogium* (a) and *Hypopitys monotropa* (b). The branch lengths denote the number of nucleotide substitutions per position. Branches corresponding to non-photosynthetic species are indicated in red, and those corresponding to photosynthetic species are indicated in green. Branches in which a transition from a photosynthetic to a non-photosynthetic lifestyle occurred are indicated half in green and half in red. *Orchis italica* and *Camellia sinensis*, which are employed as outgroups in (a) and (b), respectively, are not shown, since the mutation accumulation rate of an outgroup cannot be evaluated.



Composition of the plastid ribosome

The proteins of the plastid ribosome are encoded in both the plastid and nuclear genomes. For example, in *Arabidopsis thaliana*, 21 plastid ribosomal proteins are encoded by the plastome, and 36 are encoded by the nuclear genome. In all holo-heterotrophic plants with highly reduced plastomes, some ribosomal genes are

missing; *Pilostyles aethiopica*, in which only two ribosomal protein genes are retained, represents the most extreme known case (Bellot & Renner, 2015). This raises the question of how ribosomes are able to function in these species. There are three possibilities. First, some proteins of the plastid ribosome may simply be non-essential, and their loss may not severely affect ribosome function (Tiller

& Bock, 2014), although this does not explain such extreme cases of reduction. Alternatively, since the transfer and integration of plastid DNA into the nucleus exists in plant cells, functional copies of plastid genes can arise in the nuclear genome. There are examples in which the products of such nuclear copies are targeted to plastids and function as a part of the plastid ribosome, while the corresponding gene having been lost from the plastome (Ueda *et al.*, 2007; Jansen *et al.*, 2011; Park *et al.*, 2015). Additionally, components of mitochondrial ribosomes can be dually targeted to both plastids and mitochondria (Kubo & Arimura, 2010; Ueda *et al.*, 2008). In *E.aphyllum* (Schelkunov *et al.*, 2015), 7 ribosomal protein genes (*rpl20*, *rpl22*, *rpl23*, *rpl32*, *rpl33*, *rps15*, *rps16*) have been lost from the plastome; in *E.roseum* 6 of these 7 genes, but not *rpl20*, have been lost. Regarding *H.monotropa*, we previously considered *rps19* and *rpl22* to be pseudogenes (Logacheva *et al.*, 2016) due to the presence of a 111-bp insertion in the former and a nonsense mutation that shortens the length of the product by 17% in the latter. However, because these genes are transcribed and exhibit dN/dS values close to those of the species' photosynthetic relatives (Table S7), we assume that they are functional genes. Two genes, *rps15* and *rps16*, were completely absent from the plastome of *H.monotropa*. We do not observe any transcripts with high similarity to these plastid genes in the transcriptomes of *E.aphyllum*, *E.roseum* or *H.monotropa*. Thus, the loss of these genes from the plastome is unlikely to be compensated by transfer of plastid sequences to the nuclear genome.

Instead, we found that several transcripts that are not 1-1 orthologs of plastid-encoded ones but have more distant homology to them encode proteins that can be imported into plastids. This is the case for Rpl23 and Rps15 in *E.aphyllum* and Rps15 in *H.monotropa*. Additionally, for several proteins, the predictions made with TargetP and DualPred, two tools that we used for target prediction, were contradictory. Specifically, for a homolog of Rpl32 in *E.aphyllum* and homologs of Rpl23,

Rps15 and Rps16 in *E.roseum*, plastid targeting was predicted by only one of the two tools. In all cases except for Rpl23 of *E.aphyllum*, analysis of the transit peptides of the homologs in the photosynthetic relatives of these species suggested that these proteins were already targeted to plastids prior to the divergence of non-photosynthetic and photosynthetic species, which may have facilitated the loss of the respective plastid genes. Some of the aforementioned proteins are predicted to be targeted solely to plastids, and some are predicted to be dually targeted to plastids and mitochondria (for details see Table S8).

To determine whether the genes whose products may serve to replace the lost ribosomal proteins are encoded in mitochondrial or nuclear genomes, we used TBLASTN to align the proteins against contigs produced during the assembly of the plastomes of *E.aphyllum* and *H.monotropa* in our earlier studies (Schelkunov *et al.*, 2015, Logacheva *et al.*, 2016, respectively). We did not observe the sequences of those genes in the mitochondrial contigs and therefore conclude that they are located in the nuclear genomes.

Conclusions

In this study, we analysed and compared the transcriptomes of the mycoheterotrophic plants *Epipogium aphyllum*, *E.roseum* and *Hypopitys monotropa*. Despite the fact that *Epipogium* and *H.monotropa* are very distantly related, belonging to the Monocots and Dicots respectively, and that these species lost photosynthesis independently, we observed a remarkable level of parallelism involving the reduction and retention of similar functional groups of genes. Among the observed differences were a more profound reduction in the chlorophyll *a* synthesis pathway in *H.monotropa* and an increased rate of mutation accumulation in *Epipogium*. Overall, since there are several hundred holoheterotrophic species of flowering plants (Merckx *et al.*, 2009), with many cases of independent transitions to holo-

heterotrophic lifestyle, it is necessary to sequence and analyse more holo-heterotrophic species in addition to *Epipogium*, *H.monotropia* and *Orobancha aegyptiaca* to predict whether this parallelism is universal. Significant help in this determination may be provided by the "1000 Plants Project" (Matasci *et al.*, 2014), in which the sequencing of many holo- and hemi-heterotrophic species

is being performed. The question that remains unanswered in this study is the mode of gene loss – did it occur through deletions of large regions carrying photosynthesis-related genes, or through the accumulation of mutations in the protein-coding and regulatory elements of these genes? We expect that characterization of the nuclear genomes of non-photosynthetic plants will fill this gap.

Acknowledgements

The authors would like to thank Viktoria Shtratnikova for helpful discussion.

Authors' contributions

MIS performed the computational analysis and drafted the manuscript. AAP participated in sequencing. MDL sequenced the samples and participated in discussion and manuscript writing.

References:

Abascal F, Zardoya R, Telford MJ. 2010. TranslatorX: multiple alignment of nucleotide sequences guided by amino acid translations. *Nucleic Acids Research* **38**: W7–W13.

Barrett CF, Freudenstein JV, Li J, Mayfield-Jones DR, Perez L, Pires JC, Santos C. 2014. Investigating the path of plastid genome degradation in an early-transitional clade of heterotrophic orchids, and implications for heterotrophic angiosperms. *Molecular Biology and Evolution* **31**: 3095–3112.

Bellot S, Cusimano N, Luo S, Sun G, Zarre S, Gröger A, Tensch E, Renner SS. 2016. Assembled Plastid and Mitochondrial Genomes, as well as Nuclear Genes, Place the Parasite Family Cynomoriaceae in the Saxifragales. *Genome Biology and Evolution* **8**: 2214–2230.

Bellot S, Renner SS. 2015. The plastomes of two species in the endoparasite genus *Pilosyles* (Apodanthaceae) each retain just five or six possibly functional genes. *Genome Biology and Evolution*: evv251.

Berardini TZ, Reiser L, Li D, Mezheritsky Y, Muller R, Strait E, Huala E. 2015. The arabidopsis information resource: Making and mining the “gold standard” annotated reference plant genome: Tair:

Making and Mining the “Gold Standard” Plant Genome. *genesis* **53**: 474–485.

Bjorkman E. 1960. Monotropa Hypopitys L. - an Epiparasite on Tree Roots. *Physiologia Plantarum* **13**: 308–327.

Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**: 2114–2120.

Bromham L, Cowman PF, Lanfear R. 2013. Parasitic plants have increased rates of molecular evolution across all three genomes. *BMC evolutionary biology* **13**: 126.

Bromham L, Hua X, Lanfear R, Cowman PF. 2015. Exploring the Relationships between Mutation Rates, Life History, Genome Size, Environment, and Species Richness in Flowering Plants. *The American Naturalist* **185**: 507–524.

Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+: architecture and applications. *BMC Bioinformatics* **10**: 421.

Chang W, Schulman AH. 2008. BARE retrotransposons produce multiple groups of rarely polyadenylated transcripts from two differentially regulated promoters. *The Plant Journal* **56**: 40–50.

Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, Robles M. 2005. Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics (Oxford, England)* **21**: 3674–3676.

Cummings MP, Welschmeyer NA. 1998. Pigment composition of putatively achlorophyllous angiosperms. *Plant Systematics and Evolution* **210**: 105–111.

Eddy SR. 1998. Profile hidden Markov models. *Bioinformatics (Oxford, England)* **14**: 755–763.

Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research* **32**: 1792–1797.

Emanuelsson O, Nielsen H, Brunak S, von Heijne G. 2000. Predicting Subcellular Localization of Proteins Based on their N-terminal Amino Acid Sequence. *Journal of Molecular Biology* **300**: 1005–1016.

Fan W, Zhu A, Kozaczek M, Shah N, Pabón-Mora N, González F, Mower JP. 2016. Limited

mitogenomic degradation in response to a parasitic lifestyle in Orobanchaceae. *Scientific Reports* **6**: 36285.

Finn RD, Clements J, Arndt W, Miller BL, Wheeler TJ, Schreiber F, Bateman A, Eddy SR. 2015. HMMER web server: 2015 update. *Nucleic Acids Research* **43**: W30-38.

Finn RD, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, Potter SC, Punta M, Qureshi M, Sangrador-Vegas A, et al. 2016. The Pfam protein families database: towards a more sustainable future. *Nucleic Acids Research* **44**: D279-285.

Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q, et al. 2011. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature Biotechnology* **29**: 644–652.

Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, 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.

Hedges SB, Marin J, Suleski M, Paymer M, Kumar S. 2015. Tree of Life Reveals Clock-Like Speciation and Diversification. *Molecular Biology and Evolution* **32**: 835–845.

Hirashima M, Tanaka R, Tanaka A. 2009. Light-Independent Cell Death Induced by Accumulation of Pheophorbide a in *Arabidopsis thaliana*. *Plant and Cell Physiology* **50**: 719–729.

Huang M, Friso G, Nishimura K, Qu X, Olinares PDB, Majeran W, Sun Q, van Wijk KJ. 2013. Construction of Plastid Reference Proteomes for Maize and *Arabidopsis* and Evaluation of Their Orthologous Relationships; The Concept of Orthoproteomics. *Journal of Proteome Research* **12**: 491–504.

Jansen RK, Sasaki C, Lee S-B, Hansen AK, Daniell H. 2011. Complete plastid genome sequences of three rosids (*Castanea*, *Prunus*, *Theobroma*): evidence for at least two independent transfers of rpl22 to the nucleus. *Molecular biology and evolution* **28**: 835–847.

Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M. 2016a. KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Research* **44**: D457-462.

Kanehisa M, Sato Y, Morishima K. 2016b. BlastKOALA and GhostKOALA: KEGG Tools for

Functional Characterization of Genome and Metagenome Sequences. *Journal of Molecular Biology* **428**: 726–731.

Kikuchi S, Bedard J, Hirano M, Hirabayashi Y, Oishi M, Imai M, Takase M, Ide T, Nakai M. 2013. Uncovering the protein translocon at the chloroplast inner envelope membrane. *Science* **339**: 571–574.

Kim G, LeBlanc ML, Wafula EK, dePamphilis CW, Westwood JH. 2014. Genomic-scale exchange of mRNA between a parasitic plant and its hosts. *Science* **345**: 808–811.

Kubo N, Arimura S -i. 2010. Discovery of the rpl10 Gene in Diverse Plant Mitochondrial Genomes and Its Probable Replacement by the Nuclear Gene for Chloroplast RPL10 in Two Lineages of Angiosperms. *DNA Research* **17**: 1–9.

Li B, Dewey CN. 2011. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics* **12**: 323.

Li L, Stoeckert CJ, Roos DS. 2003. OrthoMCL: identification of ortholog groups for eukaryotic genomes. *Genome Research* **13**: 2178–2189.

Li X, Zhang T-C, Qiao Q, Ren Z, Zhao J, Yonezawa T, Hasegawa M, Crabbe MJC, Li J, Zhong Y. 2013. Complete chloroplast genome sequence of holoparasite *Cistanche deserticola* (Orobanchaceae) reveals gene loss and horizontal gene transfer from its host *Haloxylon ammodendron* (Chenopodiaceae). *PLoS one* **8**: e58747.

Lim GS, Barrett CF, Pang C-C, Davis JI. 2016. Drastic reduction of plastome size in the mycoheterotrophic *Thismia tentaculata* relative to that of its autotrophic relative *Tacca chantrieri*. *American Journal of Botany* **103**: 1129–1137.

Logacheva MD, Schelkunov MI, Nuraliev MS, Samigullin TH, Penin AA. 2014. The Plastid Genome of Mycoheterotrophic Monocot *Petrosavia stellaris* Exhibits Both Gene Losses and Multiple Rearrangements. *Genome Biology and Evolution* **6**: 238–246.

Logacheva MD, Schelkunov MI, Shtratnikova VY, Matveeva MV, Penin AA. 2016. Comparative analysis of plastid genomes of non-photosynthetic Ericaceae and their photosynthetic relatives. *Scientific Reports* **6**: 30042.

Matasci N, Hung L-H, Yan Z, Carpenter EJ, Wickett NJ, Mirarab S, Nguyen N, Warnow T, Ayyampalayam S, Barker M, et al. 2014. Data access for the 1,000 Plants (1KP) project. *GigaScience* **3**: 17.

Merckx V, Bidartondo MI, Hynson NA. 2009. Myco-heterotrophy: when fungi host plants. *Annals of Botany* **104**: 1255–1261.

Mitschke J, Fuss J, Blum T, Höglund A, Reski R, Kohlbacher O, Rensing SA. 2009. Prediction of dual protein targeting to plant organelles. *New Phytologist* **183**: 224–236.

Molina J, Hazzouri KM, Nickrent D, Geisler M, Meyer RS, Pentony MM, Flowers JM, Pelsler P, Barcelona J, Inovejas SA, et al. 2014. Possible loss of the chloroplast genome in the parasitic flowering plant *Rafflesia lagascae* (*Rafflesiaceae*). *Molecular Biology and Evolution* **31**: 793–803.

Nakai M. 2015a. YCF1: a green TIC: response to the de Vries et al. Commentary. *The Plant Cell* **27**: 1834–1838.

Nakai M. 2015b. The TIC complex uncovered: The alternative view on the molecular mechanism of protein translocation across the inner envelope membrane of chloroplasts. *Biochimica et Biophysica Acta (BBA) - Bioenergetics* **1847**: 957–967.

Park S, Jansen RK, Park S. 2015. Complete plastome sequence of *Thalictrum coreanum* (Ranunculaceae) and transfer of the rpl32 gene to the nucleus in the ancestor of the subfamily Thalicthroideae. *BMC Plant Biology* **15**: 40.

Parra G, Bradnam K, Korf I. 2007. CEGMA: a pipeline to accurately annotate core genes in eukaryotic genomes. *Bioinformatics* **23**: 1061–1067.

Rasmussen HN. 1995. *Terrestrial orchids from seed to mycotrophic plant*. Cambridge ; New York: Cambridge University Press.

Ravin NV, Gruzdev EV, Beletsky AV, Mazur AM, Prokhortchouk EB, Filyushin MA, Kochieva EZ, Kadnikov VV, Mardanov AV, Skryabin KG. 2016. The loss of photosynthetic pathways in the plastid and nuclear genomes of the non-photosynthetic mycoheterotrophic eudicot *Monotropa hypopitys*. *BMC Plant Biology* **16**: 153–161.

Saravanan V, Velan Lakshmi P. 2015. Dualpred: A Webserver for Predicting Plant Proteins Dual-

Targeted to Chloroplast and Mitochondria Using Split Protein-Relatedness-Measure Feature. *Current Bioinformatics* **10**: 323–331.

Schelkunov MI, Shtratnikova VY, Nuraliev MS, Selosse M-A, Penin AA, Logacheva MD. 2015. Exploring the limits for reduction of plastid genomes: a case study of the mycoheterotrophic orchids *Epipogium aphyllum* and *Epipogium roseum*. *Genome Biology and Evolution* **7**: 1179–1191.

Shiina T, Tsunoyama Y, Nakahira Y, Khan MS. 2005. Plastid RNA polymerases, promoters, and transcription regulators in higher plants. *International Review of Cytology* **244**: 1–68.

Smith DR, Lee RW. 2014. A plastid without a genome: evidence from the nonphotosynthetic green algal genus *Polytomella*. *Plant physiology* **164**: 1812–1819.

Stöver BC, Müller KF. 2010. TreeGraph 2: Combining and visualizing evidence from different phylogenetic analyses. *BMC Bioinformatics* **11**: 7.

Tanaka A, Tanaka R. 2006. Chlorophyll metabolism. *Current Opinion in Plant Biology* **9**: 248–255.

Tanz SK, Castleden I, Hooper CM, Vacher M, Small I, Millar HA. 2013. SUBA3: a database for integrating experimentation and prediction to define the SUBcellular location of proteins in *Arabidopsis*. *Nucleic Acids Research* **41**: D1185-1191.

Taylor L, Roberts DL. 2011. Biological Flora of the British Isles: *Epipogium aphyllum* Sw.: *Epipogium aphyllum* Sw. *Journal of Ecology* **99**: 878–890.

THE ANGIOSPERM PHYLOGENY GROUP. 2009. An update of the Angiosperm Phylogeny Group classification for the orders and families of flowering plants: APG III: APG III. *Botanical Journal of the Linnean Society* **161**: 105–121.

Tiller N, Bock R. 2014. The translational apparatus of plastids and its role in plant development. *Molecular Plant* **7**: 1105–1120.

Ueda M, Fujimoto M, Arimura S, Murata J, Tsutsumi N, Kadowaki K. 2007. Loss of the rpl32 gene from the chloroplast genome and subsequent acquisition of a preexisting transit peptide within the nuclear gene in *Populus*. *Gene* **402**: 51–56.

Ueda M, Nishikawa T, Fujimoto M, Takanashi H, Arimura S -i., Tsutsumi N, Kadowaki K -i.

2008. Substitution of the gene for chloroplast RPS16 was assisted by generation of a dual targeting signal. *Molecular Biology and Evolution* **25**: 1566–1575.

de Vries J, Sousa FL, Bölter B, Soll J, Gould SB. 2015. YCF1: a green TIC? *The Plant Cell* **27**: 1827–1833.

Westwood JH, Yoder JI, Timko MP, dePamphilis CW. 2010. The evolution of parasitism in plants. *Trends in Plant Science* **15**: 227–235.

Wicke S, Muller KF, de Pamphilis CW, Quandt D, Wickett NJ, Zhang Y, Renner SS, Schneeweiss GM. 2013. Mechanisms of functional and physical genome reduction in photosynthetic and nonphotosynthetic parasitic plants of the Broomrape Family. *The Plant Cell* **25**: 3711–3725.

Wickett NJ, Honaas LA, Wafula EK, Das M, Huang K, Wu B, Landherr L, Timko MP, Yoder J, Westwood JH, et al. 2011. Transcriptomes of the Parasitic Plant Family Orobanchaceae Reveal Surprising Conservation of Chlorophyll Synthesis. *Current Biology* **21**: 2098–2104.

Yagame T, Yamato M, Mii M, Suzuki A, Iwase K. 2007. Developmental processes of achlorophyllous orchid, *Epipogium roseum*: from seed germination to flowering under symbiotic cultivation with mycorrhizal fungus. *Journal of Plant Research* **120**: 229–236.

Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular Biology and Evolution* **24**: 1586–1591.

Yang Z, Wafula EK, Honaas LA, Zhang H, Das M, Fernandez-Aparicio M, Huang K, Bandaranayake PCG, Wu B, Der JP, et al. 2015. Comparative Transcriptome Analyses Reveal Core Parasitism Genes and Suggest Gene Duplication and Repurposing as Sources of Structural Novelty. *Molecular Biology and Evolution* **32**: 767–790.

Yang Z, Zhang Y, Wafula EK, Honaas LA, Ralph PE, Jones S, Clarke CR, Liu S, Su C, Zhang H, et al. 2016. Horizontal gene transfer is more frequent with increased heterotrophy and contributes to parasite adaptation. *Proceedings of the National Academy of Sciences* **113**: E7010–E7019.

Yekutieli D, Benjamini Y. 2001. The control of the false discovery rate in multiple testing under dependency. *The Annals of Statistics* **29**: 1165–1188.

Online supporting material

Fig. S1 Statistics regarding contamination in the studied transcriptomes. A total of 10,000 random transcripts (prior to the removal of low-coverage transcripts and searching for ORFs, but after the removal of minor isoforms) were taken from each assembly, and BLASTX alignment to NCBI NR (maximum allowed e-value of 10^{-5} , word size of 3 amino acids, low-complexity sequence filter switched off) was performed. The transcripts were classified according to their best matches. In the distribution plots, the black lines denote median values, and the boxes denote interquartile ranges.

Fig. S2 Diagram of circadian rhythms regulation.

Fig. S3 Diagram of carotenoid biosynthesis.

Fig. S4 Diagram of ubiquinone and other terpenoid-quinone biosynthesis.

Fig. S5 Diagram of thiamine metabolism.

Fig. S6 Trees of the studied species with branch lengths representing dS (rate of synonymous substitutions).

Fig. S7 Trees of the studied species with branch lengths representing dN (rate of non-synonymous substitutions).

Fig. S8 Trees of the studied species with branch lengths representing dN/dS.

Table S1 Information on transcriptome data.

Table S2 Sources of genome sequences.

Table S3 Complete list of GO terms for which the fractions of genes differed significantly among *Epipogium*, *Hypopitys monotropa* and their photosynthetic relatives * - "u" indicates underrepresentation; "o" indicates overrepresentation. Genes with a GO term are considered underrepresented in holo-heterotrophic species if their proportions relative to the total numbers of

genes with GO terms in those species are significantly lower than in photosynthetic species.

Overrepresentation corresponds to the opposite situation.

Table S4 Detailed statistics of the transcriptome assemblies. * - In a dataset deposited by the authors of the *Vaccinium macrocarpon* genome assembly, several isoforms are provided for some genes. In these cases, we referred to the longest isoform as the “major” isoform, since information on the relative expression of isoforms was not supplied.

Table S5 Analysis of the completeness of the assemblies.

Table S6 Proportions of genes whose products are targeted to various organelles relative to the total numbers of genes in the species according to TargetP analysis. Only genes with complete 5' ends and at least one assigned GO term are considered. The quality of the different assemblies differs; thus, the number of genes with completely assembled 5' ends also differs. Therefore, the numbers of proteins with transit peptides cannot be directly compared, and only proportions are provided in the table.

Table S7 Presence of genes of interest in the studied species and selective pressures acting on them.

Table S8 Results of a search for possible substitutes for ribosomal proteins whose genes have been lost from the plastomes of *Epipogium* and *Hypopitys monotropa*.