

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

4

5 Dong Zhang<sup>1,2</sup>, Fangluan Gao<sup>3</sup>, Wen X. Li<sup>1</sup>, Ivan Jakovlić<sup>4</sup>, Hong Zou<sup>1</sup>,  
6 Jin Zhang<sup>4</sup> and Gui T. Wang<sup>1,\*</sup>

7

8 <sup>1</sup>Key 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 <sup>2</sup>University of Chinese Academy of Sciences, Beijing, P. R. China,

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

14 <sup>4</sup>Bio-Transduction Lab, 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 GPL license at  
34 <https://github.com/dongzhang0725/PhyloSuite/releases>, implemented in Python and  
35 runs on Windows, Mac OSX and Linux.

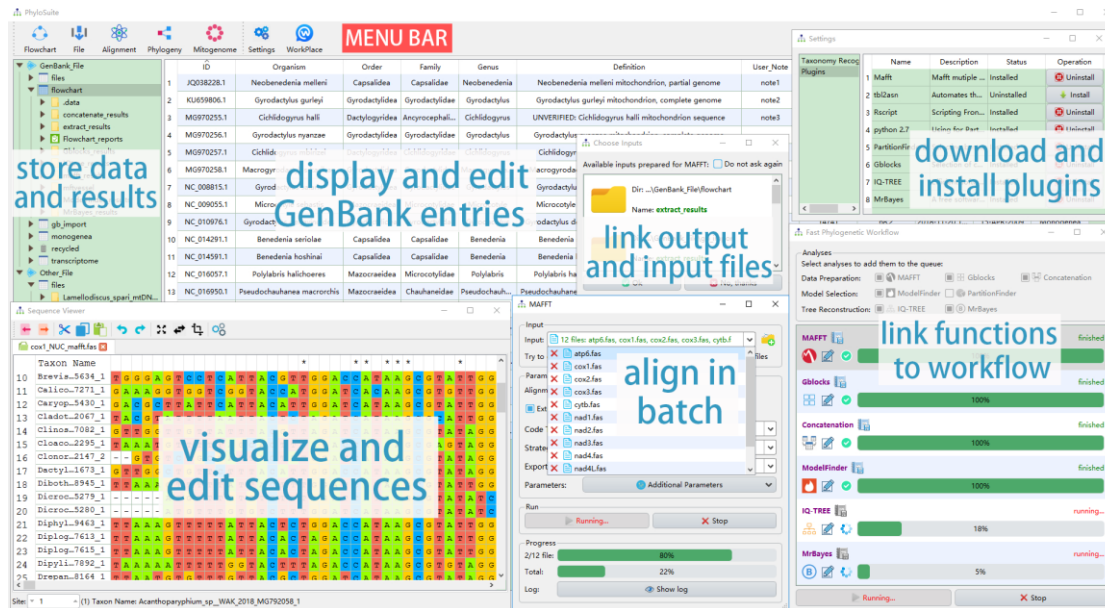
36

## 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.

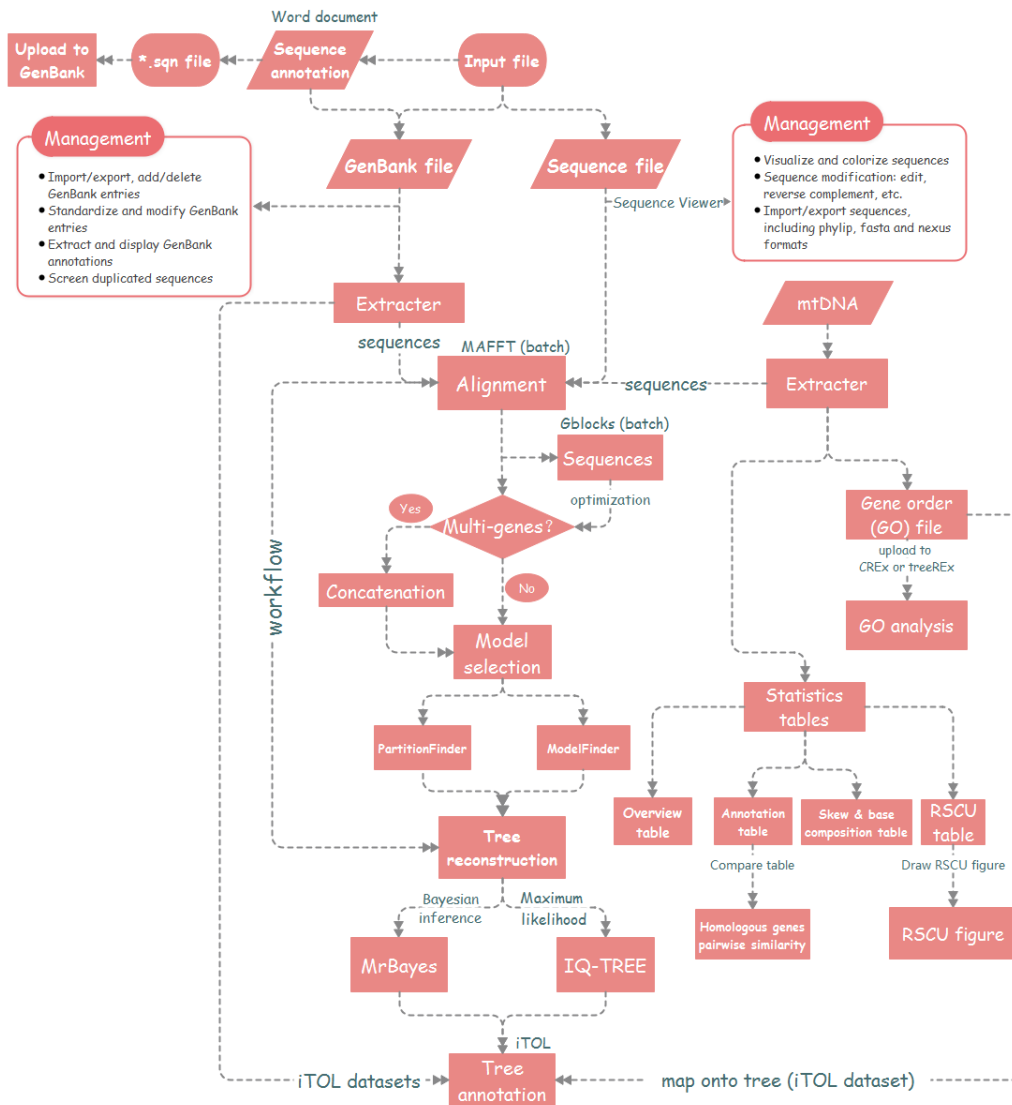
213

214

## 215 Supplementary Information

### 216 Overview

217 PhyloSuite is designed to address the global trend towards multi-gene based  
218 phylogenetic analyses (Degnan and Rosenberg, 2009; Rivera-Rivera and  
219 Montoya-Burgos, 2016): this software program incorporates and streamlines all steps  
220 included in such analyses, from data acquisition to phylogenetic tree visualization and  
221 annotation (Fig. S1).



222

223 Fig. S1. The workflow diagram of PhyloSuite.

## 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 (Fig. S3).  
 228 For example, FeatureExtract (Wernersson, 2005) and Geneious (Kearse, et al., 2012)  
 229 can extract the annotations from GenBank files, but downstream analysis is not fully  
 230 automated, so some manual data handling is required, especially for multi-gene  
 231 datasets. Armadillo (Lord, et al., 2012), EPoS (Griebel, et al., 2008) and MEGA  
 232 (Kumar, et al., 2016) do not possess the batch processing capability, which is  
 233 indispensable for multi-gene datasets. Additionally, data partitioning and best-fit  
 234 partitioning scheme estimation are also pivotal for multi-gene dataset-based  
 235 phylogenetic analyses (Blair and Murphy, 2011; Lanfear, et al., 2012), but most other  
 236 software programs lack this function, including Geneious, MEGA, Galaxy Workflow  
 237 (Oakley, et al., 2014), etc. Although, MEGA and EPoS possess the ability to use the  
 238 output of one tool directly as the input for another tool, they cannot link several  
 239 functions into a single run (workflow). Probably the closest to meeting the described  
 240 requirements is Geneious, but this is a commercial bioinformatics software, so it may  
 241 not be an ideal option (i.e. too expensive) for all scientists, especially for students.

	PhyloSuite	MitoPhAST	HomBlocks	phylogeny.fr	Galaxy Workf	Armadillo	phyloGenerato	Geneious	MEGA	EPoS
Graphical interface/Webpage	✓	×	×	✓	✓	✓	×	✓	✓	✓
Alignment	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Codon alignment	✓	✓	×	×	×	×	×	✓	✓	×
Alignment optimization	✓	✓	✓	✓	✓	×	✓	×	×	✓
Model selection	✓	✓	✓	✓	✓	✓	×	×	✓	×
Maximum likelihood tree	✓	✓	×	✓	✓	✓	✓	✓	✓	✓
Bayesian tree	✓	×	×	✓	✓	✓	✓	✓	×	✓
Data partitioning	✓	✓	✓	×	✓	×	×	✓	×	✓
Best partitioning scheme estimate	✓	✓	✓	×	×	×	×	×	×	×
Batch processing	✓	✓	✓	×	✓	×	×	✓	×	×
Tree annotation	✓	×	×	✓	×	✓	×	✓	✓	✓
Workflow	✓	✓	✓	✓	✓	✓	✓	✓	✓ <sup>a</sup>	✓ <sup>a</sup>
Free	✓	✓	✓	✓	✓	✓	✓	×	✓	✓

242  
 243 Fig. S3 Comparison of PhyloSuite with software programs with similar functions.

244 <sup>a</sup> semi-workflow.

## 245 Functions and capabilities

246 Taking a recently published mitogenomic paper (Zhu, et al., 2018) as an example,  
 247 using the ‘extract’ function user can quickly conduct most of the analyses reported in  
 248 that paper, and generate similar tables and figures: (i) mitogenome list and overall  
 249 statistics table (Fig. S4, Table 1 in Zhu et al.), (ii) annotation table (Fig. S5, Table 2 in  
 250 Zhu et al.), (iii) nucleotide composition and skewness table (Fig. S6, Table 3 in Zhu et  
 251 al.), (iv) relative synonymous codon usage (RSCU) table (Fig. S7) and figure (Fig. S8,  
 252 Fig. 2B in Zhu et al.), (v) amino acid usage statistics file (Fig. S9) used to draw Fig.  
 253 2A in Zhu et al., and (vi) reconstruct and annotate (in iTOL) phylogenetic trees (Fig. 5  
 254 in Zhu et al.) using the extracted genes. In comparison, most of the tables in that paper  
 255 were made manually by the author, which is time-consuming, tedious and error-prone.  
 256 Beyond these, several additional analyses are available: (i) gene order file is generated,  
 257 which can be used to map gene orders of mitogenomes onto the phylograms in iTOL  
 258 (Fig. S10, also see Fig. 6 in Zhang, et al. (2018)) and conduct gene order analysis with  
 259 CREx (Bernt, et al., 2007) or treeREx (Bernt, et al., 2008), (ii) statistics for individual  
 260 genes, including size, start and terminal codons, base composition and skews (Fig.  
 261 S11), (iii) general statistics table, which can be used to draw skewness and base  
 262 content figure (Fig. S12, also see Fig. 1 in Zhang, et al. (2018)), (iv) comparison of  
 263 genomic annotations and pairwise similarity calculation for homologous genes (Fig.  
 264 S13, Table 1 in Zhang, et al. (2018)).

265

Taxon	Accession number	Whole Genome				PCGs			
		Size(bp)	AT%	AT-Skew	GC-Skew	Size(bp)	AT%	AT-Skew	GC-Skew
Araneae									
Arachnida									
Agelenidae									
N/A									
Agelena silvatica	NC 033971.1	14776	74.5	-0.163	0.302	10642	73.6	-0.171	0.103
Dipluridae									
Phyxioschema suthepium	NC 020322.1	13931	67.4	-0.04	0.472	10730	66.6	-0.232	0.153
Hypochilidae									
Hypochilus thorelli	NC 010777.1	13991	70.3	-0.14	0.266	10753	69	-0.192	0.064
Liphistiidae									
Liphistius erawan	NC 020323.1	14197	67.7	0.024	-0.361	10794	66.8	-0.165	-0.09
Songthela hangzhouensis	NC 005924.1	14215	72.2	-0.023	-0.235	10765	71.5	-0.16	-0.022
Nemesiidae									
Calisoga longitarsis	NC 010780.1	14070	64	-0.146	0.365	10738	63.1	-0.257	0.121
Pholcidae									
Pholcus phalangioides	NC 020324.1	14459	65.8	-0.191	0.371	10631	65.4	-0.175	0.078
Pholcus sp. HCP-2014	KJ782458.1	14279	65.8	-0.188	0.372	10631	64.9	-0.174	0.072
Salticidae									
Carrhotus xanthogramma	KP402247.1	14563	75.1	-0.089	0.26	10809	74.1	-0.147	0.057
Selenopidae									
Selenops bursarius	NC 024878.1	14272	74.4	-0.123	0.321	10756	73.8	-0.15	0.056
Theraphosidae									
Haplopetma schmidti	NC 005925.1	13874	69.8	-0.083	0.344	10724	69.8	-0.181	0.092
Thomisidae									
Ebrechtella tricuspadata	KU852748.1	14532	76.2	-0.097	0.221	5552	74.6	-0.143	-0.067
Oxytate striatipes	NC 025557.1	14407	78.2	-0.084	0.212	10784	77.7	-0.179	0.093

266

267 Fig. S4 Mitogenome list and overall statistics table.

Gene	Position		Size	Intergenic nucleotide	Codon		Strand
	From	To			Start	Stop	
cox3	4	789	786	3	TTG	TAA	H
trnG	788	853	66	-2			H
nad3	837	1173	337	-17	ATT	T	H
trnL	1156	1218	63	-18			L
trnN	1214	1283	70	-5			H
trnA	1279	1335	57	-5			H
trnS1	1327	1381	55	-9			H
trnR	1377	1428	52	-5			H
trnE	1419	1479	61	-10			H
trnF	1460	1513	54	-20			L
nad5	1514	3142	1629		ATA	TAA	L
trnH	3147	3205	59	4			L
nad4	3207	4486	1280	1	ATA	TA	L
trnP	4755	4808	54	268			L
nad6	4817	5245	429	8	ATG	TAA	H
trnI	5247	5302	56	1			H

268

269 Fig. S5 Annotation table.

Regions	Size (bp)	T(U)	C	A	G	AT(%)	GC(%)	GT(%)	AT skew	GC skew
PCGs	9810	42.9	12.8	30.1	14.1	73	26.9	57	-0.175	0.048
1st codon position	3270	38.8	11.6	32.8	16.8	71.6	28.4	55.6	-0.084	0.183
2nd codon position	3270	44.1	15.8	23.8	16.3	67.9	32.1	60.4	-0.299	0.014
3rd codon position	3270	45.8	11.1	33.8	9.3	79.6	20.4	55.1	-0.151	-0.087
atp6	663	41	11.6	31.7	15.7	72.7	27.3	56.7	-0.129	0.149
atp8	159	36.5	9.4	41.5	12.6	78	22	49.1	0.065	0.143
cox1	1536	41.3	12.4	27.5	18.8	68.8	31.2	60.1	-0.2	0.205
cox3	786	42.1	10.7	28.4	18.8	70.5	29.5	60.9	-0.195	0.276
cytb	1134	43.5	12	28.2	16.3	71.7	28.3	59.8	-0.213	0.153
nad1	903	44.5	15.3	29.3	10.9	73.8	26.2	55.4	-0.205	-0.169
nad2	954	39.9	7.3	35.4	17.3	75.3	24.6	57.2	-0.06	0.404
nad3	337	49.3	11.9	28.8	10.1	78.1	22	59.4	-0.262	-0.081
nad4	1280	44.1	15.9	30.7	9.2	74.8	25.1	53.3	-0.18	-0.267
nad5	1629	44.9	17.4	28.5	9.1	73.4	26.5	54	-0.224	-0.312
nad6	429	40.3	5.1	36.4	18.2	76.7	23.3	58.5	-0.052	0.56
rrnL	1007	41.9	12.4	34.7	11	76.6	23.4	52.9	-0.095	-0.059
rrnS	698	41.4	10.9	34.4	13.3	75.8	24.2	54.7	-0.093	0.101
tRNAs	1222	38.3	9.7	38.6	13.4	76.9	23.1	51.7	0.004	0.163
rRNAs	1705	41.7	11.8	34.5	12	76.2	23.8	53.7	-0.094	0.007
Full genome	14344	36.6	10.5	37.1	15.8	73.7	26.3	52.4	0.007	0.205

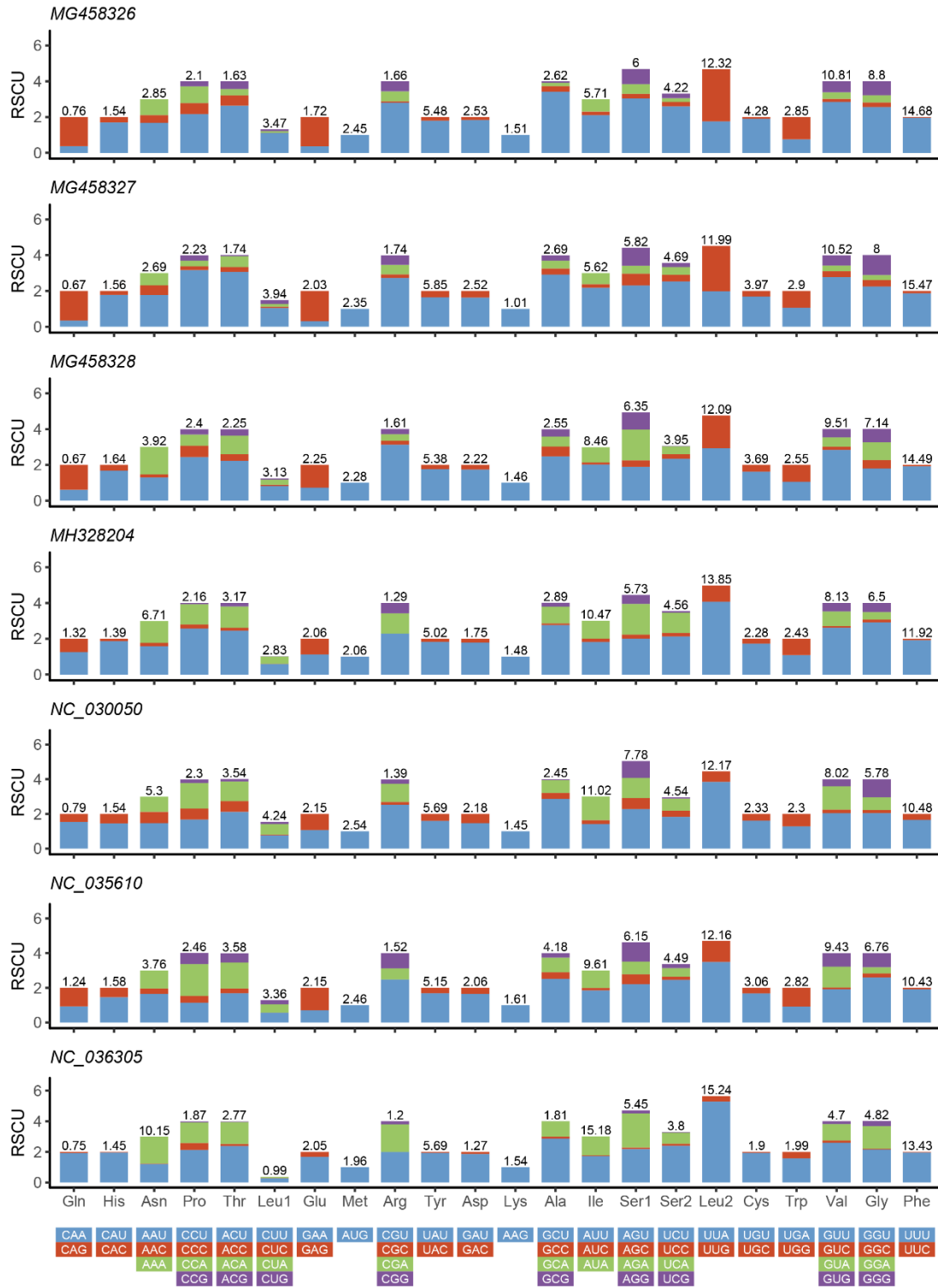
270

271 Fig. S6 Nucleotide composition and skewness table.

Sequences used: <i>Cyrtophora moluccensis</i>											
Codon Table: 5											
Domain: Data											
Codon	Count	RSCU	Codon	Count	RSCU	Codon	Count	RSCU	Codon	Count	RSCU
UUU (F)	316	1.83	UCU (S)	112	2.47	UAU (Y)	113	1.65	UGU (C)	21	1.62
UUC (F)	30	0.17	UCC (S)	28	0.62	UAC (Y)	24	0.35	UGC (C)	5	0.38
UUA (L)	273	3.81	UCA (S)	88	1.94	UAA (*)	7	1.56	UGA (W)	75	1.76
UUG (L)	36	0.5	UCG (S)	6	0.13	UAG (*)	2	0.44	UGG (W)	10	0.24
CUU (L)	65	0.91	CCU (P)	75	2.54	CAU (H)	53	1.74	CGU (R)	5	0.45
CUC (L)	7	0.1	CCC (P)	13	0.44	CAC (H)	8	0.26	CGC (R)	4	0.36
CUA (L)	43	0.6	CCA (P)	24	0.81	CAA (Q)	42	1.68	CGA (R)	29	2.64
CUG (L)	6	0.08	CCG (P)	6	0.2	CAG (Q)	8	0.32	CGG (R)	6	0.55
AUU (I)	297	1.73	ACU (T)	54	1.86	AAU (N)	103	1.6	AGU (S)	13	0.29
AUC (I)	47	0.27	ACC (T)	11	0.38	AAC (N)	26	0.4	AGC (S)	9	0.2
AUA (M)	217	1.72	ACA (T)	45	1.55	AAA (K)	73	1.51	AGA (S)	86	1.9
AUG (M)	36	0.28	ACG (T)	6	0.21	AAG (K)	24	0.49	AGG (S)	21	0.46
GUU (V)	67	1.3	GCU (A)	96	2.65	GAU (D)	50	1.79	GGU (G)	41	0.89
GUC (V)	15	0.29	GCC (A)	9	0.25	GAC (D)	6	0.21	GGC (G)	11	0.24
GUA (V)	96	1.86	GCA (A)	38	1.05	GAA (E)	54	1.54	GGA (G)	87	1.89
GUG (V)	28	0.54	GCG (A)	2	0.06	GAG (E)	16	0.46	GGG (G)	45	0.98

272

273 Fig. S7 Relative synonymous codon usage (RSCU) table.



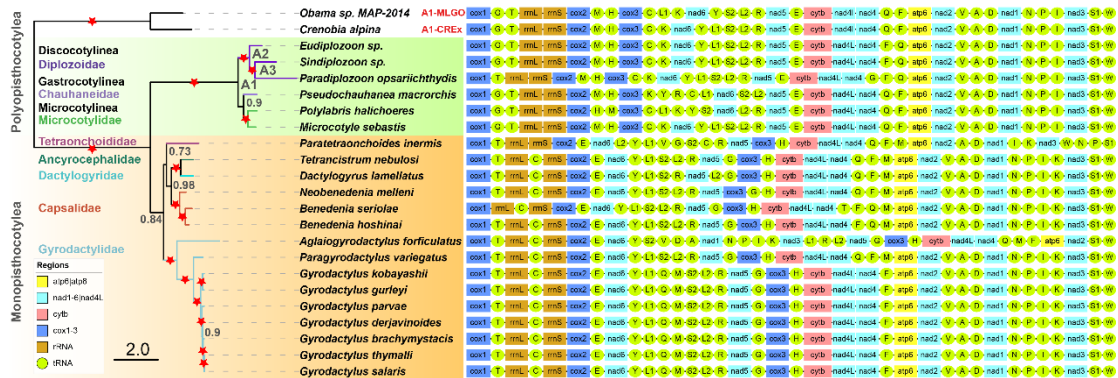
274

275 Fig. S8 Relative synonymous codon usage (RSCU) of seven species.

AA	Count	%
Phe (F)	346	10.61
Leu2 (L2)	309	9.48
Leu1 (L1)	121	3.71
Ile (I)	344	10.55
Met (M)	253	7.76
Val (V)	206	6.32
Ser2 (S2)	234	7.18
Pro (P)	118	3.62
Thr (T)	116	3.56
Ala (A)	145	4.45
Tyr (Y)	137	4.2
His (H)	61	1.87
Gln (Q)	50	1.53
Asn (N)	129	3.96
Lys (K)	97	2.98
Asp (D)	56	1.72
Glu (E)	70	2.15
Cys (C)	26	0.8
Trp (W)	85	2.61
Arg (R)	44	1.35
Ser1 (S1)	129	3.96
Gly (G)	184	5.64
codon end in A or T	2758	84.6
codon end in G or T	1739	53.34

276

277 Fig. S9 Amino acid usage statistics file.



278

279 Fig. S10 Mapping gene orders of monogenean mitogenomes onto the phylogenetic

280 tree. The figure was published in Zhang, et al. (2018).

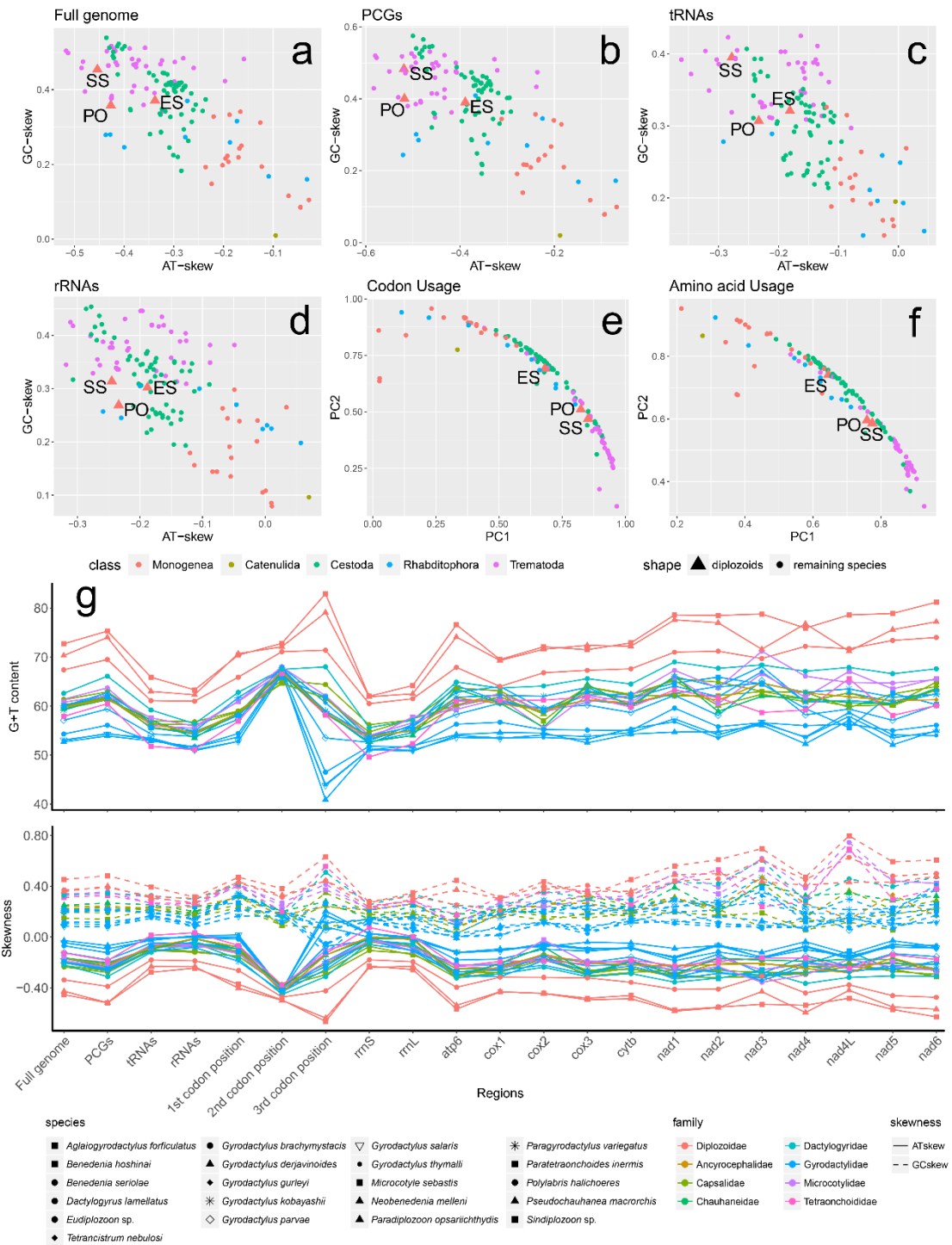
Species	P_sp	C_x	C_a	C_n	H_p	N_n	N_a	S_h	H_s	H_o
Length of PCGs (bp)										
atp6	663	663	663	663	663	663	663	661	669	660
atp8	150	153	153	162	156	159	159	150	153	153
cox1	1536	1539	1536	1533	1536	1536	1536	1533	1536	1542
cox2	640	664	666	666	666	666	666	667	669	666
Length of rRNA genes (bp)										
rrnL	1002	999	1026	1043	1028	1023	1022	1119	1048	1018
rrnS	691	693	680	696	696	699	700	698	666	691
Putative start codon										
atp6	ATG	ATA	ATA	ATA	ATA	ATA	ATA	ATG	ATG	ATT
atp8	ATT	ATT	ATT	ATT	ATA	ATT	ATT	ATA	ATG	ATT
cox1	ATA	TTA	TTA	CTG	TTA	TTA	TTA	TTG	TTA	ATT
cox2	GTG	TTG	TTG	TTG	TTG	TTG	TTG	CTG	ATG	TTG
Putative terminal codon										
atp6	TAA	TAA	TAA	TAA	TAA	TAA	TAA	T	TAG	TAA
atp8	TAA	TAA	TAA	TAA	TAA	TAA	TAA	TAA	TAA	TAG
cox1	TAG	TAA	TAA	TAA	TAA	TAA	TAA	TAA	TAG	TAA
cox2	T	T	TAA	TAA	TAA	TAA	TAA	T	TAG	TAG
AT skew										
atp6	-0.317	-0.195	-0.148	-0.171	-0.266	-0.173	-0.229	-0.122	-0.261	-0.21
atp8	-0.258	-0.023	0.083	0.015	-0.023	0.143	0.071	-0.135	0.053	-0.042
cox1	-0.316	-0.234	-0.205	-0.207	-0.222	-0.173	-0.207	-0.163	-0.255	-0.233
cox2	-0.233	-0.239	-0.117	-0.11	-0.157	-0.097	-0.149	-0.082	-0.127	-0.187
rRNAs	0.145	0.041	-0.023	0	0.04	-0.005	-0.015	0.034	0.022	0.087
rrnL	0.145	0.037	-0.027	0.002	0.005	-0.016	0.001	0.054	-0.021	0.076
rrnS	0.144	0.048	-0.018	-0.004	0.094	0.011	-0.039	0	0.089	0.103
tRNAs	0.006	0.009	0.04	0.007	0.036	0.01	0.003	-0.029	0.076	0.029
GC skew										
atp6	0.312	0.255	0.129	0.193	0.231	0.133	0.18	-0.284	0.295	0.264
atp8	0.472	0.4	0.212	0.308	-0.04	0.308	0.438	-0.59	0.487	0.515
cox1	0.293	0.206	0.165	0.202	0.245	0.193	0.216	-0.072	0.271	0.244
cox2	0.365	0.268	0.19	0.23	0.189	0.205	0.259	-0.196	0.384	0.347
rRNAs	0.145	0.041	-0.023	0	0.04	-0.005	-0.015	0.034	0.022	0.087
rrnL	-0.154	-0.039	-0.013	-0.005	-0.02	0	-0.041	0.26	-0.11	-0.161
rrnS	-0.101	0.015	0.126	0.076	0.032	0.047	0.018	0.175	-0.125	-0.09
tRNAs	0.006	0.009	0.04	0.007	0.036	0.01	0.003	-0.029	0.076	0.029
AT content										
atp6	66.7	75.7	74.4	75	76.5	78.4	73.1	73.3	71.1	73.6
atp8	64.7	86.9	78.4	84	83.9	83.6	79.9	74	74.5	77.8
cox1	61.5	69.4	69.6	68.4	69.7	72.1	68.6	67.5	66.4	69.7
cox2	63.6	72.5	70.7	73.9	73.7	75.8	71.6	67.9	68.5	70.6
rRNAs	71	80.1	76.2	80.6	79.6	80	77.7	74.3	70.3	78.2
rrnL	70.7	79.3	76.7	80.2	80.8	80.9	78.3	75.6	69.8	78.6
rrnS	71.3	81.3	75.5	81.1	77.7	78.7	76.8	72.2	71.1	77.4
tRNAs	66.8	74.4	76.6	77.2	76.9	78.6	75.9	73.3	70.1	73.4
GC content										
rRNAs	29	19.9	23.8	19.5	20.5	20	22.3	25.7	29.7	21.8
rrnL	29.3	20.6	23.3	19.9	19.2	19.2	21.7	24.4	30.3	21.3
rrnS	28.7	18.7	24.5	19	22.3	21.4	23.3	27.8	28.8	22.6
tRNAs	33.2	25.6	23.5	22.8	23.1	21.4	24.1	26.8	29.9	26.6

281

282 Fig. S11 Statistics for individual genes, including size, start and terminal codons, base

283 composition and skews. Only 4 protein-coding genes (PCGs) are shown.





284

285 Fig. S12 Skewness and base content of some flatworm mitogenomes. The figure was

286 published in Zhang, et al. (2018).

Gene	Position		Size	Strand	Identity
	From	To			
NC_005924/NC_005925/KU852748					
cox1	1/1/1217	1533/1536/2752	1533/1536/1536	H/H/H	70.18/73.96/76.17/73.44
cox2	1537/1537/2756	2203/2205/3427	667/669/672	H/H/H	60.75/64.14/66.96/63.95
trnK	2204/2204/3428	2265/2263/3486	62/60/59	H/H/H	50.00/49.21/62.30/53.83
trnD	2264/2247/3472	2318/2301/3522	55/55/51	H/H/H	48.28/61.82/53.57/54.56
atp8	2319/2293/3532	2468/2445/3675	150/153/144	H/H/H	51.88/42.61/48.41/47.63
atp6	2462/2439/-	3122/3107/-	661/669/-	H/H/-	55.46/-/-/55.46
cox3	3124/3108/-	3897/3891/-	774/784/-	H/H/-	60.84/-/-/60.84
trnG	3909/3892/5126	3962/3945/5180	54/54/55	H/H/H	50.91/57.89/62.50/57.10
nad3	3969/3945/-	4298/4266/-	330/322/-	H/H/-	52.54/-/-/52.54
trnA	4298/4349/5608	4350/4405/5665	53/57/58	H/H/H	41.38/47.62/50.00/46.33
trnR	4357/4452/5718	4413/4501/5770	57/50/53	H/H/H	26.87/50.88/58.49/45.41
trnN	4413/4316/5562	4471/4363/5617	59/48/56	H/H/H	51.67/47.62/44.64/47.98
trnS1	4472/4396/5670	4523/4459/5721	52/64/52	H/H/H	42.19/44.23/51.56/45.99
trnE	4524/4493/5759	4583/4543/5811	60/51/53	H/H/H	49.18/58.33/56.60/54.71
trnF	4572/4531/5797	4626/4583/5855	55/53/59	L/L/L	58.18/46.97/54.24/53.13
nad5	4628/4584/5856	6263/6222/7479	1636/1639/1624	L/L/L	51.61/54.32/58.86/54.93
trnH	6262/6211/7492	6318/6268/7546	57/58/55	L/L/L	49.18/50.88/51.67/50.57
nad4	6319/6267/7547	7618/7544/8822	1300/1278/1276	L/L/L	53.40/55.84/58.40/55.88
nad4L	7618/7545/8791	7899/7821/9090	282/277/300	L/L/L	46.82/47.44/57.00/50.42
trnT	7898/9478/10812	7962/9540/10866	65/63/55	H/H/H	56.72/55.38/69.84/60.65
trnP	7959/7815/9083	8019/7867/9140	61/53/58	L/L/L	55.74/47.62/58.62/53.99
nad6	8031/7860/-	8462/8288/-	432/429/-	H/H/-	48.44/-/-/48.44
cytb	8469/8289/-	9590/9423/-	1122/1135/-	H/H/-	60.09/-/-/60.09
trnS2	9589/9425/10758	9651/9478/10811	63/54/54	H/H/H	34.33/53.97/48.33/45.54

287

288 Fig. S13 Comparison of genomic annotations and pairwise similarity calculation for  
289 homologous genes.

290

## 291 Usage examples

292 PhyloSuite has been used previously to conduct analyses in a number of published  
293 papers. You may refer to the following publications for examples of the use of  
294 PhyloSuite: (Hua, et al., 2018; Li, et al., 2018; Li, et al., 2017; Liu, et al., 2017; Liu, et  
295 al., 2018; Liu, et al., 2018; Wang, et al., 2017; Wen, et al., 2017; Xi, et al., 2018;  
296 Zhang, et al., 2018; Zhang, et al., 2018; Zhang, et al., 2018; Zhang, et al., 2017;  
297 Zhang, et al., 2017; Zou, et al., 2017; Zou, et al., 2018). Note that we have merged our  
298 two older beta tools, MitoTool (<https://github.com/dongzhang0725/MitoTool>) and  
299 BioSuite (<https://github.com/dongzhang0725/BioSuite>) into PhyloSuite, so some of  
300 our older published papers may refer these two tools instead.

301

## 302 **References**

- 303 Bernt, M., Merkle, D. and Middendorf, M. (2008) An algorithm for inferring mitogenome  
304 rearrangements in a phylogenetic tree. In Nelson, C.E. and Vialette, S., (eds.), *Comparative Genomics,*  
305 *International Workshop, RECOMB-CG 2008, Proceedings of Lecture Notes in Bioinformatics.* Springer,  
306 Berlin, Vol. 5267, pp. 143-157.
- 307 Bernt, M., *et al.* (2007) CREx: inferring genomic rearrangements based on common intervals.  
308 *Bioinformatics*, 23(21), 2957-2958.
- 309 Blair, C. and Murphy, R.W. (2011) Recent trends in molecular phylogenetic analysis: where to next? *J.*  
310 *Hered.*, 102(1), 130-138.
- 311 Degnan, J.H. and Rosenberg, N.A. (2009) Gene tree discordance, phylogenetic inference and the  
312 multispecies coalescent. *Trends Ecol. Evol.*, 24(6), 332-340.
- 313 Griebel, T., Brinkmeyer, M. and Bocker, S. (2008) EPoS: a modular software framework for  
314 phylogenetic analysis. *Bioinformatics*, 24(20), 2399-2400.
- 315 Hua, C.J., *et al.* (2018) Basal position of two new complete mitochondrial genomes of parasitic  
316 Cymothoidea (Crustacea: Isopoda) challenges the monophyly of the suborder and phylogeny of the  
317 entire order. *Parasit. Vectors*, 11(1), 628.
- 318 Kearse, M., *et al.* (2012) Geneious Basic: an integrated and extendable desktop software platform for  
319 the organization and analysis of sequence data. *Bioinformatics*, 28(12), 1647-1649.
- 320 Kumar, S., Stecher, G. and Tamura, K. (2016) MEGA7: Molecular Evolutionary Genetics Analysis Version  
321 7.0 for Bigger Datasets. *Mol. Biol. Evol.*, 33(7), 1870-1874.
- 322 Lanfear, R., *et al.* (2012) PartitionFinder: combined selection of partitioning schemes and substitution  
323 models for phylogenetic analyses. *Mol. Biol. Evol.*, 29(6), 1695-1701.
- 324 Li, W.X., *et al.* (2018) Comparative mitogenomics supports synonymy of the genera *Ligula* and  
325 *Digramma* (Cestoda: Diphylobothriidae). *Parasit. Vectors*, 11(1), 324.
- 326 Li, W.X., *et al.* (2017) The complete mitochondrial DNA of three monozoic tapeworms in the  
327 Caryophyllidea: a mitogenomic perspective on the phylogeny of eucestodes. *Parasit. Vectors*, 10(1),  
328 314.
- 329 Liu, F.-F., *et al.* (2017) Tandem duplication of two tRNA genes in the mitochondrial genome of *Tagiades*  
330 *vajuna* (Lepidoptera: Hesperidae). *Eur. J. Entomol.*, 114, 407-415.
- 331 Liu, Q., *et al.* (2018) Genetic diversity of glacier-inhabiting *Cryobacterium* bacteria in China and  
332 description of *Cryobacterium zongtaii* sp. nov. and *Arthrobacter glacialis* sp. nov. *Syst. Appl. Microbiol.*
- 333 Liu, Z.Q., *et al.* (2018) Sequencing of complete mitochondrial genomes confirms synonymization of  
334 *Hyalomma asiaticum asiaticum* and *kozlovi*, and advances phylogenetic hypotheses for the Ixodidae.  
335 *PLoS One*, 13(5), e0197524.
- 336 Lord, E., *et al.* (2012) Armadillo 1.1: an original workflow platform for designing and conducting  
337 phylogenetic analysis and simulations. *PLoS One*, 7(1), e29903.
- 338 Oakley, T.H., *et al.* (2014) Osiris: accessible and reproducible phylogenetic and phylogenomic analyses  
339 within the Galaxy workflow management system. *BMC Bioinformatics*, 15(1), 230.
- 340 Rivera-Rivera, C.J. and Montoya-Burgos, J.I. (2016) LS(3): A Method for Improving Phylogenomic  
341 Inferences When Evolutionary Rates Are Heterogeneous among Taxa. *Mol. Biol. Evol.*, 33(6),  
342 1625-1634.
- 343 Wang, J.G., *et al.* (2017) Sequencing of the complete mitochondrial genomes of eight freshwater snail

344 species exposes pervasive paraphyly within the Viviparidae family (Caenogastropoda). *PLoS One*, 12(7),  
345 e0181699.

346 Wen, H.B., *et al.* (2017) The complete maternally and paternally inherited mitochondrial genomes of a  
347 freshwater mussel *Potamilus alatus* (Bivalvia: Unionidae). *PLoS one*, 12(1), e0169749.

348 Wernersson, R. (2005) FeatureExtract--extraction of sequence annotation made easy. *Nucleic Acids*  
349 *Res.*, 33(Web Server issue), W567-569.

350 Xi, B.W., *et al.* (2018) Characterization of the complete mitochondrial genome of  
351 *Parabreviscolexniepini* Xi et al., 2018 (Cestoda, Caryophyllidea). *Zookeys*, (783), 97-112.

352 Zhang, D., *et al.* (2018) Mitochondrial genomes of two diplectanids (Platyhelminthes: Monogenea)  
353 expose paraphyly of the order Dactylogyridea and extensive tRNA gene rearrangements. *Parasit.*  
354 *Vectors*, 11(1), 601.

355 Zhang, D., *et al.* (2018) Homoplasy or plesiomorphy? Reconstruction of the evolutionary history of  
356 mitochondrial gene order rearrangements in the subphylum Neodermata (under review). *Int. J.*  
357 *Parasitol.*

358 Zhang, D., *et al.* (2018) Three new Diplozoidae mitogenomes expose unusual compositional biases  
359 within the Monogenea class: implications for phylogenetic studies. *BMC Evol. Biol.*, 18(1), 133.

360 Zhang, D., *et al.* (2017) Sequencing, characterization and phylogenomics of the complete  
361 mitochondrial genome of *Dactylogyrus lamellatus* (Monogenea: Dactylogyridae). *J. Helminthol.*, 1-12.

362 Zhang, D., *et al.* (2017) Sequencing of the complete mitochondrial genome of a fish-parasitic flatworm  
363 *Paratetraonchoides inermis* (Platyhelminthes: Monogenea): tRNA gene arrangement reshuffling and  
364 implications for phylogeny. *Parasit. Vectors*, 10(1), 462.

365 Zhu, H.F., *et al.* (2018) Complete mitochondrial genome of the crab spider *Ebrechtella tricuspidata*  
366 (Araneae: Thomisidae): A novel tRNA rearrangement and phylogenetic implications for Araneae.  
367 *Genomics*.

368 Zou, H., *et al.* (2017) The complete mitochondrial genome of parasitic nematode *Camallanus cotti*:  
369 extreme discontinuity in the rate of mitogenomic architecture evolution within the Chromadorea class.  
370 *BMC Genomics*, 18(1), 840.

371 Zou, H., *et al.* (2018) The complete mitochondrial genome of *Cymothoa indica* has a highly rearranged  
372 gene order and clusters at the very base of the Isopoda clade. *PLoS One*, 13(9), e0203089.

373

374

375

376