# *Rhopalocnemis phalloides* has one of the most reduced and mutated plastid genomes known

**Mikhail I Schelkunov** [Corresp., 1, 2] , **Maxim S Nuraliev** [3, 4] , **Maria D Logacheva** [1, 5]

[1] Skolkovo Institute of Science and Technology, Moscow, Russia

[2] Institute for Information Transmission Problems, Moscow, Russia

[3] Faculty of Biology, Moscow State University, Moscow, Russia

[4] Joint Russian–Vietnamese Tropical Scientific and Technological Center, Cau Giay, Hanoi, Vietnam

[5] A.N. Belozersky Research Institute of Physico-Chemical Biology, Moscow State University, Moscow, Russia

Corresponding Author: Mikhail I Schelkunov
Email address: shelkmike@gmail.com

Although most plant species are photosynthetic, several hundred species have lost the ability to photosynthesize and instead obtain nutrients via various types of heterotrophic feeding. Their plastid genomes markedly differ from the plastid genomes of photosynthetic plants. In this work, we describe the sequenced plastid genome of the heterotrophic plant *Rhopalocnemis phalloides*, which belongs to the family Balanophoraceae and feeds by parasitising other plants. The genome is highly reduced (18,622 base pairs versus approximately 150 kilobase pairs in autotrophic plants) and possesses an extraordinarily high AT content, 86.8%, which is inferior only to AT contents of plastid genomes of *Balanophora*, a genus from the same family. The gene content of this genome is quite typical of heterotrophic plants, with all of the genes related to photosynthesis having been lost. The remaining genes are notably distorted by a high mutation rate and the aforementioned AT content. The high AT content has led to sequence convergence between some of the remaining genes and their homologues from AT-rich plastid genomes of protists. Overall, the plastid genome of *R. phalloides* is one of the most unusual plastid genomes known.

1  *Rhopalocnemis phalloides* has one of the most reduced
2  and mutated plastid genomes known
3

4  Mikhail I. Schelkunov[1, 2], Maxim S. Nuraliev[3, 4], Maria D. Logacheva[1, 5]

5  [1] Skolkovo Institute of Science and Technology, Moscow, Russia

6  [2] Institute for Information Transmission Problems, Moscow, Russia

7  [3] Faculty of Biology, Moscow State University, Moscow, Russia

8  [4] Joint Russian–Vietnamese Tropical Scientific and Technological Center, Cau Giay, Hanoi,
9  Vietnam

10  [5] A.N. Belozersky Research Institute of Physico-Chemical Biology, Moscow State University,
11  Moscow, Russia

12

13  Corresponding Author:

14  Mikhail Schelkunov

15  Email address: shelkmike@gmail.com

16 **Abstract**

17 Although most plant species are photosynthetic, several hundred species have lost the ability to
18 photosynthesize and instead obtain nutrients via various types of heterotrophic feeding. Their
19 plastid genomes markedly differ from the plastid genomes of photosynthetic plants. In this work,
20 we describe the sequenced plastid genome of the heterotrophic plant *Rhopalocnemis phalloides*,
21 which belongs to the family Balanophoraceae and feeds by parasitising other plants. The genome
22 is highly reduced (18,622 base pairs versus approximately 150 kilobase pairs in autotrophic
23 plants) and possesses an extraordinarily high AT content, 86.8%, which is inferior only to AT
24 contents of plastid genomes of *Balanophora*, a genus from the same family. The gene content of
25 this genome is quite typical of heterotrophic plants, with all of the genes related to
26 photosynthesis having been lost. The remaining genes are notably distorted by a high mutation
27 rate and the aforementioned AT content. The high AT content has led to sequence convergence
28 between some of the remaining genes and their homologues from AT-rich plastid genomes of
29 protists. Overall, the plastid genome of *R. phalloides* is one of the most unusual plastid genomes
30 known.

31

32 **Introduction**

33 Though plants are generally considered photosynthetic organisms, there are several hundred
34 plant species that have lost the ability to photosynthesize during the course of evolution
35 (Westwood et al., 2010; Merckx et al., 2013). They feed either by parasitising other plants or by
36 obtaining nutrients from fungi.. In addition to the completely heterotrophic plants, there are also
37 plants that combine the ability to photosynthesize with the heterotrophic lifestyle. They are
38 termed partial heterotrophs (or hemi-heterotrophs, or mixotrophs) in contrast to the former,
39 which are termed complete heterotrophs (or holo-heterotrophs).

40 The completely heterotrophic plants show a high degree of similarity, though there were several
41 dozen cases of independent transition to complete heterotrophy. For example, these plants all
42 either lack leaves or have very reduced leaves. These plants are non-green because of the
43 absence (or at least highly reduced amounts (Cummings & Welschmeyer, 1998)) of chlorophyll.
44 Additionally, a common feature of many completely heterotrophic angiosperms is that they
45 spend most of their lifetimes underground, since without the need to photosynthesize their only
46 reason to appear aboveground is for flowering and seed dispersal.

47 Genomic studies of heterotrophic plants are mostly focused on plastid genomes, since 1) most of
48 the plastid genes are related to photosynthesis, and thus changes in the plastid genomes are
49 expected to be more prominent compared to mitochondrial and nuclear genomes, and 2) plastid
50 genomes are smaller than nuclear and mitochondrial ones and usually have higher copy numbers
51 and are thus easier to sequence (Daniell et al., 2016; Gualberto & Newton, 2017; Sakamoto &
52 Takami, 2018). The main feature of the plastid genomes of complete heterotrophs is the loss of
53 genes responsible for photosynthesis and respective shortening of the genomes, from
54 approximately 150 kbp (typical of autotrophic plants) to, in the most extreme known case, 12
55 kbp (Bellot & Renner, 2015; Graham, Lam & Merckx, 2017; Wicke & Naumann, 2018). The
56 remaining genes are the ones with functions not related to photosynthesis. Usually they are *accD*

57  (a gene whose product participates in fatty acid synthesis—one of the plastid functions besides
58  photosynthesis), *clpP* (encodes a component of a complex responsible for degradation of waste
59  proteins in plastids), *ycf1* (thought to encode a component of the translocon—the complex which
60  imports proteins from cytoplasm into plastids), *ycf2* (a conserved gene present in almost all
61  plants, but with unknown function) and various genes required for translation of the
62  aforementioned ones, namely genes that code for protein and RNA components of the plastid
63  ribosome and for tRNAs. One of the tRNA-coding genes, *trnE-UUC*, also has an additional
64  function, with its product participating in haem synthesis (Kumar et al., 1996).

65  In addition to the expected shortening of the genome, there are some peculiar and still
66  unexplained features in the plastid genomes of heterotrophic plants, namely their increased
67  mutation accumulation rate (Bromham, Cowman & Lanfear, 2013; Wicke & Naumann, 2018)
68  and increased AT content (Wicke & Naumann, 2018). In the most extreme cases, plastid
69  genomes of heterotrophic plants may accumulate mutations approximately 100 times faster than
70  their closest autotrophic relatives (Bellot & Renner, 2015). The most obvious explanation, the
71  relaxation of selection, is refuted by the fact that dN/dS (a common measure of selective
72  pressure) is usually not increased in the plastid genes of heterotrophic plants, except for
73  photosynthesis-related genes during their pseudogenization, but the mutation accumulation rate
74  is high even after the loss of all such genes (Logacheva, Schelkunov & Penin, 2011; Barrett et
75  al., 2014; Schelkunov et al., 2015; Lam, Soto Gomez & Graham, 2015; Wicke et al., 2016;
76  Naumann et al., 2016). AT content is increased from approximately 65% in autotrophic species
77  (Smith, 2012) to 88.4% in the most prominent case among heterotrophic species (Su et al.,
78  2019), also because of an unknown reason.

79  Genes not related to photosynthesis, such as *accD* and *infA*, are sometimes transferred to the
80  nuclear genome (Millen et al., 2001; Rousseau-Gueutin et al., 2013; Liu et al., 2016). Therefore,
81  when all genes with functions not related to translation are transferred to the nuclear genome,
82  there will be no reasons to keep the translation apparatus in plastids, and the genes responsible
83  for translation will also be lost. Thus, the plastid genome is potentially able to disappear entirely.
84  Indeed, two putative cases of the complete plastid genome loss are known: one in algae of the
85  genus *Polytomella* (Smith & Lee, 2014) and the other one in the parasitic plant *Rafflesia*
86  *lagascae* (Molina et al., 2014); the second case is disputable (Krause, 2015).

87  The initial aim of the present study was to prove that the completely heterotrophic plant
88  *Rhopalocnemis phalloides* had also lost its plastid genome completely. *Rhopalocnemis*
89  *phalloides* is a parasitic plant from the family Balanophoraceae (order Santalales) which occurs
90  in Asia and feeds by obtaining nutrients from roots of various plants. Initially we sequenced
91  approximately 10 million pairs of reads on the HiSeq 2000 platform and observed no contigs
92  with similarity to typical plastid genes, while there were obvious mitochondrial contigs. Based
93  on our experience in studying plastid genomes of heterotrophic plants, mitochondrial contigs
94  usually have lower sequencing coverage than plastid contigs; thus, the plastid genome is always
95  easier to assemble. This led us to suppose that the plastid genome in *R. phalloides* may have
96  been completely lost. To verify this, we sequenced approximately 200 million pairs of additional
97  reads. What we found is that the plastid genome *is* in fact present, but its tremendous AT content

98   (86.8%) hampered PCR, which is one of the usual steps in library preparation of Illumina, and
99   thus the sequencing coverage of the genome was much lower than one might have expected. This
100  article is dedicated to the analysis of this plastid genome.

101

## Materials & Methods

103

## Sample collection and sequencing

105

106  The specimen of *R. phalloides* was collected during an expedition of the Russian-Vietnamese
107  Tropical Centre in Kon Tum province, Vietnam, in May 2015 (voucher information: Southern
108  Vietnam, Kon Tum prov., Kon Plong distr., Thach Nham protected forest, 17 km N of Mang Den
109  town, in open forest, N 14° 45' 15" E 108° 17' 40", elev. 1400 m, Nuraliev M.S., Kuznetsov
110  A.N., Kuznetsova S.P., No. 1387, 18.04.2015). The studied material was preserved in silica gel
111  and in RNAlater. The voucher is deposited at the Moscow University Herbarium (MW) (Seregin,
112  2018) with the barcode MW0755444.

113  DNA was extracted from an inflorescence using a CTAB-based method (Doyle, 1987), and the
114  DNA library was prepared using the NEBNext DNA Ultra II kit (New England Biolabs).
115  Sequencing was performed with a NextSeq 500 sequencing machine (Illumina) in the paired end
116  mode, producing 387,351,294 reads (193,675,647 read pairs), each 150 bp long.

117  RNA was extracted from an inflorescence using the RNeasy Mini kit (Qiagen). Plastid
118  transcripts are usually not polyadenylated, so the method of poly(A) RNA selection was not
119  applicable in our study. Instead, we used a protocol based on depletion of ribosomal RNA with
120  the Plant Leaf Ribo Zero kit (Illumina). The RNA-seq library was prepared using the NEBNext
121  RNA Ultra II kit (New England Biolabs) and sequenced on a HiSeq 2500 sequencing machine
122  (Illumina) with TruSeq reagents v.4 in the paired end mode, producing 54,794,466 reads
123  (27,397,233 read pairs), 125 bp each.

124

## Genome assembly and annotation

126

127  Both DNA-seq and RNA-seq reads were trimmed by Trimmomatic 0.32 (Bolger, Lohse &
128  Usadel, 2014) in the palindromic mode, removing bases with quality less than 3 from the 3' ends
129  of reads, and fragments starting from 4-base-long windows with average quality less than 15
130  (SLIDINGWINDOW:4:15). Reads that, after trimming, had average quality less than 20 or
131  length shorter than 30 bases were removed.

132  The assembly was performed from DNA-seq reads by two tools. First, it was made by CLC
133  Assembly Cell 4.2 (https://www.qiagenbioinformatics.com/products/clc-assembly-cell/) with the
134  default parameters. Second, it was made by Spades 3.9.0 (Bankevich et al., 2012). Because the

135   performance of Spades is slow when running on large number of reads, prior to starting its
136   assembly we removed from reads k-mers with coverage less than 50× by Kmernator 1.2.0
137   (https://github.com/JGIBioinformatics/Kmernator). This allowed us to eliminate most reads
138   belonging to the nuclear genome (and, potentially, some reads belonging to low-covered plastid
139   regions), thus highly reducing the number of reads. The Spades assembly was run on this
140   reduced read set, with the "--only-assembler" and "--careful" options. To determine the read
141   coverage of contigs in these two assemblies, we aligned to them reads by CLC Assembly Cell
142   4.2, requiring at least 80% of the length of each read to align with a sequence similarity of at
143   least 98%.

144   To find contigs potentially belonging to plastid and mitochondrial genomes, we aligned by
145   BLASTN and TBLASTN from BLAST 2.3.0+ suit (Camacho et al., 2009) proteins and non-
146   coding RNA (ncRNA) genes from reference species. As the references, we used sequences from
147   the plastid genomes of *Balanophora reflexa* (NCBI accession KX784266), *Balanophora*
148   *laxiflora* (NCBI accession KX784265), *Viscum album* (NCBI accession NC_028012), *Osyris*
149   *alba* (NCBI accession NC_027960), *Arabidopsis thaliana* (NCBI accession NC_000932),
150   *Nicotiana tabacum* (NCBI accession NC_018041) and mitochondrial genomes of *Viscum album*
151   (NCBI accession NC_029039), *Citrullus lanatus* (NCBI accession GQ856147), *Mimulus*
152   *guttatus* (NCBI accession NC_018041). *Balanophora*, *Viscum album* and *O. alba* were used
153   because they, like *R. phalloides*, belong to Santalales. Other species were chosen because they
154   belong to other orders of eudicots. Alignment was performed with the maximum e-value of $10^{-3}$
155   and low complexity filter switched off. The word size was 7 for BLASTN and 3 for TBLASTN.
156   Here and later, the local BLAST was used with the parameter "max_target_seqs" set to $10^9$ to
157   avoid the problem discussed by Shah et al. (2018), who state that BLAST results may be
158   improper when this parameter is set to a small value.

159   Five contigs containing plastid genes were found in the CLC assembly and three contigs in the
160   Spades assembly. After aligning contigs of these two assemblies to each other (BLASTN,
161   maximum e-value $10^{-3}$, word size 7, low complexity filter switched off), it appeared that the
162   places at which the CLC contigs were broken by gaps corresponded to continuous places in the
163   Spades contigs and, vice versa, gaps in the Spades contigs corresponded to continuous places in
164   the CLC contigs. This allowed us, by joining the contigs of these two assemblies, to create a
165   circular sequence corresponding to the plastid genome. To check the assembly, we mapped reads
166   (in CLC Assembly Cell 4.2, requiring at least 80% of the length of each read to align with a
167   sequence similarity of at least 98%) to the resultant sequence and verified (by eye, in CLC
168   Genomics Workbench 7.5.1, https://www.qiagenbioinformatics.com/products/clc-genomics-
169   workbench/) that there were no places uncovered by reads and no places where the insert size
170   abruptly decreased or increased. Such places of abrupt increase or decrease of the insert size may
171   indicate regions with assembly errors, consisting of sequence insertions or deletions,
172   respectively. As read mapping is complicated on the edges of a sequence, we also performed
173   such analysis on a reoriented version of the plastid genome, in which the sequence was broken in
174   the middle and the ends were joined. These analyses indicated that the assembly contained no
175   errors.

176  To find genes in the plastid genome, we used a complex strategy, because highly mutated genes
177  may be hard to notice. We used the following methods:

178  1.  The alignment of reference protein-coding and ncRNA-coding genes by BLASTN and
179      TBLASTN, as described above.
180  2.  Open reading frames were scanned by InterProScan 5.11 (Jones et al., 2014) using the
181      InterPro 51.0 (Finn et al., 2017) database with the default parameters. "Open reading
182      frames" here were any sequences at least 20 codons long uninterrupted by stop codons.
183      Not requiring an ORF to begin from a start-codon allowed for the detection of exons in
184      multi-exonic genes.
185  3.  The genome was scanned by Infernal 1.1.2 (Nawrocki & Eddy, 2013) with RNA models
186      from Rfam 12.2 database (Nawrocki et al., 2015) to predict ncRNA-coding genes. The
187      maximum allowed e-value was set to $10^{-3}$.
188  4.  To predict rRNA-coding genes, RNAmmer 1.2 server (Lagesen et al., 2007) was used in
189      bacterial mode and eukaryotic mode.
190  5.  The genome was scanned by tRNAscan-SE 1.23 (Lowe & Eddy, 1997) with the default
191      parameters, in the organellar (models trained on plastid and mitochondrial tRNAs) and
192      also in the general (models trained on tRNAs from all three genomes) mode, to predict
193      tRNA-coding genes.
194  6.  The genome was annotated by DOGMA (Wyman, Jansen & Boore, 2004) and Verdant
195      (McKain et al., 2017).
196  7.  When determining which ATG codon was a true start codon, we compared the sequence
197      of a gene with sequences of its homologs from the aforementioned reference species.
198  8.  To determine exon borders, RNA-seq reads with a minimum length of 100 bp (to
199      minimize false mappings) were mapped to the genome by CLC Assembly Cell 4.2,
200      requiring at least 50% of each read's length to map with a sequence similarity of at least
201      90%. Exon borders were found by eye in CLC Genomics Workbench 7.5.1 as regions of
202      genes in which there were many partially mapped reads. The exon borders of the
203      reference species were used for comparison.
204  9.  To check for RNA editing that could create new start or stop codons, we mapped RNA-
205      seq reads with a minimum length of 100 bp by CLC Assembly Cell 4.2, requiring at least
206      80% of each read's length to map with a sequence similarity of at least 90%. Mismatches
207      between the reads and the genome were inspected by eye in CLC Genomics Workbench
208      7.5.1.
209  10. After annotating the genes, we additionally verified that there were no remaining regions
210      with high sequence complexity, relatively low AT content or high coverage by RNA-seq
211      reads where no genes were predicted. Regions with high sequence complexity were
212      predicted in the genome by CLC Genomics Workbench 7.5.1 using K2 algorithm
213      (Wootton & Federhen, 1993) with a window size of 101 bp. The AT content plot was
214      created by a custom script with 200-bp-long windows. RNA-seq reads with a minimum
215      length of 100 bp were mapped by CLC Assembly Cell 4.2, requiring at least 80% of each
216      read's length to map with a sequence similarity of at least 90%.

217  After completing gene prediction, the plastid genome was reoriented to start from the first
218  position of *rps14*, as this is the first gene in the canonical representation of the plastid genome of
219  *A. thaliana* which is also present in the plastid genome of *R. phalloides*.

220

221  ## Estimation of contamination amount

222

223  The nuclear genome size could be overestimated if, in addition to the own DNA of *R. phalloides*,
224  contaminating DNA was sequenced. For example, this contamination may originate from
225  endophytic bacteria and fungi. To estimate the amount of contamination, 1000 random DNA-seq
226  read pairs, taken after the trimming, were aligned by BLAST to NCBI databases. Taxonomies of
227  their best matches were used as proxies for the reads' source taxonomies. To increase the
228  sensitivity of the search, the analysis was performed as follows:

229  1. All reads were aligned to NCBI NT (the database current as of September 18, 2017) by
230  BLASTN from BLAST 2.3.0+ suite with the maximum allowed e-value of $10^{-3}$ and the
231  word size of 7 bp. To decrease the number of false-positive matches, hard masking of
232  low-complexity regions ("-soft_masking false" option) was used.
233  2. All reads were aligned to NCBI NR (the database current as of September 18, 2017) by
234  BLASTX from BLAST 2.3.0+ suite with the maximum allowed e-value of $10^{-3}$ and the
235  word size of 3 bp. Hard masking in BLASTX is enabled by default.
236  3. If at least one of two reads in a pair had matches to NT, the taxonomy of the match with
237  the lowest e-value was considered the taxonomy of the read pair. If the read pair had no
238  matches in NT, the taxonomy of the match to NR with the lowest e-value was considered
239  the taxonomy of the read. Therefore, the alignment to NT had higher priority than the
240  alignment to NR. This was done to take into account synonymous positions of genes,
241  where possible, and thus increase the precision of the taxonomic assignment of read
242  pairs.

243

244  ## Other analyses

245

246  To determine the phylogenetic placement of *R. phalloides* within Balanophoraceae, we utilised
247  the alignment of genes from 186 species (180 species of Santalales plus 6 outgroup species)
248  created by Su et al. (2015). *Rhopalocnemis phalloides* was not studied in that article. Seven
249  genes were used for the phylogenetic analysis in that work: plastid *accD*, *matK*, *rbcL*; nuclear
250  18S rDNA, 26S (also known as 25S) rDNA and *RPB2*; and mitochondrial *matR*. As *matK* and
251  *rbcL* are absent from the plastid genome of *R. phalloides*, we were unable to use them. *accD* of
252  *R. phalloides* contains many mutations and thus can be aligned improperly, so we did not use it
253  either. Mitochondrial *matR* is disrupted in *R. phalloides* by several frameshifting indels. Owing
254  to the large size of the nuclear genome of *R. phalloides* (see the paragraph "Other genomes of *R.*
255  *phalloides*"), *RPB2* had a low coverage, and its sequence could not be obtained from the

256    available DNA-seq reads. The sequences of 18S rDNA and 26S rDNA were easier to determine,
257    as they had many copies in the nuclear genome and thus their coverage was higher. To find their
258    sequences among the contigs, we aligned 18S rDNA and 26S rDNA of *A. thaliana* by BLASTN
259    with the default parameters to the contigs of the Spades assembly. The sequences of 18S rDNA
260    and 26S rDNA were added to the alignment of Su et al. (2015) using MAFFT 7.402 (Katoh &
261    Standley, 2013) with options --addfragments and --maxiterate 1000. The phylogenetic tree was
262    built with RAxML 8.2.4 (Stamatakis, 2014), utilising 20 starting stepwise-addition parsimony
263    trees, employing GTR+Gamma model, with the same 6 outgroup species as in the work of Su et
264    al. (2015) (*Antirrhinum majus*, *A. thaliana*, *Camellia japonica*, *Cornus florida*, *Myrtus communis*
265    and *Spinacia oleracea*). The required number of bootstrap pseudoreplicates was determined by
266    RAxML automatically with the extended majority-rule consensus tree criterion ("autoMRE").
267    The tree was visualised with FigTree 1.4.3 (http://tree.bio.ed.ac.uk/software/figtree/).

268    To compare the substitution rate in the plastid genome of *R. phalloides* with substitution rates in
269    other species of Santalales, we used common protein-coding genes of *R. phalloides*, *B. reflexa*,
270    *A. thaliana* (used as the outgroup) and the only four species of Santalales with published plastid
271    genomes as of 2017: *V. album*, *O. alba*, *Champereia manillana* and *Schoepfia jasminodora*.
272    Their protein-coding gene alignment was created by TranslatorX 1.1 (Abascal, Zardoya &
273    Telford, 2010) based on an alignment of the corresponding amino acid sequences performed by
274    Muscle 3.8.31 (Edgar, 2004) with the default parameters. *ycf1*, *ycf2* and *accD* of *R. phalloides*
275    differed from the homologous genes of other species so much that a reliable alignment was not
276    possible. Alignments of other genes were then concatenated into a single alignment and provided
277    to Gblocks server (Castresana, 2000), which removed poorly aligned regions from the alignment.
278    Gblocks was run in the codon mode, with the default parameters. Substitution rates and selective
279    pressure were evaluated by codeml from PAML 4.7 (Yang, 2007) with the F3×4 codon model,
280    starting dN/dS value of 0.5 and starting transition/transversion rate of 2. The phylogenetic tree
281    provided to PAML was a subtree of the large phylogenetic tree of Santalales, produced as
282    described above. Additionally, the analysis of substitution rates and selective pressure was
283    performed by BppSuite 2.3.2 (Guéguen et al., 2013). To the best of our knowledge, this is the
284    only tool that is capable of phylogenetic analyses of protein-coding sequences that takes into
285    account different codon frequencies in different sequences (Guéguen & Duret, 2017), whereas
286    PAML uses a single averaged codon frequency for all sequences. This is important, because the
287    codon frequencies in *R. phalloides* and *B. reflexa* highly differ from the codon frequencies in the
288    mixotrophic Santalales of comparison. The program bppml from BppSuite was run using a
289    nonhomogeneous ("one_per_branch") model, the substitution model was YN98, the codon
290    model F3×4, starting dN/dS values of 0.5 and starting transition/transversion rates of 2. Starting
291    branch lengths were 0.1 substitutions per codon. The parameter estimation was performed by the
292    full-derivatives method with optimization by the Newton-Raphson method
293    ("optimization=FullD(derivatives=Newton)"), using parameters transformation
294    ("optimization.reparametrization=yes").

295    To check for similarity between the genes and the proteins of the *R. phalloides* plastid genome
296    and sequences from other species, we performed BLASTN and BLASTP alignment against

297    NCBI NT and NR databases, respectively, on the NCBI website
298    (https://blast.ncbi.nlm.nih.gov/Blast.cgi) on March 4, 2018 with the default parameters.

299    To build the phylogenetic tree of *rrn16*, sequences from *R. phalloides*, *O. alba*, *V. album*, *N.*
300    *tabacum* and *A. thaliana* were supplemented with sequences from *Corynaea crassa* (NCBI
301    accession U67744), *Balanophora japonica* (NCBI accession KC588390), *Nitzschia* sp. IriIs04
302    (NCBI accession AB899709), *Leucocytozoon caulleryi* (NCBI accession AP013071) and
303    *Plasmodium cynomolgi* (NCBI accession AB471804). These particular sequences of *Nitzschia*,
304    *Leucocytozoon* and *Plasmodium* were randomly chosen among the SAR ("Stramenopiles,
305    Alveolata, Rhizaria", a clade of protists) sequences that produced matches to the *rrn16* of *R.*
306    *phalloides* in the analysis described in the previous paragraph. The *rrn16* sequences were aligned
307    by Muscle 3.8.31. Poorly aligned regions were removed by Gblocks server. The phylogenetic
308    tree was built by RAxML 8.2.4 with the parameters described above. The phylogenetic tree of
309    the species was taken from the TimeTree database (timetree.org). The subtree of
310    Balanophoraceae in the tree was taken from the general tree of Santalales, created as described
311    above. The plastids of species from the SAR clade originated from a secondary endosymbiosis
312    with a red alga, but this endosymbiosis occurred in SAR once in the root, and thus (taking into
313    account that red algae is an outgroup to Embryophyta) the phylogenetic tree of plastids of the
314    studied species coincide with the phylogenetic tree of nuclei. The trees were drawn by
315    TreeGraph 2.14.0-771 beta (Stöver & Müller, 2010).

316    Codon usage and amino acid usage of the common protein-coding genes of *R. phalloides* and
317    species of comparison were calculated by CodonW 1.4.2 (Peden, 1999). Frequencies of 21-bp-
318    long k-mers were calculated for the trimmed DNA-seq reads by Jellyfish 2.1.2 (Marçais &
319    Kingsford, 2011), not using the Bloom filter (to count the number of low-frequency k-mers
320    precisely).

321    The list of plastid genomes with their lengths and AT contents was obtained from the NCBI
322    database (https://www.ncbi.nlm.nih.gov/genome/browse/#!/organelles/). Information on whether
323    a specific plant species is completely heterotrophic was obtained by literature analysis.

324

## Results & Discussion

326

### The gene content of the *R. phalloides* plastid genome

328

329    The plastid genome of *R. phalloides* is circular-mapping and has a length of 18,622 bp long. Its
330    map is represented in Fig. 1, in a linear form, for convenience. The protein-coding gene content
331    is quite typical for highly reduced plastid genomes of completely heterotrophic plants (Graham,
332    Lam & Merckx, 2017; Wicke & Naumann, 2018). The plastid genome of *R. phalloides* possesses
333    the genes *accD*, *clpP*, *ycf1*, *ycf2* (mentioned in the Introduction) and 9 genes encoding protein
334    components of the ribosome. Additionally, it codes for *rrn16* and *rrn23*, RNA components of the
335    plastid ribosome. Like several other highly reduced plastid genomes, it lacks *rrn4.5* and *rrn5*—

336  genes coding for two other RNAs of the ribosome—which poses the interesting puzzle of how
337  the ribosome works without these genes in the plastid genome. One possibility is that these genes
338  were transferred either to the mitochondrial or to the nuclear genome and are now transcribed
339  there and imported to the plastids from the cytoplasm. The other possibility is that the ribosome
340  is capable of working without them, akin to how it can work without some ribosomal proteins
341  (Tiller & Bock, 2014). We plan to clarify this question in an upcoming article dedicated to the
342  analysis of the *R. phalloides* transcriptome.

343  The tRNA-coding gene content of the *R. phalloides* plastid genome is also puzzling. The
344  standard method to predict tRNA-coding genes is the program tRNAscan-SE. It has a dedicated
345  "organellar" mode in which tRNA models were trained on mitochondrial- and plastid-encoded
346  tRNA sequences and structures. It also has a "general" mode whose models are based on nuclear-
347  encoded tRNAs. In the organellar mode, the tool predicts 64 tRNA-coding genes, which is much
348  more than the approximately 30 tRNA-coding genes encoded in plastomes of typical autotrophic
349  species (Wicke et al., 2011). In the general mode, the tool predicts zero tRNA-coding genes. Our
350  experience in working with different plastid genomes suggests that results of predictions in these
351  two modes usually coincide. Of the 64 predicted tRNA-coding genes, 61 have introns, and the
352  mean AT content of the 64 genes is 94%. Therefore, we supposed that most of them, if not all,
353  were false-positive predictions. They could originate from the ease with which sequences of low
354  complexity form secondary structures—these spuriously generated cloverleaf-like structures may
355  have deceived the algorithms of tRNAscan-SE. Seventeen of the predicted tRNA-coding genes
356  were for isoleucine tRNAs, and 11 were for lysine. This further attested to the false-positive
357  nature of these genes, as false-positively predicted tRNA-coding genes in an AT-rich genome are
358  expected to have AT-rich anticodons, and the anticodons of isoleucine and lysine tRNAs are two
359  of the most AT-rich of all amino acid anticodons. Of the three tRNA-coding genes without
360  introns, one has an AT content of 76%, another 92% and the third 96%. Because of the relatively
361  low AT, the first seems to be a possible candidate for a true gene. Its AT content was not only
362  the lowest among the three predictions that do not have introns but also among all 64 predicted
363  tRNA-coding genes. This is a *trnL* gene with anticodon TAA (UAA). Nevertheless, we could not
364  confidently determine whether this gene was a false-positive prediction so we did not use it for
365  any analyses. The only predicted *trnE* gene had an AT content of 99% and is thus very likely to
366  be a false prediction. Therefore, the plastid genome of *R. phalloides* probably lost its *trnE*, like
367  the plastid genomes of completely heterotrophic plants from the genus *Pilostyles* (Bellot &
368  Renner, 2015), although earlier *trnE* was deemed indispensable because of its function in haem
369  synthesis (Howe & Smith, 1991; Barbrook, Howe & Purton, 2006). In the plastid genomes of
370  *Balanophora*, the other genus of Balanophoraceae for which completely sequenced plastid
371  genomes are available, *trnE* is present but is supposed to participate in haem synthesis only,
372  having lost its function in translation (Su et al., 2019). The predicted *trnE* of *R. phalloides* is
373  located in the intergenic region between *ycf1* and *rrn23*, not where it is in *Balanophora*, which is
374  an additional argument for the false-positive nature of this prediction. Potentially, in *R.
375  phalloides, trnE* could have been transferred to the nuclear or the mitochondrial genome,
376  transcribed there and imported into the plastids from the cytoplasm.

377  Overall, the gene content of the plastid genome of *R. phalloides* is similar to the gene content of
378  the plastid genomes of *Balanophora*. All differences between them lie in differential losses of
379  genes participating in translation, except for the loss of *trnE* in *Rhopalocnemis*, that function also
380  in haem synthesis. Compared to *Balanophora*, *R. phalloides* lacks *rps2*, *rps4*, *rps11*, *rpl14*, *trnE*,
381  *rrn4.5*, while *Balanophora* lacks *rpl16* and *rpl36* which are present in *R. phalloides*.

382  Most plastid genes of *R. phalloides* are shorter than their homologues in close mixotrophic
383  relatives, although not as short as homologues in *Balanophora* (Table S1). The compaction of
384  non-coding regions in the plastid genome in *R. phalloides* is also not as pronounced as in
385  *Balanophora*, with 78.5% being coding (i.e. non-intergenic and non-intronic) in the plastid
386  genome of *R. phalloides* and 94.2% and 95.2% in the plastid genomes of *B. reflexa* and *B.*
387  *laxiflora*, respectively. Several genes of *R. phalloides* overlap. Namely, *ycf1* and *ycf2* overlap; so
388  does *rpl36* which overlaps with both *rps7* and *rpl16* (Fig. 1). The intron loss in *R. phalloides* is
389  also not as pronounced as in *Balanophora*, with four cis-spliced and one trans-spliced introns
390  remaining in *R. phalloides*, whereas only the trans-spliced intron remains in *Balanophora*, with
391  all cis-spliced introns lost.

392  The gene order in the plastid genome of *R. phalloides* is neither collinear with the gene order of
393  *Balanophora*, nor it is collinear with the typical gene order of photosynthetic plants (Table S2).
394  Namely, the plastid genome of *R. phalloides* has 7 collinear blocks with the plastid genomes of
395  *Balanophora* and 4 collinear blocks with the plastid genome of *A. thaliana*.

396

## The plastid genome of *R. phalloides* has a very high AT content

398

399  One of the most interesting features of the *R. phalloides* plastid genome is its AT content of
400  86.8%. Among plant genomes, it is surpassed only by the plastid genomes of *B. reflexa* and *B.*
401  *laxiflora*, two close relatives of *R. phalloides*, which have an AT content of 88.4% and 87.8%,
402  respectively (Su et al., 2019). Among prokaryotes and eukaryotes other than plants, there are also
403  several other known genomes with higher AT content, all belonging to mitochondria or
404  apicoplasts, with the record held by the mitochondrial genome of a fungus, *Nakaseomyces*
405  *bacillisporus* CBS 7720 (Bouchier et al., 2009), with the AT content of 89.1% (according to the
406  NCBI site, information current as of June 14, 2018).

407  The increased AT content is a common feature of plastid genomes of completely heterotrophic
408  plants (Fig. 2, Table S3); to the best of our knowledge, it remains unexplained. It correlates with
409  the degree of plastid genome reduction, with plants whose plastid genomes are the most AT rich
410  having simultaneously some of the smallest plastid genomes.

411  It was the high AT content which prevented us from detecting the plastid genome of *R.*
412  *phalloides* from the initial assembly made of approximately 10 million read pairs. High AT
413  content hampers PCR (Benjamini & Speed, 2012), and as library preparation for Illumina
414  sequencing machines usually involves PCR, coverage of AT-rich regions is decreased. When we
415  assembled the genome using an insufficient number of reads, the genome's sequence was broken

416    into multiple contigs containing regions with relatively low AT content. The breaks occurred in
417    the regions in which the AT content was higher; therefore, the coverage in those regions was
418    decreased the most. The obtained sequences were not enough to determine whether the plastid
419    genome was present because (as we describe further below) the sequences were usually similar
420    to those from taxons other than plants, owing to the high AT content and high mutation
421    accumulation rate. Therefore, we initially thought that these short contigs were horizontal
422    transfers located in the mitochondrial genome. Increasing the number of reads allowed us to
423    obtain the full sequence of the plastid genome of *R. phalloides*.

424    The sequencing coverage in the *R. phalloides* plastid genome ranges from approximately 3,000×
425    in the least AT-rich regions to 17 in the most AT-rich regions (Fig. S1). The AT content and the
426    sequencing coverage correlate with a Spearman's correlation coefficient of -0.93. Read insert size
427    also depends on the AT content, with the least AT-rich regions covered by reads with an insert
428    size of approximately 300 bp and the most AT-rich regions with an insert size of approximately
429    200 bp (Fig. S2); Spearman's correlation coefficient was -0.69. We suppose that the coverage
430    drop associated with high AT content could be the reason why the authors of a work dedicated to
431    an analysis of the *Lophophytum mirabile* (also a completely heterotrophic plant from the same
432    family as *R. phalloides*) did not observe contigs with plastid genes (Sanchez-Puerta et al., 2017).
433    Additionally, in *R. lagascae*, which was reported to have no plastid genome (Molina et al.,
434    2014), it may potentially be present but be unnoticed due to its high AT content. *Rafflesia
435    lagascae* genome assembly was performed using approximately 400 million Illumina reads, the
436    same amount we used for the assembly of *R. phalloides*. Therefore, if *R. lagascae* indeed
437    possesses a plastid genome, it should be much more AT-rich than the plastid genome of *R.
438    phalloides*.

439    The AT content is high in protein-coding genes (the average value weighted by length is 88.1%),
440    as well as ncRNA-coding genes (the average value weighted by length is 77.5%) and non-coding
441    regions (the average value weighted by length is 93.8%). In protein-coding genes, this led not
442    only to a shift in codon frequencies towards AT-rich codons (Table S4) but also to a shift in
443    amino acid frequencies in proteins, with amino acids encoded by AT-rich codons used more
444    (Fig. 3, Table S5). For example, isoleucine, the amino acid with the most AT-rich codons, is
445    used two times more often in the proteins encoded in the plastid genome of *R. phalloides* than in
446    homologous proteins of phylogenetically close mixotrophic species. Similarly, glycine, whose
447    codons are among the most GC-rich, is used two times more rarely. Plastid sequences of
448    *Balanophora* experience the same effects. Additionally, the genetic code in the plastid genomes
449    of *Balanophora* is supposed to be non-canonical, utilising TAG (which is a stop codon in most
450    genetic codes) as the tryptophan codon instead of the typical TGG. In contrast, the plastid
451    genome of *R. phalloides* uses TGG for tryptophan, whereas the TAG codon is not used at all,
452    even as a stop codon.

453    Interestingly, such high AT content has led to convergence of gene sequences of *R. phalloides*
454    with sequences from phylogenetically distant AT-rich species. When aligning sequences of
455    genes and proteins of *R. phalloides* to sequences from NCBI NT and NR databases, respectively,
456    the best matches are often sequences from distantly related heterotrophic plants whose plastid

457  genomes also have high AT content (Table S6). There are also many matches to sequences from
458  various protists and some matches to sequences of animals and bacteria. Not all the matches are
459  to homologous sequences, with some resulting from accidental similarity to non-coding
460  sequences.

461  We thoroughly investigated one of the prominent cases of convergence—the *rrn16* gene. This
462  particular gene was selected for the analysis because it is the only gene whose sequences are
463  known for 3 genera of Balanophoraceae (*Rhopalocnemis*, *Corynaea*, *Balanophora*). This
464  allowed us to check whether the convergence with distant species could also be observed in other
465  genera of Balanophoraceae. BLASTN alignment of *rrn16* of *R. phalloides* to NCBI NT produces
466  two best hits to other species of Balanophoraceae, namely *C. crassa* and *B. japonica* (for which
467  only this plastid gene is sequenced and available in GenBank), whereas the next several dozen
468  matches were to protists from the genera *Plasmodium*, *Nitzschia* and *Leucocytozoon*, belonging
469  to SAR. Our initial hypothesis was a horizontal transfer from SAR to a common ancestor of the
470  aforementioned Balanophoraceae. This was supported by the fact that a phylogenetic analysis of
471  *rrn16* places the sequences of Balanophoraceae within SAR with a bootstrap support value of
472  100 (Fig. 4). A simple counterargument is that *Plasmodium*, *Nitzschia* and *Leucocytozoon*,
473  though all belonging to SAR, are, in fact, quite distant phylogenetically from each other (with
474  *Nitzschia* belonging to Stramenopiles, and *Plasmodium* and *Leucocytozoon* to Alveolata), and
475  thus the fact that they appear in BLAST results together suggests some sort of bias. What is
476  common for the species of the genera whose *rrn16* produces best matches to *rrn16* of *R.*
477  *phalloides* is that they have extremely high AT content, close to that of *R. phalloides*. This led us
478  to guess that the similarity originates not from the phylogenetic relatedness of *rrn16* of *R.*
479  *phalloides* to *rrn16* of these species but from convergence because of their high AT content. To
480  verify this, we rebuilt the phylogenetic tree, excluding from the multiple alignment all columns
481  with A or T in *R. phalloides*, thus eliminating the possible convergence originating from the high
482  AT content. In the resulting tree, *R. phalloides* was placed among plants (Fig. 4C), which is
483  correct, confirming that the placement in SAR was owing to the AT content. The removal of
484  columns with A and T in *B. japonica* or *C. crassa* led to similar results: in the produced trees
485  (Fig. S3) these species were situated either within plants (in the case of *B. japonica*) or between
486  plants and SAR (in the case of *C. crassa*). An alternative explanation for the seeming
487  phylogenetic closeness of *rrn16* of these three species of Balanophoraceae to *rrn16* of SAR can
488  be long branch attraction, but it is a characteristic problem of the maximum parsimony method
489  and it affects phylogenetic trees built with the maximum likelihood method to a lesser degree
490  (Kück et al., 2012). Additionally, the similarity of *rrn16* orthologues can potentially be a result
491  of misalignment, but the alignment was good, and the convergence was clearly observed in the
492  alignment (Fig. S4).

493  Overall, our results suggest that phylogenetic analyses of heterotrophic plants (and, in general, of
494  any species whose genomes have highly biased nucleotide composition) should be performed
495  cautiously, as even bootstrap support values of 100 do not guarantee reliable phylogenetic
496  reconstruction in such cases.

497

498 ## Natural selection and substitution rate in the plastid genome of *R.*
499 *phalloides*
500

501 The nucleotide substitution rate is known to be increased in plastid genomes of heterotrophic
502 plants, ranging from a hardly detectable increase in plants that have lost their photosynthetic
503 ability recently (Barrett, Wicke & Sass, 2018) to a nearly 100-fold increase with respect to the
504 closest photosynthetic species in the most reduced plastid genomes (Bellot & Renner, 2015). To
505 the best of our knowledge, the reason for this increase is not yet known (and will be discussed in
506 more details in the section "Why is the AT content so high?").

507 To compare the substitution rate in *R. phalloides* with the rates in its closest mixotrophic
508 relatives, one should first determine the phylogenetic placement of *R. phalloides* relative to the
509 species of comparison. The placement of the family Balanophoraceae has long been debated,
510 with some scientists stating that it does not even belong to Santalales (Kuijt, 1968; Cronquist,
511 1981; Takhtadzhian, 2009). A recent work, which utilised sequences of 7 genes for phylogeny
512 evaluation, suggested that Balanophoraceae indeed belong to Santalales (Su et al., 2015).
513 Moreover, the results of that work suggest polyphyly of Balanophoraceae, which consist of two
514 clades: "Balanophoraceae A" and "Balanophoraceae B". A common feature of Balanophoraceae
515 A is that they have highly increased substitution rates, and a common feature of
516 Balanophoraceae B is that their substitution rates are approximately the same as in autotrophic
517 and mixotrophic Santalales. Although it analysed 11 species of Balanophoraceae, that study did
518 not analyse *R. phalloides*. To estimate the phylogenetic relationships of *R. phalloides*, we added
519 the sequences of its nuclear 18S rDNA and 26S rDNA to the alignment of sequences from 186
520 species used in that article and rebuilt the tree. As one may have expected, *R. phalloides* is
521 placed in Balanophoraceae A, with a bootstrap support value of 100 (Fig. S5). It is sister to a
522 group of *C. crassa* and *Helosis cayennensis*.

523 To evaluate substitution rates, dN and dS in the plastid genome of *R. phalloides*, we used
524 common protein-coding genes of this genome, the plastid genome of *Balanophora reflexa* and
525 plastid genomes of several other species of Santalales, available as of 2017. The genes *ycf1*, *ycf2*
526 and *rps7* were excluded from the analysis because their sequences in *R. phalloides* could not be
527 reliably aligned with homologous sequences of other species owing to the high amount of
528 accumulated mutations. The analysis by PAML showed that the number of nucleotide
529 substitutions in the plastid genome of *R. phalloides* since the divergence from common ancestor
530 with the mixotrophic Santalales of comparison is, on average, 21 times higher than in the plastid
531 genomes of those mixotrophic Santalales (Fig. 5). This number should be treated with caution,
532 as:

533    1. The model of nucleotide substitutions used in PAML utilises the equilibrium codon
534       frequencies, equal for all branches. This is definitely not the case in the studied
535       Santalales, as the codon frequencies in the plastid genome of *R. phalloides* highly differ
536       from those in plastid genomes of mixotrophic Santalales (Table S4).
537       We are aware of a single tool for phylogenetic analyses that can take into account
538       different codon frequencies in different sequences. This is a program collection BppSuite.

539         However, the analysis of these data by BppSuite provided a value of approximately
540         44,000 instead of 21, which was probably owing to an algorithmic mistake.
541   2.  Non-synonymous substitutions quickly reach saturation, and thus the number of non-
542      synonymous substitutions is underestimated for long branches (dos Reis & Yang, 2013).
543      The same is true for synonymous substitutions (Vanneste, Van de Peer & Maere, 2013).
544   3.  We removed columns in the alignment with many differences between species using the
545      program Gblocks, because such columns may result from misalignment. As regions of
546      genes with positive or weak negative selection accumulate mutations faster, such regions
547      can also be potentially removed by Gblocks, leading to underestimation of substitution
548      rates.
549   4.  We failed to produce reliable alignments for the genes ycf1, ycf2 and rps7, consequently
550      the substitution rate in these genes may be higher than in others. Therefore, the exclusion
551      of these genes from the analysis may lead to underestimation of the true substitution rate.

552  The substitution rate analysis for *B. reflexa* provided a very similar result to that of *R. phalloides*.
553  The dN/dS values on the branches of *Balanophoraceae* were slightly lower than on the branches
554  of mixotrophic Santalales. Because the estimation of the dN and dS values could be imprecise,
555  the values of dN/dS should also be treated cautiously. In the future, the problems associated with
556  the analysis of long branches can be reduced by increasing the taxon sampling for
557  *Balanophoraceae*, thus decreasing the branch lengths. Although the precise value of dN/dS on
558  the branch of *R. phalloides* is hard to estimate, the selection acting on its genes is definitely non-
559  neutral, as open reading frames of all the genes are intact. If we denote the probability that there
560  is a specific codon in a specific position as P(X), and the AT content of a gene as α, then the
561  probability that a random codon is a stop is

562  $P(Stop)=P(TAA)+P(TGA)+P(TAG)=(\alpha/2)\times(\alpha/2)\times(\alpha/2)+(\alpha/2)\times((1-\alpha)/2)\times(\alpha/2)+(\alpha/2)\times(\alpha/2)\times((1-$
563  $\alpha)/2)=\alpha^2/4-\alpha^3/8.$

564  As the weighted (by length) average AT content in protein-coding genes of *R. phalloides* is 88%,
565  the probability of a random codon being a stop, as follows from this equation, is approximately
566  11%. This means that because stop codons are AT-rich, in a random sequence with such a high
567  AT content as in *R. phalloides*, every 9th codon will be a stop. Therefore, a strong negative
568  selection must be acting on the genes to keep open reading frames unbroken.

569

570  ## Other genomes of *R. phalloides*

571

572  Sequencing of approximately 400 million paired-end reads could have been enough to assemble
573  the mitochondrial and the nuclear genomes of *R. phalloides*. The alignment by BLASTN and
574  TBLASTN of ncRNAs and proteins, respectively, encoded in mitochondrial genomes of the
575  reference species to the contigs of *R. phalloides* revealed several dozen matching contigs with
576  coverages of approximately 5,000× and lengths of approximately 1,000–5,000 bp. They are
577  probably short mitochondrial chromosomes, similar to those observed in the plant *L. mirabile*
578  (Sanchez-Puerta et al., 2017), also from Balanophoraceae, whose mitochondrial genome

579   putatively consists of 54 small circular chromosomes. We do not plan to investigate the
580   mitochondrial genome of *R. phalloides* in detail and are ready to provide the mitochondrial
581   contigs upon request.

582   Known sizes of nuclear genomes of plants from Santalales vary from approximately 200 Mbp in
583   *Santalum album* (Mahesh et al., 2018) to approximately 100 Gbp in *Viscum album* (Zonneveld,
584   2010). For example, if the nuclear genome size in *R. phalloides* is 500 Mbp, 400 million 150-bp-
585   long reads will produce a coverage of approximately $400\times150/500=120\times$, which is enough for a
586   draft assembly. To estimate the nuclear genome size, we built a k-mer frequency histogram (Fig.
587   S6). The peak of the distribution, corresponding to the k-mer coverage of the nuclear genome,
588   was difficult to determine, but it was below the k-mer coverage value of $2\times$. As the k-mer size
589   (21) was much lower than the read size (150), the read coverage was approximately equal to the
590   k-mer coverage. Therefore, the nuclear genome size could be estimated to be at least
591   $400\times150/2=30,000$ Mbp. Potentially, the genome size could be overestimated if there is a lot of
592   contamination (for example by DNA of endophytic bacteria and fungi), but a taxonomic analysis
593   of reads suggests that contamination in unlikely to be high (Table S7). The assembly of a 30,000-
594   Mbp-long genome is impossible using only the reads produced in the current study. Instead of
595   the complete nuclear genome assembly, we plan to study it by means of transcriptome assembly,
596   which is the subject of our next work.

597

## Why is the AT content so high?

599

600   The increase in the AT content in the plastid genomes of heterotrophic plants, as well as the
601   increase in their substitution rates, are known and much-discussed phenomena (Bromham,
602   Cowman & Lanfear, 2013; Wicke et al., 2016; Hadariová et al., 2018; Wicke & Naumann,
603   2018). However, their origin is still unknown. The simplest hypothesis for the increase in the
604   substitution rate could be the relaxation of selection acting on genes. However, plastid genes of
605   heterotrophic plants usually show no signs of relaxed selection, except for photosynthesis-related
606   genes during pseudogenization. Interestingly, a high AT content and substitution rate have also
607   been observed in plastids of non-photosynthetic protists (such as *Plasmodium*) (Oborník et al.,
608   2009), which lost the genes required for photosynthesis after the transition to a heterotrophic
609   lifestyle. Additionally, both of these phenomena have been observed in genomes of
610   endosymbiotic bacteria (McCutcheon & Moran, 2011), which may be dozens of times shorter
611   than genomes of their free-living relatives owing to the loss of genes required, for example, for
612   biosynthesis of substances that are now provided to the symbiont by its host. Therefore, these
613   two phenomena are probably not only unrestricted to plants but are not even related to the loss of
614   photosynthesis.

615   A phenomenon that can simultaneously result in both an increase of AT content and an increase
616   of substitution rate is the reduction in genome recombination intensity. A plastid genome is
617   capable of recombining both within itself (the recombination of two copies of the inverted
618   repeat) (Zhu et al., 2016; Li et al., 2016) and between two copies of a genome (Maréchal &

619  Brisson, 2010). The recombination is an important step in repair, both in plastids (Zampini et al.,
620  2017) and in bacteria (Cox, 1998), so the reduction in recombination will increase the
621  substitution rate. Also, gene conversion in plastid (Wu & Chaw, 2015; Zhitao Niu et al., 2017) as
622  well as in bacterial (Lassalle et al., 2015) genomes is GC-biased, although earlier gene
623  conversion in plastid genomes was supposed to be AT-biased (Khakhlova & Bock, 2006). This
624  means that if there is a mismatch between an adenine or a thymine on one strand versus a
625  guanine or a cytosine on the other during a recombination, it is more likely that the guanine or
626  the cytosine will be kept, while the adenine or the thymine will be removed and replaced by a
627  cytosine or a guanine, which is complementary to the base on the other strand. Therefore,
628  recombination aids in increasing the GC content in plastid and bacterial genomes, and a decrease
629  in recombination will make a genome more AT-rich. The link between the low recombination
630  rate and the high AT content has already been proposed for endosymbiotic bacteria with small
631  genomes (Lassalle et al., 2015).

632  Recently, it was shown that in transcriptomes of the heterotrophic plants *Epipogium aphyllum*, *E.*
633  *roseum* and *Hypopitys monotropa*, the transcript of the protein RECA1, which is required for
634  recombination of the plastid genomes, is absent (Schelkunov, Penin & Logacheva, 2018). This
635  may support the above hypothesis. However, the direct reason for the loss of RECA1 is not
636  known. A potential explanation for the loss could be that during the transition from a
637  mixotrophic to a heterotrophic lifestyle, plastid enzymes related to photosynthesis accumulate
638  mutations, and since a mutated enzyme may be harmful for the organism, it is evolutionarily
639  adaptive to accumulate the mutations very fast to quickly achieve complete disruption of a gene,
640  instead of having a semi-degraded gene encoding a harmful protein. This effect, consisting of
641  elimination of pseudogenes at rates faster than neutral, has already been shown to take place in
642  bacteria (Kuo & Ochman, 2010). Therefore, in the period directly following the loss of
643  photosynthesis, it may be beneficial for the plant to disturb the plastid recombination and thus
644  disturb the repair. In fact, this process may start even before the loss of photosynthesis, because
645  in plastid genomes of mixotrophic plants *ndh* genes often undergo pseudogenization (Wicke et
646  al., 2011), and their quick removal may require an increased mutation accumulation rate. Such an
647  increase in the mutation accumulation rate may require pseudogenization of genes of DNA
648  replication, recombination and repair (DNA-RRR), such as *RECA1*, and once they are
649  pseudogenised, it will be hard for a plant to return to the normal repair intensity in the plastid
650  genome, making the transition to high mutation accumulation rates irreversible.

651  It is known that the mutation accumulation rate in heterotrophic plants (Bromham, Cowman &
652  Lanfear, 2013), including Balanophoraceae (Su & Hu, 2012; Su et al., 2015), is also increased in
653  nuclear and mitochondrial genomes, although to a lesser extent than in plastid genomes. These
654  phenomena are also still unexplained. The nuclear genome contains more than a hundred
655  (Schelkunov, Penin & Logacheva, 2018) genes that encode proteins, working in multisubunit
656  complexes with proteins, encoded in the plastid genome. These are the genes encoding proteins
657  of the electron-transfer chain, the plastid-encoded RNA polymerase (PEP), the plastid ribosome
658  and others. When a species loses its photosynthetic ability, the nuclear-encoded genes of the
659  electron-transfer chain are no longer under selective pressure and start to accumulate mutations.
660  Therefore, their proteins may become harmful and may require quick elimination. Thus, the

661 increase in the nuclear mutation accumulation rate, which may speed up the accumulation of
662 disruptive mutations in these genes, may also by selectively beneficial. The increase in the
663 mutation accumulation rate in the mitochondrial genome could potentially be explained by the
664 fact that many DNA-RRR proteins are common for the plastid and the mitochondrial genomes
665 (Shedge et al., 2007; Carrie & Small, 2013). Therefore, if it is selectively beneficial to increase
666 the mutation accumulation rate in the plastid genome, the mitochondrial genome may also be
667 affected.

668 This hypothesis of accelerated junk removal may be tested by studying plastid and nuclear
669 genomes of many related heterotrophic species and checking whether the crumbling genes
670 accumulate mutations at rates faster than neutral shortly after the loss of photosynthesis and
671 whether some of the DNA-RRR genes deteriorate at the same time.

672

## Conclusions and future studies

674 The plastid genome of *R. phalloides* profoundly differs from plastid genomes of typical plants,
675 including the massive gene loss, the increased substitution rate and the high AT content. By
676 decreasing sequencing coverage, such high AT content may "hide" plastid genomes of some
677 heterotrophic plants, making these genomes harder to find by means of high-throughput
678 sequencing. Alterations in the nuclear genome, accompanying these changes in the plastid
679 genome, are an interesting issue. Our next work will be dedicated to the study of the nuclear
680 genome of *R. phalloides* by means of transcriptome sequencing.

681

## Acknowledgements

686

## References

688 Abascal F, Zardoya R, Telford MJ. 2010. TranslatorX: multiple alignment of nucleotide

689     sequences guided by amino acid translations. *Nucleic Acids Research* 38:W7–W13. DOI:

690     10.1093/nar/gkq291.

691 Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko

692     SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev

693     MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and its applications

694   to single-cell sequencing. *Journal of Computational Biology* 19:455–477. DOI:

695   10.1089/cmb.2012.0021.

696 Barbrook AC, Howe CJ, Purton S. 2006. Why are plastid genomes retained in non-

697   photosynthetic organisms? *Trends in Plant Science* 11:101–108. DOI:

698   10.1016/j.tplants.2005.12.004.

699 Barrett CF, Freudenstein JV, Li J, Mayfield-Jones DR, Perez L, Pires JC, Santos C. 2014.

700   Investigating the path of plastid genome degradation in an early-transitional clade of

701   heterotrophic orchids, and implications for heterotrophic angiosperms. *Molecular Biology*

702   *and Evolution* 31:3095–3112. DOI: 10.1093/molbev/msu252.

703 Barrett CF, Wicke S, Sass C. 2018. Dense infraspecific sampling reveals rapid and independent

704   trajectories of plastome degradation in a heterotrophic orchid complex. *New Phytologist*

705   218:1192–1204. DOI: 10.1111/nph.15072.

706 Bellot S, Renner SS. 2015. The plastomes of two species in the endoparasite genus *Pilostyles*

707   (Apodanthaceae) each retain just five or six possibly functional genes. *Genome Biology*

708   *and Evolution*:evv251. DOI: 10.1093/gbe/evv251.

709 Benjamini Y, Speed TP. 2012. Summarizing and correcting the GC content bias in high-

710   throughput sequencing. *Nucleic Acids Research* 40:e72. DOI: 10.1093/nar/gks001.

711 Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence

712   data. *Bioinformatics* 30:2114–2120. DOI: 10.1093/bioinformatics/btu170.

713 Bouchier C, Ma L, Créno S, Dujon B, Fairhead C. 2009. Complete mitochondrial genome

714   sequences of three Nakaseomyces species reveal invasion by palindromic GC clusters

715   and considerable size expansion. *FEMS yeast research* 9:1283–1292. DOI:

716   10.1111/j.1567-1364.2009.00551.x.

717   Bromham L, Cowman PF, Lanfear R. 2013. Parasitic plants have increased rates of molecular

718        evolution across all three genomes. *BMC evolutionary biology* 13:126.

719   Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.

720        BLAST+: architecture and applications. *BMC Bioinformatics* 10:421. DOI:

721        10.1186/1471-2105-10-421.

722   Carrie C, Small I. 2013. A reevaluation of dual-targeting of proteins to mitochondria and

723        chloroplasts. *Biochimica et Biophysica Acta (BBA) - Molecular Cell Research* 1833:253–

724        259. DOI: 10.1016/j.bbamcr.2012.05.029.

725   Castresana J. 2000. Selection of conserved blocks from multiple alignments for their use in

726        phylogenetic analysis. *Molecular Biology and Evolution* 17:540–552.

727   Cox MM. 1998. A broadening view of recombinational DNA repair in bacteria. *Genes to Cells*

728        3:65–78. DOI: 10.1046/j.1365-2443.1998.00175.x.

729   Cronquist A. 1981. *An integrated system of classification of flowering plants*. New York:

730        Columbia University Press.

731   Cummings MP, Welschmeyer NA. 1998. Pigment composition of putatively achlorophyllous

732        angiosperms. *Plant Systematics and Evolution* 210:105–111. DOI: 10.1007/BF00984730.

733   Daniell H, Lin C-S, Yu M, Chang W-J. 2016. Chloroplast genomes: diversity, evolution, and

734        applications in genetic engineering. *Genome Biology* 17. DOI: 10.1186/s13059-016-

735        1004-2.

736   Doyle J. 1987. A rapid DNA isolation procedure for small quantities of fresh leaf tissue.

737        *Phytochem Bull* 19:11–15.

738   Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high

739        throughput. *Nucleic Acids Research* 32:1792–1797. DOI: 10.1093/nar/gkh340.

740    Finn RD, Attwood TK, Babbitt PC, Bateman A, Bork P, Bridge AJ, Chang H-Y, Dosztányi Z,

741         El-Gebali S, Fraser M, Gough J, Haft D, Holliday GL, Huang H, Huang X, Letunic I,

742         Lopez R, Lu S, Marchler-Bauer A, Mi H, Mistry J, Natale DA, Necci M, Nuka G,

743         Orengo CA, Park Y, Pesseat S, Piovesan D, Potter SC, Rawlings ND, Redaschi N,

744         Richardson L, Rivoire C, Sangrador-Vegas A, Sigrist C, Sillitoe I, Smithers B, Squizzato

745         S, Sutton G, Thanki N, Thomas PD, Tosatto SCE, Wu CH, Xenarios I, Yeh L-S, Young

746         S-Y, Mitchell AL. 2017. InterPro in 2017—beyond protein family and domain

747         annotations. *Nucleic Acids Research* 45:D190–D199. DOI: 10.1093/nar/gkw1107.

748    Graham SW, Lam VKY, Merckx VSFT. 2017. Plastomes on the edge: the evolutionary

749         breakdown of mycoheterotroph plastid genomes. *New Phytologist* 214:48–55. DOI:

750         10.1111/nph.14398.

751    Gualberto JM, Newton KJ. 2017. Plant Mitochondrial Genomes: Dynamics and Mechanisms of

752         Mutation. *Annual Review of Plant Biology* 68:225–252. DOI: 10.1146/annurev-arplant-

753         043015-112232.

754    Guéguen L, Duret L. 2017. Unbiased estimate of synonymous and non-synonymous substitution

755         rates with non-stationary base composition. *Molecular Biology and Evolution*. DOI:

756         10.1093/molbev/msx308.

757    Guéguen L, Gaillard S, Boussau B, Gouy M, Groussin M, Rochette NC, Bigot T, Fournier D,

758         Pouyet F, Cahais V, Bernard A, Scornavacca C, Nabholz B, Haudry A, Dachary L,

759         Galtier N, Belkhir K, Dutheil JY. 2013. Bio++: efficient extensible libraries and tools for

760         computational molecular evolution. *Molecular Biology and Evolution* 30:1745–1750.

761         DOI: 10.1093/molbev/mst097.

762    Hadariová L, Vesteg M, Hampl V, Krajčovič J. 2018. Reductive evolution of chloroplasts in

763           non-photosynthetic plants, algae and protists. *Current Genetics* 64:365–387. DOI:

764           10.1007/s00294-017-0761-0.

765    Howe CJ, Smith A. 1991. Plants without chlorophyll. *Nature* 349:109–109. DOI:

766           10.1038/349109c0.

767    Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell

768           A, Nuka G, Pesseat S, Quinn AF, Sangrador-Vegas A, Scheremetjew M, Yong S-Y,

769           Lopez R, Hunter S. 2014. InterProScan 5: genome-scale protein function classification.

770           *Bioinformatics* 30:1236–1240. DOI: 10.1093/bioinformatics/btu031.

771    Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:

772           improvements in performance and usability. *Molecular Biology and Evolution* 30:772–

773           780. DOI: 10.1093/molbev/mst010.

774    Khakhlova O, Bock R. 2006. Elimination of deleterious mutations in plastid genomes by gene

775           conversion. *The Plant Journal* 46:85–94. DOI: 10.1111/j.1365-313X.2006.02673.x.

776    Krause K. 2015. Grand-scale theft: Kleptoplasty in parasitic plants? *Trends in Plant Science*

777           20:196–198. DOI: 10.1016/j.tplants.2015.03.005.

778    Kück P, Mayer C, Wägele J-W, Misof B. 2012. Long Branch Effects Distort Maximum

779           Likelihood Phylogenies in Simulations Despite Selection of the Correct Model. *PLoS*

780           *ONE* 7:e36593. DOI: 10.1371/journal.pone.0036593.

781    Kuijt J. 1968. Mutual Affinities of Santalalean Families. *Brittonia* 20:136. DOI:

782           10.2307/2805616.

783  Kumar AM, Schaub U, Söll D, Ujwal ML. 1996. Glutamyl-transfer RNA: at the crossroad

784       between chlorophyll and protein biosynthesis. *Trends in Plant Science* 1:371–376. DOI:

785       10.1016/S1360-1385(96)80311-6.

786  Kuo C-H, Ochman H. 2010. The Extinction Dynamics of Bacterial Pseudogenes. *PLoS Genetics*

787       6:e1001050. DOI: 10.1371/journal.pgen.1001050.

788  Lagesen K, Hallin P, Rødland EA, Staerfeldt H-H, Rognes T, Ussery DW. 2007. RNAmmer:

789       consistent and rapid annotation of ribosomal RNA genes. *Nucleic Acids Research*

790       35:3100–3108. DOI: 10.1093/nar/gkm160.

791  Lam VKY, Soto Gomez M, Graham SW. 2015. The highly reduced plastome of

792       mycoheterotrophic *Sciaphila* (Triuridaceae) is colinear with its green relatives and is

793       under strong purifying selection. *Genome Biology and Evolution* 7:2220–2236. DOI:

794       10.1093/gbe/evv134.

795  Lassalle F, Périan S, Bataillon T, Nesme X, Duret L, Daubin V. 2015. GC-Content Evolution in

796       Bacterial Genomes: The Biased Gene Conversion Hypothesis Expands. *PLOS Genetics*

797       11:e1004941. DOI: 10.1371/journal.pgen.1004941.

798  Li F-W, Kuo L-Y, Pryer KM, Rothfels CJ. 2016. Genes Translocated into the Plastid Inverted

799       Repeat Show Decelerated Substitution Rates and Elevated GC Content. *Genome Biology*

800       *and Evolution* 8:2452–2458. DOI: 10.1093/gbe/evw167.

801  Liu T-J, Zhang C-Y, Yan H-F, Zhang L, Ge X-J, Hao G. 2016. Complete plastid genome

802       sequence of *Primula sinensis* (Primulaceae): structure comparison, sequence variation

803       and evidence for *accD* transfer to nucleus. *PeerJ* 4:e2101. DOI: 10.7717/peerj.2101.

804    Logacheva MD, Schelkunov MI, Penin AA. 2011. Sequencing and analysis of plastid genome in

805        mycoheterotrophic orchid *Neottia nidus-avis*. *Genome Biology and Evolution* 3:1296–

806        1303. DOI: 10.1093/gbe/evr102.

807    Lowe TM, Eddy SR. 1997. tRNAscan-SE: a program for improved detection of transfer RNA

808        genes in genomic sequence. *Nucleic Acids Research* 25:955–964.

809    Mahesh HB, Subba P, Advani J, Shirke MD, Loganathan RM, Chandana SL, Shilpa S,

810        Chatterjee O, Pinto SM, Prasad TSK, Gowda M. 2018. Multi-Omics Driven Assembly

811        and Annotation of the Sandalwood ( *Santalum album* ) Genome. *Plant Physiology*

812        176:2772–2788. DOI: 10.1104/pp.17.01764.

813    Marçais G, Kingsford C. 2011. A fast, lock-free approach for efficient parallel counting of

814        occurrences of k-mers. *Bioinformatics* 27:764–770. DOI: 10.1093/bioinformatics/btr011.

815    Maréchal A, Brisson N. 2010. Recombination and the maintenance of plant organelle genome

816        stability. *New Phytologist* 186:299–317. DOI: 10.1111/j.1469-8137.2010.03195.x.

817    McCutcheon JP, Moran NA. 2011. Extreme genome reduction in symbiotic bacteria. *Nature*

818        *Reviews Microbiology* 10:13–26. DOI: 10.1038/nrmicro2670.

819    McKain MR, Hartsock RH, Wohl MM, Kellogg EA. 2017. Verdant: automated annotation,

820        alignment and phylogenetic analysis of whole chloroplast genomes. *Bioinformatics*

821        33:130–132. DOI: 10.1093/bioinformatics/btw583.

822    Merckx VSFT, Freudenstein JV, Kissling J, Christenhusz MJM, Stotler RE, Crandall-Stotler B,

823        Wickett N, Rudall PJ, Maas-van de Kamer H, Maas PJM. 2013. Taxonomy and

824        Classification. In: Merckx V ed. *Mycoheterotrophy*. New York, NY: Springer New York,

825        19–101. DOI: 10.1007/978-1-4614-5209-6_2.

826    Millen RS, Olmstead RG, Adams KL, Palmer JD, Lao NT, Heggie L, Kavanagh TA, Hibberd

827        JM, Gray JC, Morden CW, Calie PJ, Jermiin LS, Wolfe KH. 2001. Many parallel losses

828        of infA from chloroplast DNA during angiosperm evolution with multiple independent

829        transfers to the nucleus. *The Plant Cell* 13:645–658.

830    Molina J, Hazzouri KM, Nickrent D, Geisler M, Meyer RS, Pentony MM, Flowers JM, Pelser P,

831        Barcelona J, Inovejas SA, Uy I, Yuan W, Wilkins O, Michel C-I, LockLear S,

832        Concepcion GP, Purugganan MD. 2014. Possible loss of the chloroplast genome in the

833        parasitic flowering plant *Rafflesia lagascae* (*Rafflesiaceae*). *Molecular Biology and*

834        *Evolution* 31:793–803. DOI: 10.1093/molbev/msu051.

835    Naumann J, Der JP, Wafula EK, Jones SS, Wagner ST, Honaas LA, Ralph PE, Bolin JF, Maass

836        E, Neinhuis C, Wanke S, dePamphilis CW. 2016. Detecting and characterizing the highly

837        divergent plastid genome of the nonphotosynthetic parasitic plant *Hydnora visseri*

838        (Hydnoraceae). *Genome Biology and Evolution* 8:345–363. DOI: 10.1093/gbe/evv256.

839    Nawrocki EP, Burge SW, Bateman A, Daub J, Eberhardt RY, Eddy SR, Floden EW, Gardner

840        PP, Jones TA, Tate J, Finn RD. 2015. Rfam 12.0: updates to the RNA families database.

841        *Nucleic Acids Research* 43:D130–D137. DOI: 10.1093/nar/gku1063.

842    Nawrocki EP, Eddy SR. 2013. Infernal 1.1: 100-fold faster RNA homology searches.

843        *Bioinformatics* 29:2933–2935. DOI: 10.1093/bioinformatics/btt509.

844    Oborník M, Janouškovec J, Chrudimský T, Lukeš J. 2009. Evolution of the apicoplast and its

845        hosts: From heterotrophy to autotrophy and back again. *International Journal for*

846        *Parasitology* 39:1–12. DOI: 10.1016/j.ijpara.2008.07.010.

847    Peden JF. 1999. Analysis of codon usage. PhD Thesis Thesis. UK: University of Nottingham.

848    dos Reis M, Yang Z. 2013. Why Do More Divergent Sequences Produce Smaller

849        Nonsynonymous/Synonymous Rate Ratios in Pairwise Sequence Comparisons? *Genetics*

850        195:195–204. DOI: 10.1534/genetics.113.152025.

851    Rousseau-Gueutin M, Huang X, Higginson E, Ayliffe M, Day A, Timmis JN. 2013. Potential

852        functional replacement of the plastidic acetyl-CoA carboxylase subunit (accD) gene by

853        recent transfers to the nucleus in some angiosperm lineages. *Plant Physiology* 161:1918–

854        1929. DOI: 10.1104/pp.113.214528.

855    Sakamoto W, Takami T. 2018. Chloroplast DNA Dynamics: Copy Number, Quality Control and

856        Degradation. *Plant and Cell Physiology* 59:1120–1127. DOI: 10.1093/pcp/pcy084.

857    Sanchez-Puerta MV, García LE, Wohlfeiler J, Ceriotti LF. 2017. Unparalleled replacement of

858        native mitochondrial genes by foreign homologs in a holoparasitic plant. *New Phytologist*

859        214:376–387. DOI: 10.1111/nph.14361.

860    Schelkunov MI, Penin AA, Logacheva MD. 2018. RNA-seq highlights parallel and contrasting

861        patterns in the evolution of the nuclear genome of fully mycoheterotrophic plants. *BMC*

862        *Genomics* 19. DOI: 10.1186/s12864-018-4968-3.

863    Schelkunov MI, Shtratnikova VY, Nuraliev MS, Selosse M-A, Penin AA, Logacheva MD. 2015.

864        Exploring the limits for reduction of plastid genomes: a case study of the

865        mycoheterotrophic orchids *Epipogium aphyllum* and *Epipogium roseum*. *Genome*

866        *Biology and Evolution* 7:1179–1191. DOI: 10.1093/gbe/evv019.

867    Seregin A. 2018. Moscow University Herbarium (MW). DOI: 10.15468/cpnhcc.

868    Shah N, Nute MG, Warnow T, Pop M. 2018. Misunderstood parameter of NCBI BLAST

869        impacts the correctness of bioinformatics workflows. *Bioinformatics*. DOI:

870        10.1093/bioinformatics/bty833.

871    Shedge V, Arrieta-Montiel M, Christensen AC, Mackenzie SA. 2007. Plant Mitochondrial

872         Recombination Surveillance Requires Unusual RecA and MutS Homologs. *THE PLANT*

873         *CELL ONLINE* 19:1251–1264. DOI: 10.1105/tpc.106.048355.

874    Smith DR. 2012. Updating Our View of Organelle Genome Nucleotide Landscape. *Frontiers in*

875         *Genetics* 3. DOI: 10.3389/fgene.2012.00175.

876    Smith DR, Lee RW. 2014. A plastid without a genome: evidence from the nonphotosynthetic

877         green algal genus *Polytomella*. *Plant physiology* 164:1812–1819. DOI:

878         10.1104/pp.113.233718.

879    Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of

880         large phylogenies. *Bioinformatics* 30:1312–1313. DOI: 10.1093/bioinformatics/btu033.

881    Stöver BC, Müller KF. 2010. TreeGraph 2: Combining and visualizing evidence from different

882         phylogenetic analyses. *BMC Bioinformatics* 11:7. DOI: 10.1186/1471-2105-11-7.

883    Su H-J, Barkman TJ, Hao W, Jones SS, Naumann J, Skippington E, Wafula EK, Hu J-M, Palmer

884         JD, dePamphilis CW. 2019. Novel genetic code and record-setting AT-richness in the

885         highly reduced plastid genome of the holoparasitic plant Balanophora. *Proceedings of the*

886         *National Academy of Sciences of the United States of America* 116:934–943. DOI:

887         10.1073/pnas.1816822116.

888    Su H-J, Hu J-M. 2012. Rate heterogeneity in six protein-coding genes from the holoparasite

889         Balanophora (Balanophoraceae) and other taxa of Santalales. *Annals of Botany*

890         110:1137–1147. DOI: 10.1093/aob/mcs197.

891    Su H-J, Hu J-M, Anderson FE, Der JP, Nickrent DL. 2015. Phylogenetic relationships of

892         Santalales with insights into the origins of holoparasitic Balanophoraceae. *Taxon* 64:491–

893         506. DOI: 10.12705/643.2.

894     Takhtadzhian AL. 2009. *Flowering plants*. New York: Springer.

895     Tiller N, Bock R. 2014. The translational apparatus of plastids and its role in plant development.

896          *Molecular Plant* 7:1105–1120. DOI: 10.1093/mp/ssu022.

897     Vanneste K, Van de Peer Y, Maere S. 2013. Inference of genome duplications from age

898          distributions revisited. *Molecular Biology and Evolution* 30:177–190. DOI:

899          10.1093/molbev/mss214.

900     Westwood JH, Yoder JI, Timko MP, dePamphilis CW. 2010. The evolution of parasitism in

901          plants. *Trends in Plant Science* 15:227–235. DOI: 10.1016/j.tplants.2010.01.004.

902     Wicke S, Müller KF, dePamphilis CW, Quandt D, Bellot S, Schneeweiss GM. 2016. Mechanistic

903          model of evolutionary rate variation en route to a nonphotosynthetic lifestyle in plants.

904          *Proceedings of the National Academy of Sciences* 113:9045–9050. DOI:

905          10.1073/pnas.1607576113.

906     Wicke S, Naumann J. 2018. Molecular Evolution of Plastid Genomes in Parasitic Flowering

907          Plants. In: *Advances in Botanical Research*. Elsevier, 315–347. DOI:

908          10.1016/bs.abr.2017.11.014.

909     Wicke S, Schneeweiss GM, dePamphilis CW, Müller KF, Quandt D. 2011. The evolution of the

910          plastid chromosome in land plants: gene content, gene order, gene function. *Plant

911          Molecular Biology* 76:273–297. DOI: 10.1007/s11103-011-9762-4.

912     Wootton JC, Federhen S. 1993. Statistics of local complexity in amino acid sequences and

913          sequence databases. *Computers & Chemistry* 17:149–163. DOI: 10.1016/0097-

914          8485(93)85006-X.

915 Wu C-S, Chaw S-M. 2015. Evolutionary Stasis in Cycad Plastomes and the First Case of

916         Plastome GC-Biased Gene Conversion. *Genome Biology and Evolution* 7:2000–2009.

917         DOI: 10.1093/gbe/evv125.

918 Wyman SK, Jansen RK, Boore JL. 2004. Automatic annotation of organellar genomes with

919         DOGMA. *Bioinformatics* 20:3252–3255. DOI: 10.1093/bioinformatics/bth352.

920 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular Biology and*

921         *Evolution* 24:1586–1591. DOI: 10.1093/molbev/msm088.

922 Zampini É, Truche S, Lepage É, Tremblay-Belzile S, Brisson N. 2017. Plastid Genome Stability

923         and Repair. In: Li X-Q ed. *Somatic Genome Variation in Animals, Plants, and*

924         *Microorganisms*. Hoboken, NJ, USA: John Wiley & Sons, Inc., 119–163. DOI:

925         10.1002/9781118647110.ch7.

926 Zhitao Niu, Qingyun Xue, Hui Wang, Xuezhu Xie, Shuying Zhu, Wei Liu, Xiaoyu Ding. 2017.

927         Mutational Biases and GC-Biased Gene Conversion Affect GC Content in the Plastomes

928         of Dendrobium Genus. *International Journal of Molecular Sciences* 18:2307. DOI:

929         10.3390/ijms18112307.

930 Zhu A, Guo W, Gupta S, Fan W, Mower JP. 2016. Evolutionary dynamics of the plastid inverted

931         repeat: the effects of expansion, contraction, and loss on substitution rates. *New*

932         *Phytologist* 209:1747–1756. DOI: 10.1111/nph.13743.

933 Zonneveld BJM. 2010. New Record Holders for Maximum Genome Size in Eudicots and

934         Monocots. *Journal of Botany* 2010:1–4. DOI: 10.1155/2010/527357.

935

# Figure 1

Map of the *Rhopalocnemis phalloides* plastid genome showing various features.

The circular-mapping plastid genome is represented linearly for convenience. Green arrows are rRNA-coding genes, red arrows are ribosomal protein-coding genes, and blue arrows are genes coding proteins with other functions. Grey arcs represent splicing. Blue columns show non-coding regions.
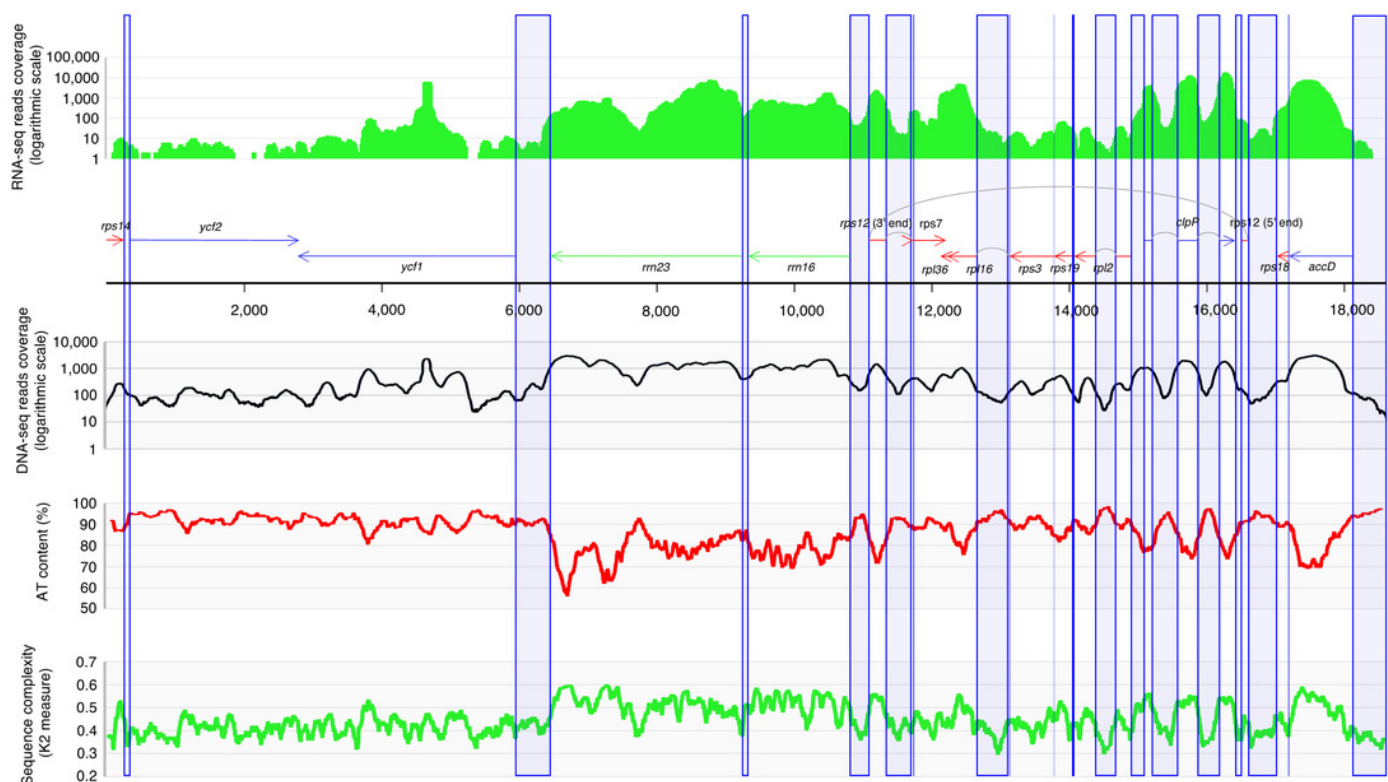
# Figure 2

AT content and lengths of the plastid genomes of Embryophyta.

Red dots denote completely heterotrophic plants and black dots mixotrophic and completely autotrophic.

PeerJ

# Figure 3

Amino acid frequencies in the plastid proteins of *Rhopalocnemis phalloides* and *Balanophora reflexa* are affected by the high AT content.
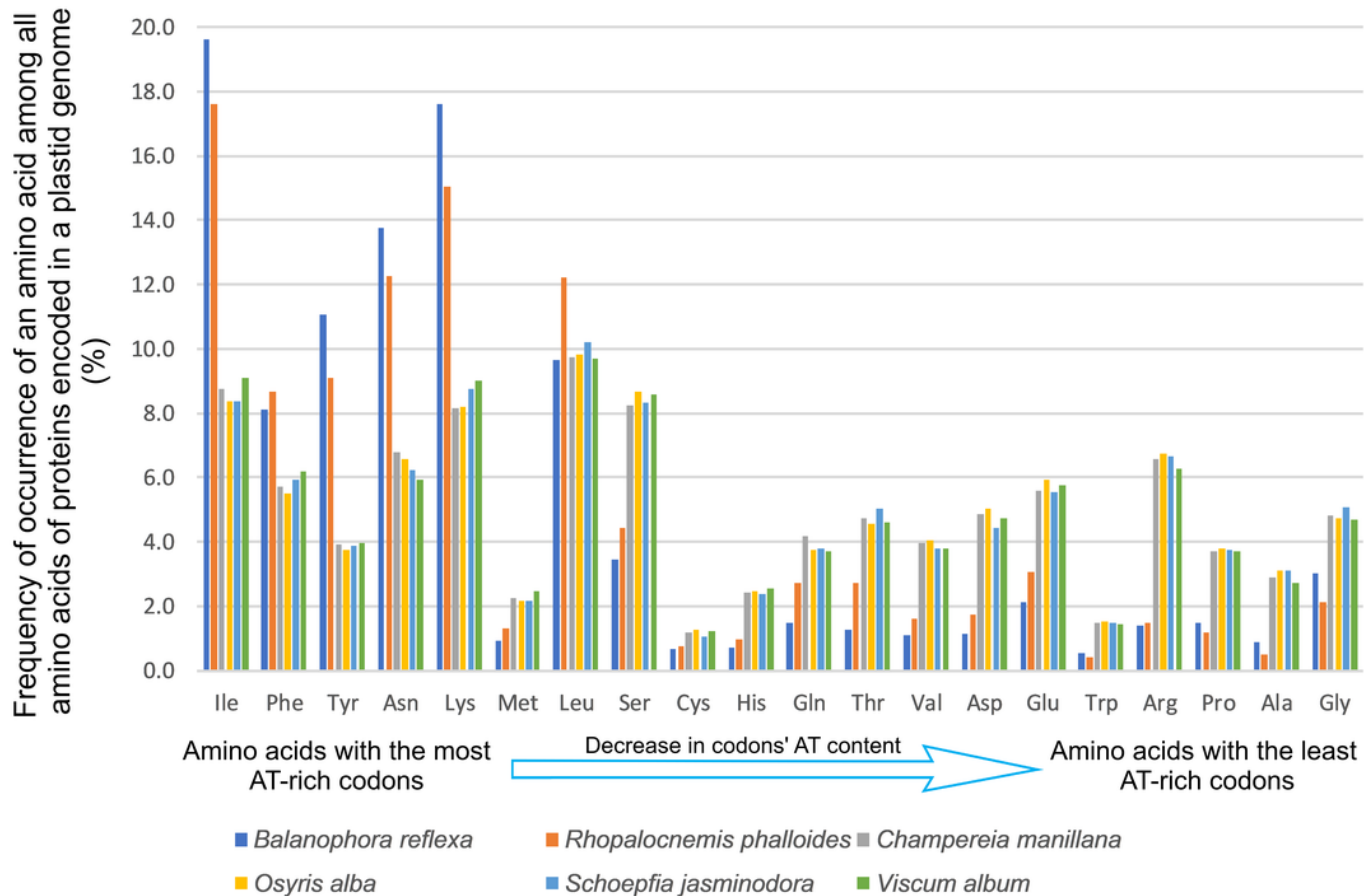
# Figure 4

*rrn16* of *Rhopalocnemis phalloides* shows convergence with *rrn16* from SAR owing to the high AT content.

(A) Phylogenetic tree of species. (B) Phylogenetic tree of *rrn16*. (C) Phylogenetic tree of *rrn16*, built by alignment columns in which *Rhopalocnemis phalloides* has guanine or cytosine. Species with names in orange rectangles are non-photosynthetic plants from Balanophoraceae, species with names in green rectangles are photosynthetic plants and species with names in purple rectangles are from SAR. The numbers on the branches are bootstrap support values. The second and the third tree are unrooted.
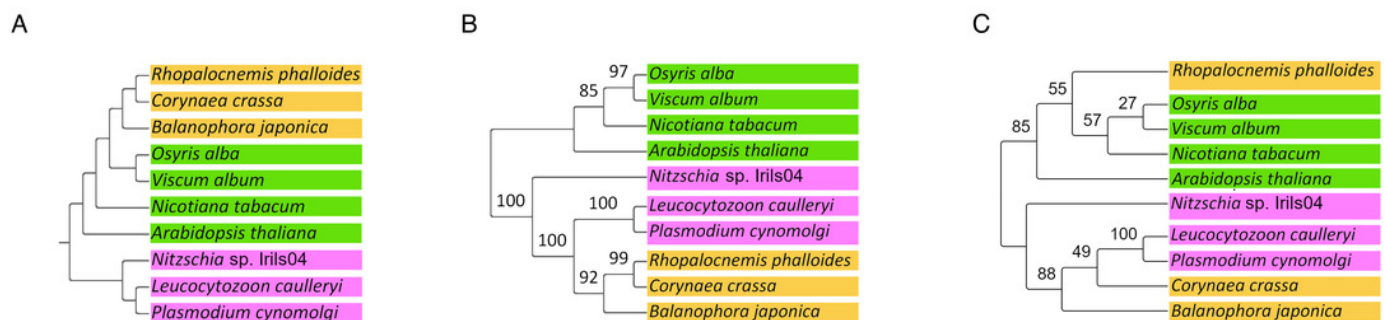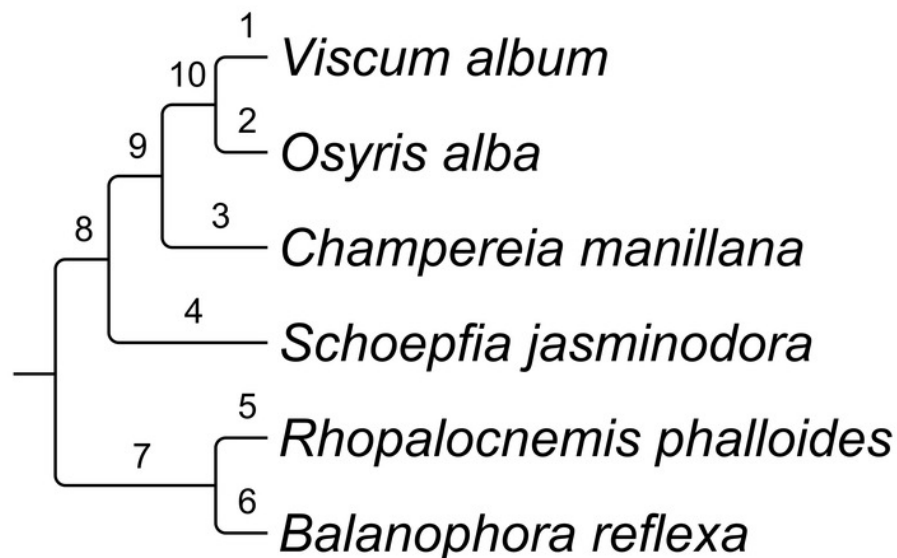
# Figure 5

Evolutionary parameters of the phylogeny of Santalales.

*Arabidopsis thaliana*, used as the outgroup, is not shown. The total length of the alignment, used for the analysis, was 3,363 bp after removal of poorly aligned regions by Gblocks. * dN/dS on this branch cannot be calculated owing to a very small dS value.

| Branch number | Substitutions per position | dN | dS | dN/dS |
|---|---|---|---|---|
| 1 | 0.069 | 0.034 | 0.219 | 0.156 |
| 2 | 0.024 | 0.010 | 0.082 | 0.128 |
| 3 | 0.034 | 0.016 | 0.108 | 0.153 |
| 4 | 0.065 | 0.024 | 0.240 | 0.098 |
| 5 | 0.747 | 0.207 | 3.067 | 0.068 |
| 6 | 0.588 | 0.179 | 2.345 | 0.076 |
| 7 | 0.593 | 0.189 | 2.326 | 0.081 |
| 8 | 0.012 | 0.014 | 0.000 | N/A* |
| 9 | 0.005 | 0.002 | 0.018 | 0.095 |
| 10 | 0.004 | 0.002 | 0.015 | 0.128 |