

1 **PhyloSuite: an integrated and scalable desktop platform for**
2 **streamlined molecular sequence data management and**
3 **evolutionary phylogenetics studies**

4

5 Dong Zhang^{1,2}, Fangluan Gao³, Wen X. Li¹, Ivan Jakovlić⁴, Hong Zou¹,

6 Jin Zhang⁴ and Gui T. Wang^{1,*}

7

8 ¹Key Laboratory of Aquaculture Disease Control, Ministry of Agriculture, and State
9 Key Laboratory of Freshwater Ecology and Biotechnology, Institute of Hydrobiology,
10 Chinese Academy of Sciences, Wuhan, P. R. China,

11 ²University of Chinese Academy of Sciences, Beijing, P. R. China,

12 ³Institute of Plant Virology, Fujian Agriculture and Forestry University, Fuzhou
13 350002, Fujian, P. R. China,

14 ⁴Bio-Transduction Lab, Wuhan Institute of Biotechnology, Wuhan, P. R. China

15

16 * **Corresponding author:** gtwang@ihb.ac.cn (GTW)

17

18 **Abstract**

19 Multi-gene and genomic datasets have become commonplace in the field of
20 phylogenetics, but many of the existing tools are not designed for such datasets,
21 which makes the analysis time-consuming and tedious. We therefore present
22 PhyloSuite, a user-friendly workflow desktop platform dedicated to streamlining
23 molecular sequence data management and evolutionary phylogenetics studies. It
24 employs a plugin-based system that integrates a number of useful phylogenetic and

25 bioinformatic tools, thereby streamlining the entire procedure, from data acquisition
26 to phylogenetic tree annotation, with the following features: (i) point-and-click and
27 drag-and-drop graphical user interface, (ii) a workspace to manage and organize
28 molecular sequence data and results of analyses, (iii) GenBank entries extraction and
29 comparative statistics, (iv) a phylogenetic workflow with batch processing capability,
30 (v) elaborate bioinformatic analysis for mitochondrial genomes. The aim of
31 PhyloSuite is to enable researchers to spend more time playing with scientific
32 questions, instead of wasting it on conducting standard analyses. The compiled binary
33 of PhyloSuite is available under the LGPL license at
34 <https://github.com/dongzhang0725/PhyloSuite/releases>, implemented in Python and
35 runs on Windows, Mac OSX and Linux.

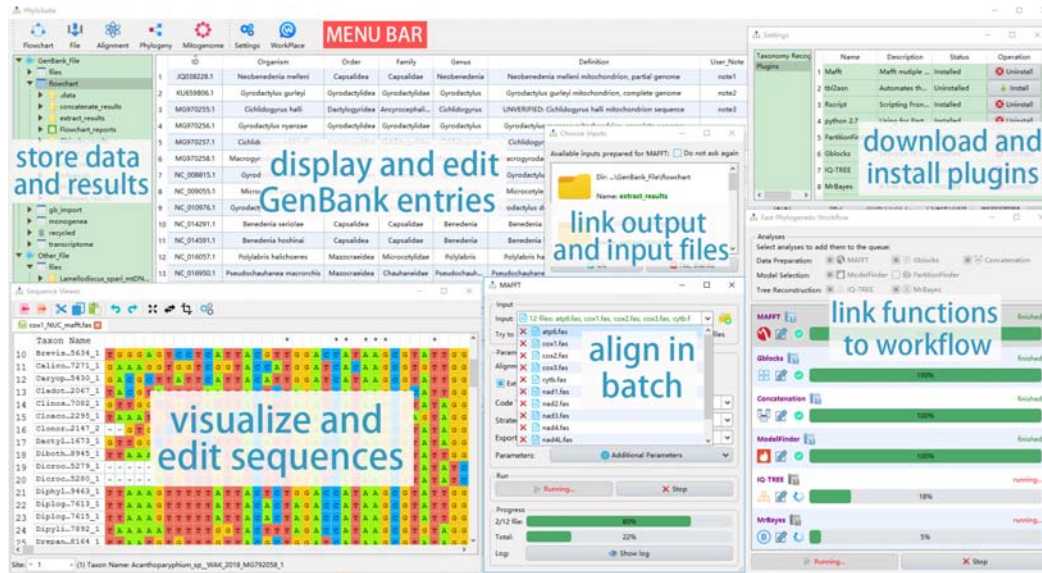
36

37 **Introduction**

38 Advancements in next-generation sequencing technologies (Metzker, 2009) have
39 resulted in a huge increase in the amount of genetic data available through public
40 databases. While this opens a multitude of research possibilities, retrieving and
41 managing such large amounts of data may be difficult and time-consuming for
42 researchers who are not computer-savvy. A standard analytical procedure for
43 phylogenetic analysis is: selecting and downloading GenBank entries, extracting
44 target genes (for multi-gene datasets, such as organelle genomes) and/or mining other
45 data, sequence alignment, alignment optimization, concatenation of alignments (for
46 multi-gene datasets), selection of best-fit partitioning schemes and evolutionary

47 models, phylogeny reconstruction, and finally visualization and annotation of the
48 phylogram. This can be very time-consuming if different programs have to be
49 employed for different steps, especially as they often have different input file format
50 requirements, and sometimes even require manual file tweaking. Therefore,
51 multifunctional, workflow-type software packages are becoming increasingly needed
52 by a broad range of evolutionary biologists (Smith, 2015). Specifically, as single-gene
53 datasets are rapidly being replaced by multi-gene or genomic datasets as a tool of
54 choice for phylogenetic reconstruction (Degnan and Rosenberg, 2009; Rivera-Rivera
55 and Montoya-Burgos, 2016), automated gene extraction from genomic data and batch
56 manipulation in some of the above steps, like alignment, are becoming a necessity.

57 Although there are several tools in existence, designed to streamline this process
58 by incorporating some or all of the steps mentioned above, none of these
59 incorporate all of the above functions in a manner suitable for current trends in
60 phylogenetic analyses (see detailed comparison in Supplementary data). Therefore,
61 we present PhyloSuite, a versatile tool designed to incorporate all of the functions
62 described above, including a series of different phylogenetic analysis algorithms, into
63 a single workflow that does not require programming skills, has an intuitive graphical
64 user interface (GUI), workspace, batch mode, extensive plugins support, inbuilt
65 updating function, etc. (Fig. 1). This tool aims to be accessible to all scientists,
66 streamline the phylogenetic analysis procedure, and allow scientists to focus on
67 solving scientific questions rather than waste time on toying with different scientific
68 software programs.



69

70 Fig. 1. The interface and the main functions of PhyloSuite

71 Implementation

72 PhyloSuite is a user-friendly stand-alone GUI-based software written in Python 3.6.7

73 and packaged and tested on Windows, Mac OSX and Linux. The functions are (Fig. 1,

74 Supplementary data): (i) retrieving, extracting, organizing and managing molecular

75 sequence data, including GenBank entries, nucleotide and amino acid sequences, and

76 sequences annotated in Word documents; (ii) batch alignment of sequences with

77 MAFFT (Katoh and Standley, 2013), for which we added a codon alignment

78 (translation align) mode; (iii) batch optimization of ambiguously aligned regions

79 using Gblocks (Talavera and Castresana, 2007); (iv) batch conversion of alignment

80 formats (FASTA, PHYLIP, PAML, AXT and NEXUS); (v) concatenation of multiple

81 alignments into a single dataset and preparation of a partition file for downstream

82 analyses; (vi) selection of the best-fit evolutionary model and/or partitioning scheme

83 using ModelFinder (Kalyaanamoorthy, et al., 2017) or PartitionFinder (Lanfear, et al.,

84 2017); (vii) phylogeny reconstruction using IQ-TREE (maximum likelihood) (Nguyen,

85 et al., 2015) and/or MrBayes (Bayesian inference) (Ronquist, et al., 2012); (viii)
86 linking the functions from (ii) to (vii) into a workflow; (ix) annotating phylogenetic
87 trees in the iTOL webtool (Letunic and Bork, 2016) using datasets generated by the (i)
88 function; (x) comprehensive bioinformatic analysis of mitochondrial genomes
89 (mitogenomes); (xi) visualization and editing of sequences using a MEGA-like
90 sequence viewer; (xii) storing, organizing and visualizing data and results of each
91 analysis in the PhyloSuite workspace.

92 **Genetic data management**

93 PhyloSuite provides a flexible GenBank entries processing function (see
94 Supplementary data). GenBank files can be imported either directly, or via a list of
95 IDs, which PhyloSuite will automatically download from the GenBank. Almost all of
96 the information in the annotation section of a GenBank record can be extracted and
97 displayed in the GUI. Additionally, the information can be standardized in batch using
98 a corresponding function or edited manually in the GUI, ambiguously annotated
99 mitogenomic tRNA genes can be semi-automatically reannotated using ARWEN
100 (Laslett and Canback, 2008), and taxonomic data can be automatically retrieved from
101 WoRMS and NCBI Taxonomy databases. The 'extract' function allows users to
102 extract genes in batches, as well as generate an assortment of statistics and dataset
103 files (iTOL datasets). The extracted results can be used for downstream analyses
104 without additional manipulation. The nucleotide and amino acid sequences can be
105 visualized and edited in a MEGA-like explorer equipped with common functions
106 (reverse complement, etc.). Importantly, PhyloSuite can parse the sequence

107 annotations recorded in a Word document via the inbuilt 'comment' function, and
108 generate a GenBank file and an *.sqn file for direct submission to the GenBank. This
109 function provides a novel and simple way to annotate genetic sequences, which shall
110 benefit researchers who are not computer-savvy.

111 **Phylogenetic analysis workflow**

112 By allowing users to combine seven plugin programs/functions and execute them
113 sequentially, PhyloSuite streamlines the evolutionary phylogenetics analysis (see
114 Supplementary data). The standard execution order of these functions is: MAFFT,
115 Gblocks, Concatenation, ModelFinder or PartitionFinder2, MrBayes and/or IQ-TREE.
116 The results of upstream functions are directly prepared as the input for downstream
117 functions, so only the first function of each workflow requires an input file(s).
118 Functions can also be used in a non-standard order and/or separately, in which case
119 PhyloSuite will automatically search for available input files (results of other tools) in
120 the workspace. Before starting the workflow, PhyloSuite will summarize the
121 parameters of each function, allowing the user to check and modify them, or
122 autocorrect conflicting parameters, such as sequence types. Once a workflow is
123 finished, PhyloSuite will describe the settings of each function as well as present the
124 references for each plugin program in the GUI.

125 **Bioinformatics analysis for mitogenomic data**

126 PhyloSuite was originally designed for, and its major comparative advantages are in,
127 the mitochondrial genomics analyses. There is a specialized configuration available
128 for the extraction of mitogenomic features. In addition to gene extraction, PhyloSuite

129 will generate a dozen of statistics and dataset files useful for downstream analyses
130 (see Supplementary data). The ‘itol’ dataset can be used to annotate the obtained
131 phylogram (colorize lineages, map gene orders, etc.). The gene order file can be used
132 to conduct gene order analysis with CREx (Bernt, et al., 2007) or treeREx (Bernt, et
133 al., 2008). The tables generated include the list of mitogenomes and overall statistics,
134 annotation, nucleotide composition and skewness, relative synonymous codon usage
135 (RSCU) and amino acid usage. The RSCU figure (see Fig. 3 in Zhang, et al. (2017))
136 can be drawn using the RSCU table and ‘Draw RSCU figure’ function. The
137 annotation table can be used to compare genomic annotations and calculate pairwise
138 similarity of homologous genes with ‘Compare table’ function (see Table 1 in Zhang,
139 et al. (2018)). In the future, PhyloSuite aims to gradually extend these analyses to
140 other small genomes (organelles, viruses, etc.).

141 **Discussion**

142 PhyloSuite links the management of genetic sequence data and a series of
143 phylogenetic analysis tools, thereby simplifying and speeding up multi-gene based
144 phylogenetic analyses, from data acquisition to phylogram annotation. In summary,
145 highlights of PhyloSuite include: (i) a user-friendly workspace to visualize, organize,
146 manipulate and store sequence data and results; (ii) flexible GenBank entries
147 processing (standardization, reannotation, etc.); (iii) batch data processing capability
148 and workflow; (iv) a state of the art mitogenomic bioinformatics analysis. Although
149 PhyloSuite is designed primarily to allow non-computer-savvy users to drag-and-drop
150 and point-and-click their way through the phylogenetic analysis, experienced

151 scientists will also find it useful to streamline their research, store and manage results,
152 and increase productivity. It will especially benefit evolutionary biologists who wish
153 to test the effects of different datasets and analytical methods on the phylogenetic
154 reconstruction.

155

156 **Acknowledgements**

157 The authors would like to thank Dr. Xiao-Qin Xia for modifying the manuscript and
158 Mr. Cheng-En Zheng for technical assistance in the software development.

159

160 **Funding**

161 This work was supported by the National Natural Science Foundation of China
162 [31872604]; the Earmarked Fund for China Agriculture Research System
163 [CARS-45-15]; and the Major Scientific and Technological Innovation Project of
164 Hubei Province [2015ABA045].

165

166 *Conflict of Interest:* none declared.

167

168

169 **References**


170 Bernt, M., Merkle, D. and Middendorf, M. (2008) An algorithm for inferring
171 mitogenome rearrangements in a phylogenetic tree. In Nelson, C.E. and Vialette, S.,
172 (eds.), *Comparative Genomics, International Workshop, RECOMB-CG 2008,*
173 *Proceedings of Lecture Notes in Bioinformatics.* Springer, Berlin, Vol. 5267, pp.
174 143-157.
175 Bernt, M., *et al.* (2007) CREx: inferring genomic rearrangements based on common

- 176 intervals. *Bioinformatics*, 23(21), 2957-2958.
- 177 Degnan, J.H. and Rosenberg, N.A. (2009) Gene tree discordance, phylogenetic
178 inference and the multispecies coalescent. *Trends Ecol Evol*, 24(6), 332-340.
- 179 Kalyaanamoorthy, S., *et al.* (2017) ModelFinder: fast model selection for accurate
180 phylogenetic estimates. *Nat Methods*, 14(6), 587-589.
- 181 Katoh, K. and Standley, D.M. (2013) MAFFT multiple sequence alignment software
182 version 7: improvements in performance and usability. *Molecular biology and
183 evolution*, 30(4), 772-780.
- 184 Lanfear, R., *et al.* (2017) PartitionFinder 2: new methods for selecting partitioned
185 models of evolution for molecular and morphological phylogenetic analyses. *Mol Biol
186 Evol*, 34(3), 772-773.
- 187 Laslett, D. and Canback, B. (2008) ARWEN: a program to detect tRNA genes in
188 metazoan mitochondrial nucleotide sequences. *Bioinformatics*, 24(2), 172-175.
- 189 Letunic, I. and Bork, P. (2016) Interactive tree of life (iTOL) v3: an online tool for the
190 display and annotation of phylogenetic and other trees. *Nucleic Acids Res*, 44(W1),
191 W242-245.
- 192 Metzker, M.L. (2009) Sequencing technologies — the next generation. *Nature
193 Reviews Genetics*, 11(1), 31-46.
- 194 Nguyen, L.T., *et al.* (2015) IQ-TREE: a fast and effective stochastic algorithm for
195 estimating maximum-likelihood phylogenies. *Mol Biol Evol*, 32(1), 268-274.
- 196 Rivera-Rivera, C.J. and Montoya-Burgos, J.I. (2016) LS(3): A Method for Improving
197 Phylogenomic Inferences When Evolutionary Rates Are Heterogeneous among Taxa.
198 *Mol Biol Evol*, 33(6), 1625-1634.
- 199 Ronquist, F., *et al.* (2012) MrBayes 3.2: efficient Bayesian phylogenetic inference and
200 model choice across a large model space. *Systematic biology*, 61(3), 539-542.
- 201 Smith, D.R. (2015) Buying in to bioinformatics: an introduction to commercial
202 sequence analysis software. *Brief Bioinform*, 16(4), 700-709.
- 203 Talavera, G. and Castresana, J. (2007) Improvement of phylogenies after removing
204 divergent and ambiguously aligned blocks from protein sequence alignments.
205 *Systematic Biology*, 56(4), 564-577.
- 206 Zhang, D., *et al.* (2018) Three new Diplozoidae mitogenomes expose unusual
207 compositional biases within the Monogenea class: implications for phylogenetic
208 studies. *BMC Evol Biol*, 18(1), 133.
- 209 Zhang, D., *et al.* (2017) Sequencing of the complete mitochondrial genome of a
210 fish-parasitic flatworm *Paratetraonchoides inermis* (Platyhelminthes: Monogenea):
211 tRNA gene arrangement reshuffling and implications for phylogeny. *Parasites &
212 Vectors*, 10(1), 462.

213


214

224 **Comparison with extant software programs**

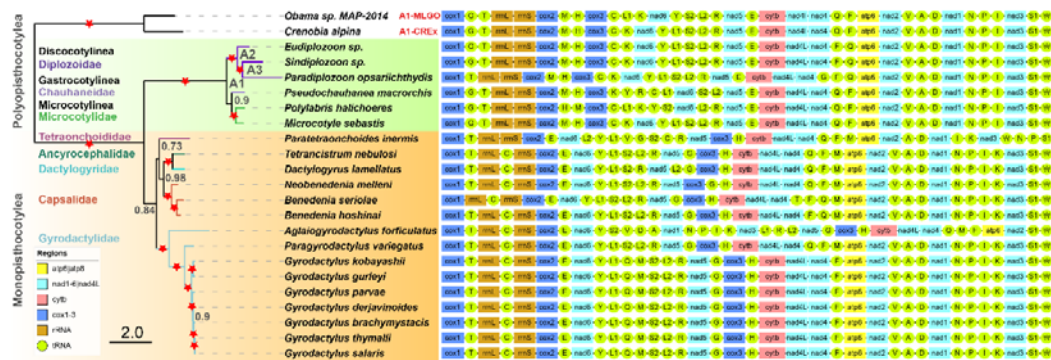
225 Although the extant software programs possess some of the abilities of PhyloSuite,
226 none of them incorporate all functions necessary for a streamlined multi-gene
227 phylogenetic analysis, from data retrieval to the phylogenetic tree annotation (click
228  here to access the attached Table S1). For example, FeatureExtract (Wernersson, 2005)
229 and Geneious (Kearse, et al., 2012) can extract the annotations from GenBank files,
230 but downstream analysis is not fully automated, so some manual data handling is
231 required, especially for multi-gene datasets. Armadillo (Lord, et al., 2012), EPoS
232 (Griebel, et al., 2008) and MEGA (Kumar, et al., 2016) do not possess the batch
233 processing capability, which is indispensable for multi-gene datasets. Additionally,
234 data partitioning and best-fit partitioning scheme estimation are also pivotal for
235 multi-gene dataset-based phylogenetic analyses (Blair and Murphy, 2011; Lanfear, et
236 al., 2012), but most other software programs lack this function, including Geneious,
237 MEGA, Galaxy Workflow (Oakley, et al., 2014), etc. Although, MEGA and EPoS
238 possess the ability to use the output of one tool directly as the input for another tool,
239 they cannot link several functions into a single run (workflow). Probably the closest to
240 meeting the described requirements is Geneious, but this is a commercial
241 bioinformatics software, so it may not be an ideal option (i.e. too expensive) for all
242 scientists, especially for students.

243 **Functions and capabilities**

244 Taking a recently published mitogenomic paper (Zhu, et al., 2018) as an example,
245 using the 'extract' function user can quickly conduct most of the analyses reported in

246 that paper, and generate similar tables and figures: (i) mitogenome list and overall
 247 statistics table (Table 1 in Zhu et al.), (ii) annotation table (Table 2), (iii) nucleotide
 248 composition and skewness table (Table 3), (iv) relative synonymous codon usage
 249 (RSCU) table and figure (Fig. 2B), (v) amino acid usage statistics file used to draw
 250 Fig. 2A, and (vi) reconstruct and annotate (in iTOL) phylogenetic trees (Fig. 5) using
 251 the extracted genes. In comparison, most of the tables in that paper were made
 252 manually by the author, which is time-consuming, tedious and error-prone. Beyond
 253 these, several additional analyses are available: (i) gene order file is generated, which
 254 can be used to map gene orders of mitogenomes onto the phylograms in iTOL (Fig.
 255 S2, also see Fig. 6 in Zhang, et al. (2018)), (ii) statistics for individual genes,
 256 including size, start and terminal codons, base composition and skews (click here to
 257  access the attached Table S2), (iii) general statistics table, which can be used to draw
 258 skewness and base content figure (Fig. S3, also see Fig. 1 in Zhang, et al. (2018)), (iv)
 259 comparison of genomic annotations and pairwise similarity calculation for
 260 homologous genes, such as Table 1 in Zhang, et al. (2018).

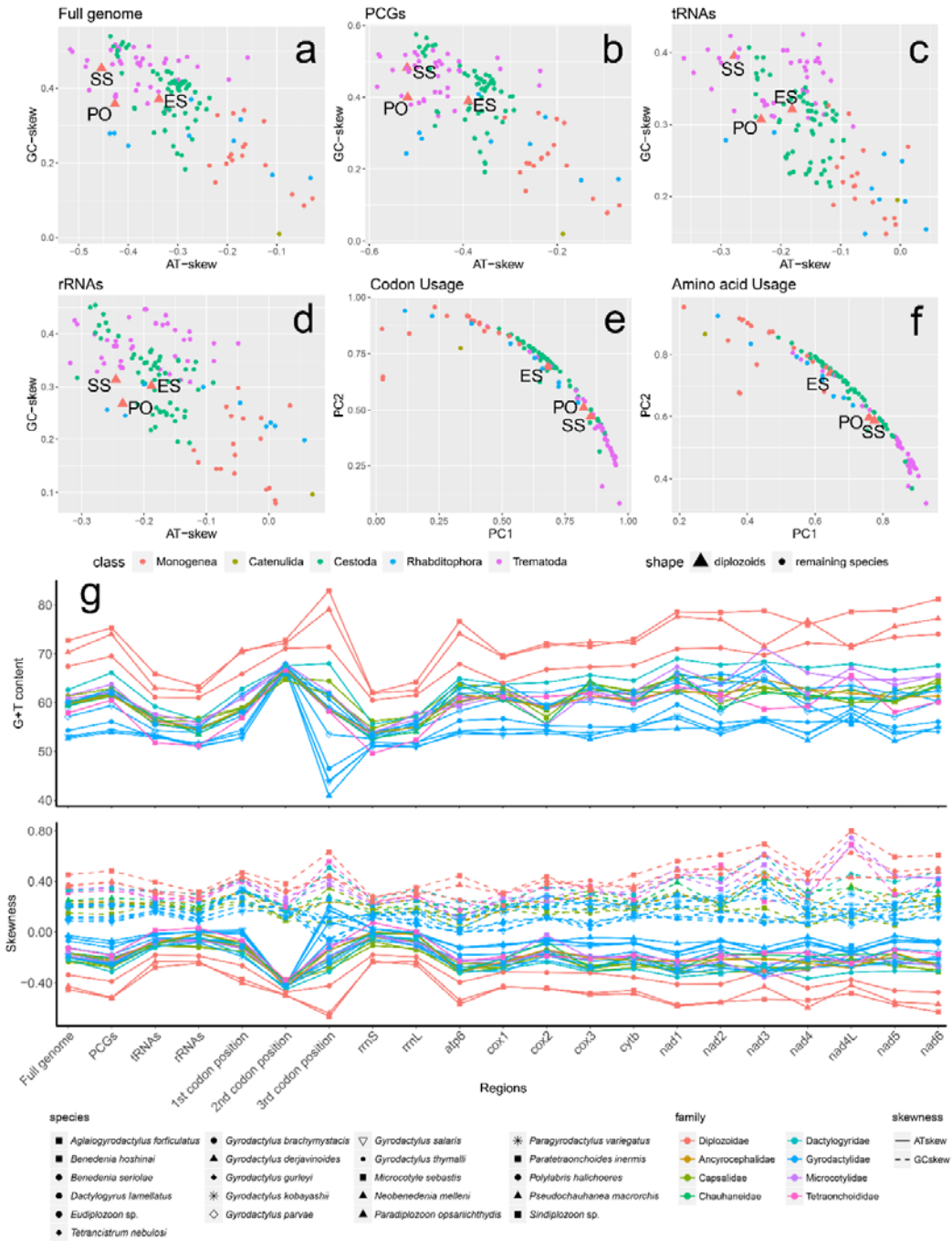
261



262

263 Fig. S2 Mapping gene orders of monogenean mitogenomes onto the phylogenetic

264 tree.



265

266 Fig. S3 Skewness and base content of some flatworm mitogenomes.

267

268

269 References

270 Blair, C. and Murphy, R.W. (2011) Recent trends in molecular phylogenetic analysis:

- 271 where to next? *J Hered*, 102(1), 130-138.
- 272 Griebel, T., Brinkmeyer, M. and Bocker, S. (2008) EPoS: a modular software
273 framework for phylogenetic analysis. *Bioinformatics*, 24(20), 2399-2400.
- 274 Kearse, M., *et al.* (2012) Geneious Basic: an integrated and extendable desktop
275 software platform for the organization and analysis of sequence data. *Bioinformatics*,
276 28(12), 1647-1649.
- 277 Kumar, S., Stecher, G. and Tamura, K. (2016) MEGA7: Molecular Evolutionary
278 Genetics Analysis Version 7.0 for Bigger Datasets. *Mol Biol Evol*, 33(7), 1870-1874.
- 279 Lanfear, R., *et al.* (2012) PartitionFinder: combined selection of partitioning schemes
280 and substitution models for phylogenetic analyses. *Molecular biology and evolution*,
281 29(6), 1695-1701.
- 282 Lord, E., *et al.* (2012) Armadillo 1.1: an original workflow platform for designing and
283 conducting phylogenetic analysis and simulations. *PLoS One*, 7(1), e29903.
- 284 Oakley, T.H., *et al.* (2014) Osiris: accessible and reproducible phylogenetic and
285 phylogenomic analyses within the Galaxy workflow management system. *BMC*
286 *Bioinformatics*, 15(1), 230.
- 287 Wernersson, R. (2005) FeatureExtract--extraction of sequence annotation made easy.
288 *Nucleic Acids Res*, 33(Web Server issue), W567-569.
- 289 Zhang, D., *et al.* (2018) Three new Diplozoidae mitogenomes expose unusual
290 compositional biases within the Monogenea class: implications for phylogenetic
291 studies. *BMC Evol Biol*, 18(1), 133.
- 292 Zhu, H.F., *et al.* (2018) Complete mitochondrial genome of the crab spider *Ebrechtella*
293 *tricuspidata* (Araneae: Thomisidae): A novel tRNA rearrangement and phylogenetic
294 implications for Araneae. *Genomics*.