

Chromosome-level assembly of the mustache toad genome using third-generation DNA sequencing and Hi-C analysis

--Manuscript Draft--

Manuscript Number:	GIGA-D-19-00099R2	
Full Title:	Chromosome-level assembly of the mustache toad genome using third-generation DNA sequencing and Hi-C analysis	
Article Type:	Data Note	
Funding Information:	National Key Research and Development Program of China (2017YFC0505202)	Dr. Dingqi Rao
	National Natural Science Foundation of China (30270175)	Dr. Dingqi Rao
	National Natural Science Foundation of China (30870278)	Dr. Dingqi Rao
	National Natural Science Foundation of China (31372165)	Dr. Dingqi Rao
Abstract:	<p>Background The mustache toad, <i>Vibrissaphora ailaonica</i>, is endemic to China and belongs to the Megophryidae family. Like other mustache toad species, <i>V. ailaonica</i> males temporarily develop keratinized nuptial spines on their upper jaw during each breeding season, which fall off at the end of the breeding season. This feature is likely to be result of the reversal of sexual dimorphism in body size, with males being larger than females. A high-quality reference genome for the mustache toad would be invaluable to investigate the genetic mechanism underlying these repeatedly developing keratinized spines.</p> <p>Findings To construct the mustache toad genome, we generated 225 Gb of short reads and 277 Gb of long reads using Illumina and Pacific Biosciences (PacBio) sequencing technologies, respectively. Sequencing data were assembled into a 3.53-Gb genome assembly, with a contig N50 length of 821 Kb. We also used high-throughput chromosome conformation capture (Hi-C) technology to identify contacts between contigs, then assembled contigs into scaffolds and assembled a genome with 13 chromosomes and a scaffold N50 length of 412.42 Mb. Based on the 26,227 protein-coding genes annotated in the genome, we analyzed phylogenetic relationships between the mustache toad and other chordate species. The mustache toad has a relatively higher evolutionary rate and separated from a common ancestor of the marine toad, bullfrog, and Tibetan frog 206.1 million years ago. Furthermore, we identified 201 expanded gene families in the mustache toad, which were mainly enriched in immune pathway, keratin filament, and metabolic processes.</p> <p>Conclusions Using Illumina, PacBio, and Hi-C technologies, we constructed the first high-quality chromosome-level mustache toad genome. This work not only offers a valuable reference genome for functional studies of mustache toad traits, but also provides important chromosomal information for wider genome comparisons.</p>	
Corresponding Author:	Dingqi Rao, Ph.D Kunming Institute of Zoology Chinese Academy of Sciences Kunming, Yunnan CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Kunming Institute of Zoology Chinese Academy of Sciences	
Corresponding Author's Secondary		

Institution:	
First Author:	Yongxin Li
First Author Secondary Information:	
Order of Authors:	Yongxin Li
	Yandong Ren
	Dongru Zhang
	Hui Jiang
	Zhongkai Wang
	Xueyan Li
	Dingqi Rao, Ph.D
Order of Authors Secondary Information:	
Response to Reviewers:	All the comments in this version have been corrected.
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
<p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
Resources	Yes
<p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum</p>	

Standards Reporting Checklist?	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	Yes

[Click here to view linked References](#)

Li, Ren, Zhang et al.

1 **Chromosome-level assembly of the mustache toad genome using third-generation**

2 **DNA sequencing and Hi-C analysis**

3 Running title: The mustache toad genome

4 Yongxin Li^{1,2,‡}, Yandong Ren^{2,‡}, Dongru Zhang^{1,‡}, Hui Jiang³, Zhongkai Wang²,

5 Xueyan Li¹, Dingqi Rao^{1,*}

6

7 ¹ State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of

8 Zoology, Chinese Academy of Sciences, Kunming 650223, China; ² Center for

9 Ecological and Environmental Sciences, Northwestern Polytechnical University,

10 Xi'an 710072, China; ³ National Engineering Laboratory of Marine Germplasm

11 Resources Exploration and Utilization, Zhejiang Ocean University, Zhoushan 316022,

12 China

13

14 [‡]Equal contribution

15 ^{*}Correspondence address. Dingqi Rao. State Key Laboratory of Genetic Resources

16 and Evolution, Kunming Institute of Zoology, Chinese Academy of Sciences,

17 Kunming 650223, China. Tel: +86-0871-65128516; E-mail: kizar@mail.kiz.ac.cn;

18 ORCID: 0000-0003-2680-1503

19

20 **Abstract**

21 **Background:** The mustache toad, *Vibrissaphora ailaonica*, is endemic to China and

22 belongs to the Megophryidae family. Like other mustache toad species, *V. ailaonica*

23 males temporarily develop keratinized nuptial spines on their upper jaw during each
24 breeding season, which fall off at the end of the breeding season. This feature is likely
25 to be result of the reversal of sexual dimorphism in body size, with males being larger
26 than females. A high-quality reference genome for the mustache toad would be
27 invaluable to investigate the genetic mechanism underlying these repeatedly
28 developing keratinized spines. **Findings:** To construct the mustache toad genome, we
29 generated 225 Gb of short reads and 277 Gb of long reads using Illumina and Pacific
30 Biosciences (PacBio) sequencing technologies, respectively. Sequencing data were
31 assembled into a 3.53-Gb genome assembly, with a contig N50 length of 821 Kb. We
32 also used high-throughput chromosome conformation capture (Hi-C) technology to
33 identify contacts between contigs, then assembled contigs into scaffolds and
34 assembled a genome with 13 chromosomes and a scaffold N50 length of 412.42 Mb.
35 Based on the 26,227 protein-coding genes annotated in the genome, we analyzed
36 phylogenetic relationships between the mustache toad and other chordate species. The
37 mustache toad has a relatively higher evolutionary rate and separated from a common
38 ancestor of the marine toad, bullfrog, and Tibetan frog 206.1 million years ago.
39 Furthermore, we identified 201 expanded gene families in the mustache toad, which
40 were mainly enriched in immune pathway, keratin filament, and metabolic processes.
41 **Conclusions:** Using Illumina, PacBio, and Hi-C technologies, we constructed the first
42 high-quality chromosome-level mustache toad genome. This work not only offers a
43 valuable reference genome for functional studies of mustache toad traits, but also
44 provides important chromosomal information for wider genome comparisons.

45

46 **Keywords:** Mustache toad; Genome assembly; Evolution; PacBio; Hi-C

47

48 **Data Description**

49 The mustache toad, *Vibrissaphora ailaonica* (NCBI: txid:428466), is an amphibian
50 belonging to the Megophryidae family that is endemic to China (including the
51 China–Vietnam border) [1–3]. This mustache toad species exhibits many interesting
52 features, including unique keratinized spines along the upper jaw [1, 4–6]. These
53 spines grow repeatedly in sexually mature males during the breeding season, and fall
54 off at the end of this process [5–8] (Figure 1). This morphological difference between
55 males and females is further highlighted by their sexual dimorphism in body size
56 (males are significantly larger than females). The spines (and body size) may be used
57 as a weapon for sexually mature individuals to compete for nests and mating
58 opportunities [7, 9, 10]. Another unique aspect of the mustache toad is that breeding
59 occurs during the cold season, whereas most frogs and toads breed in the warmer
60 months [1]. However, despite the importance of the mustache toad in terms of
61 dynamic spine development and sexual dimorphism in body size, few genomic
62 resources exist for this species. In fact, to date, no next-generation sequencing (NGS)
63 data have been reported in the *Vibrissaphora* genus. The lack of genome sequence and
64 transcriptome data for *V. ailaonica* has hindered identification of functional genes
65 related to their attractive and dynamic appearance (e.g., spine and body size). The
66 shortage of amphibian genomes represented in the Genome 10K project makes it

67 necessary to analyze other important genomes to study phylogenetic relationships in
68 amphibians on a larger scale [11].

69 In this study, we combined genomic sequencing data from Illumina short reads,
70 PacBio long reads, and Hi-C data, to generate the first chromosome-level reference
71 genome for the mustache toad. The completeness and continuity of the genome were
72 comparable with that of other important amphibian species. The high-quality
73 reference genome generated in this study will facilitate research on population genetic
74 traits and functional gene identification related to important characteristics of the
75 mustache toad.

76

77 **Analyses and Methods**

78 **Sampling and sequencing**

79 During the breeding season (in February), a male mustache toad (*V. ailaonica*) with
80 keratinized nuptial spines on its upper jaw was caught for sequencing from Ailao
81 Mountain (Figure 1). To obtain sufficient high-quality DNA for the PacBio Sequel
82 platform (Pacific Biosciences, USA), the mustache toad was dissected, and fresh liver
83 tissue was used for DNA extraction using phenol/chloroform extraction. DNA quality
84 was checked by agarose gel electrophoresis, and high integrity DNA molecules were
85 obtained. Other tissues, including spines, brain, stomach, intestine, liver, lung, spleen,
86 blood, and tongue, were snap-frozen in liquid nitrogen for 10 min. These 9
87 organs/tissues were stored at -80°C for RNA-seq analysis. Isolated total RNA was
88 used to isolate intact poly (A) + RNA using the NEBnext Ultra-Directional RNA

89 Library Prep kit (NEB, protocol B) for library construction. The mRNA was further
90 fragmented and randomly primed during first-strand synthesis by reverse transcription.
91 This procedure was followed by second-strand synthesis with DNA polymerase I to
92 create double-stranded cDNA fragments using Transcriptor First Strand cDNA
93 Synthesis Kit (Roche).

94 For the Hi-C experiments, collected blood was used for library construction. The
95 blood sample (150 μ l) was cross-linked for 10 min with formaldehyde (1% final
96 concentration), after which glycine (0.2 M final concentration) was added for 5 min to
97 stop the cross-linking process. The sample was then stored until required for further
98 analysis.

99 Extracted DNA was sequenced using the Illumina and PacBio Sequel platforms. Short
100 reads generated from the Illumina platform were used to estimate genome size and to
101 correct errors in the assembled genome, and the PacBio long reads were used for
102 genome assembly. To this end, five libraries with insertion lengths of 220 bp or
103 500 bp were sequenced on an Illumina HiSeq 2500 platform, generating 150-bp
104 paired-end reads. A 20-Kb library was constructed using the PacBio platform,
105 according to the manufacturer's protocols. Finally, we obtained 225.03 Gb of Illumina
106 short reads and 277.15 Gb of PacBio long reads (Table 1, Additional Table S1,
107 Additional Table S2). The average N50 length of subreads was 14.78 Kb, providing
108 ultra-long genomic sequences for the following assembly and analysis (Additional
109 Table S2).

110 RNA-seq samples were obtained by mixing an equal amount of RNA extracted from

111 each tissue that had been stored and used for library construction. After sequencing on
112 the Illumina HiSeq 4000 platform, we obtained 14.18 Gb of sequencing data (Table 1,
113 Additional Table S3). Four Hi-C libraries were constructed using the same sample
114 with same parameters, and sequenced on the Illumina HiSeq X-ten platform, which
115 generated 378.78 Gb of clean data (Table 1, Additional Table S4).

116

117 **Genome characteristics estimation**

118 Illumina short reads were filtered for quality as follows. First, adaptors were removed
119 from the sequencing reads. Then, read pairs were excluded if any one read had more
120 than 10% 'N', and read pairs with more than 50% low quality bases were removed.

121 Finally, PCR duplicates produced during library construction were removed.

122 Filtered reads were used to estimate genome size and other characteristics. Using the
123 *k*-mer method, we calculated the 17-mer depth frequency distribution in the mustache
124 toad. Genome size was estimated by:

$$125 \quad G = \text{TKN}_{17\text{-mer}} / \text{PKFD}_{17\text{-mer}}$$

126 where $\text{TKN}_{17\text{-mer}}$ is the total *k*-mer number, and $\text{PKFD}_{17\text{-mer}}$ is the peak *k*-mer
127 frequency depth of 17-mer.

128 We estimated a genome size of 3.52 Gb (peak = 54) and found heterozygous, repeated
129 sequence peaks, suggesting that the mustache toad genome exhibits complex genome
130 assembly (Figure 2).

131

132 **Genome assembly using PacBio long reads and Hi-C data**

133 Based on 38 single-molecule real-time cells, and using the PacBio Sequel platform,
134 we generated 277.15 Gb of subreads (Table 1, Additional Table S2). The average and
135 N50 length of subreads was 9.65 Kb and 14 78 Kb, respectively (Additional Table S2).
136 All long reads were assembled using wtdbg software [12] (WTDBG, RRID:
137 SCR_017225). As a result, we obtained a 3.95-Gb genome assembly, with a contig
138 N50 length of 739.54 Kb. However, although the size of the genome assembly was
139 comparable with the estimated *k*-mer result, the end result was a slightly larger. This
140 may be associated with the complexity of the mustache toad genome (which has a
141 high rate of heterozygosity rate and repetitive sequences). Redundancy in the genome
142 assembly was removed using Redundans software (v0.13c) [13], with an identity of
143 0.7 and overlap of 0.7. This resulted in a genome assembly of 3.58 Gb and a contig
144 N50 length of 834.90 Kb. To ensure that all contigs removed were not real sequences,
145 we used BUSCO (Benchmarking Universal Single-Copy Orthologs) [14] and the
146 mapping ratio of Illumina reads in both the raw genome and the redundancy-filtered
147 genome. Results of these checks indicated that the parameters used in the
148 redundancy-filtered step were appropriate for this study (Additional Table S5). To
149 further improve the quality and accuracy of our genome assembly, Illumina short
150 reads were used to polish the genome using Pilon software (Pilon, RRID:
151 SCR_014731, v1.21) [15] at the single-base level.

152 Hi-C data were used to improve the connection integrity of the contigs (15,899
153 contigs). We obtained 378.78 Gb of Hi-C sequencing data, which was first filtered
154 using Hic-Pro (v2.10.0) [16] (Table 1, Additional Table S4), and then mapped to the

155 polished mustache toad genome [17]. The locations and directions of the contigs were
156 determined by 3D *de novo* assembly (3d-DNA) software (v180419) [18], with default
157 parameters. Most contigs were then successfully clustered and anchored in 13 groups
158 (Figure 3) [19]. Finally, we obtained the first chromosome-level, high quality
159 mustache toad assembly (3.53 Gb) with a scaffold N50 length of 412.42 Mb, which
160 provides a solid genomic resource to assist further study of the mustache toad (Table
161 2).

162

163 **Genome assembly evaluation**

164 The quality of a genome assembly is directly related to the accuracy and completeness
165 of protein-coding gene prediction. Therefore, we evaluated the assembled mustache
166 toad genome using three methods. First, the assembled genome was compared against
167 the core gene set in BUSCO (BUSCO, RRID:SCR_015008, v2.0) [14]. We found 245
168 (80.8%) and 833 (85.1%) conserved core genes in the mustache toad genome using
169 the eukaryote and metazoan databases, respectively (Table 3). When we further
170 considered the fragmented BUSCO genes found in the genome, there were 272
171 (89.7%) and 881 (90.1%) conserved core genes in the eukaryote and metazoan
172 databases, respectively. These results indicated that the assembled mustache toad
173 genome is comparable with published amphibian genomes (Table 3).

174 Second, all filtered short reads generated from the Illumina platform were aligned to
175 the genome using BWA (Burrows–Wheeler Aligner) software (BWA,
176 RRID:SCR_010910, v0.7.12) [20]; 1,778 million clean reads could be mapped to the

177 genome, accounting for 97.78% of total clean reads (Additional Table S6).
178 Third, RNA-seq reads were *de novo*-assembled using Bridger software (Bridger,
179 RRID:SCR_017039, version: r2014-12-01) [21], with redundant transcripts removed
180 by TGICL [22]. This resulted in 19,876 transcripts (Additional Table S7). These
181 transcripts were then aligned to the genome, with 17,878 transcripts (89.95%) found
182 in the assembled genome, and 94.52% of transcripts being longer than 1 Kb
183 (Additional Table S8). Analysis of N50 length and BUSCO results revealed that the
184 mustache toad genome was comparable with that of other published amphibian
185 genomes (Tables 2–4), indicating that our assembled mustache toad genome exhibited
186 high completeness and accuracy.
187 The GC distribution of the mustache toad genome, and that of other vertebrate species,
188 was calculated using the slide window method. GC distributions were similar, with an
189 average GC content of 43.68% in the mustache toad, and 36.60% to 44.49% in other
190 species (Additional Figure S1).

191

192 **Genome annotation**

193 Tandem Repeat Finder (TRF, v4.04) [23] was used to identify repetitive elements, and
194 RepeatModeler software (RepeatModeler, RRID:SCR_015027, v1.0.4) was used to
195 detect transposable elements (TEs) in the mustache toad genome. Then, the *de novo*
196 library of repeats produced by RepeatModeler analysis and the rebase
197 (RepBase16.02) database were used for RepeatMasker (RepeatMasker,
198 RRID:SCR_012954, version: open-4.0) [24] analysis to identify homologous repeats.

199 RepeatProteinMask was used to query the TE protein database at the protein level.
200 Lastly, we identified 2.45 Gb of repeat sequences, accounting for 69.48% of the
201 estimated genome size (Additional Table S9). Among these repeat sequences, 60.87%
202 (2.15 Gb) was predicted by the *de novo* method (Table 5).
203 After repeat sequence annotation, we masked all repeats, except for the tandem repeat
204 sequences, for protein-coding gene annotation. Augustus software (Augustus,
205 RRID:SCR_008417, v2.5.5) [25] was used to *de novo*-predict coding genes using a
206 zebrafish (*Danio rerio*) dataset as the training species. For the homology-based
207 method, protein sequences of chordate species, including *D. rerio* (GCF_000002035.6)
208 [26], *Nanorana parkeri* (GCF_000935625.1) [27], *Homo sapiens*
209 (GCF_000001405.38) [28], *Gallus gallus* (GCF_000002315.5) [29], *Pelodiscus*
210 *sinensis* (GCF_000230535.1) [30], *Xenopus laevis* (GCF_001663975.1) [31], and
211 *Petromyzon marinus* [32], were downloaded and aligned against the mustache toad
212 genome using the TBLASTN module (TBLASTN, RRID:SCR_011822, BLAST
213 version: 2.3.0). The transcripts assembled by RNA-seq reads were first translated into
214 amino acids and then aligned to the genome using TBLASTN software for gene
215 annotation. EVidenceModeler (EVidenceModeler, RRID:SCR_014659, version:
216 r2012-06-25) [33] was used to integrate results from the three methods, and genes
217 with poor transcriptome evidence support were filtered out. Finally, 26,227
218 high-quality protein-coding genes were predicted in the mustache toad genome. The
219 distributions of mRNA, coding sequences, exon and intron lengths were comparable
220 with those of closely related species (Figure 4).

221 Gene functional annotation can help to elucidate gene function. Thus, we aligned all
222 26,227 protein-coding genes to protein databases, including InterProScan, Kyoto
223 Encyclopedia of Genes and Genomes (KEGG), SwissProt, and TrEMBL. Results
224 showed that most of the genes obtained could be annotated from these functional
225 databases (Table 6).

226

227 **Phylogenetic tree and divergence time analysis**

228 To reveal phylogenetic relationships between the mustache toad and other closely
229 related species, we identified the single-copy genes among these species. First, protein
230 sequences, including those of *D. rerio* (GCF_000002035.6) [26], *N. parkeri*
231 (GCF_000935625.1) [27], *H. sapiens* (GCF_000001405.38) [28], *G. gallus*
232 (GCF_000002315.5) [29], *Anolis carolinensis* (GCF_000090745.1) [34], *Xenopus*
233 *tropicalis* (GCF_000004195.3) [31], *Rhinella marina* (GigaDB) [35], *Rana*
234 *catesbeiana* (GCA_002284835.2) [36], *Ambystoma mexicanum* [37, 38], and
235 *Alligator sinensis* (GCF_000455745.1) [39], were downloaded from the National
236 Center for Biotechnology Information (NCBI). The longest transcript of each gene in
237 each species was selected. BLASTP (BLASTP, RRID:SCR_001010, BLAST version:
238 2.2.24) was then used to align these protein sequences from the 11 species (including
239 the mustache toad), with an e-value of 1e-5. Homology relationships (including
240 orthologs and paralogs) were then determined using OrthoMCL software (v1.4) [40].
241 Genes with only one copy in the species were identified as single-copy genes. In total,
242 238 genes were identified (Figure 5). Detailed statistics about gene families are shown

243 in Additional Table S10.

244 The 238 single-copy genes were aligned using MUSCLE software (MUSCLE,
245 RRID:SCR_011812, v3.8.31) [41, 42], and concatenated to supergenes for
246 maximum-likelihood-based phylogenetic analyses. We performed phylogenetic
247 analysis, with zebrafish as the outgroup, using RAxML software (RAxML,
248 RRID:SCR_006086, v8.2.3) [43], with the parameter ‘-m’ for PROTGAMMAAUTO.
249 Results indicated that the mustache toad has a close relationship with the ancestor of
250 the marine toad (*R. marina*), bullfrog (*R. catesbeiana*), and Tibetan frog (*N. parkeri*),
251 with topological relationships in other clades found to be the same as reported
252 previously (Figure 6). To further investigate the divergence time of these species,
253 especially toad and frogs, the MCMCTREE model (part of PAML software package;
254 PAML, RRID:SCR_014932, v4.8) [44] was used with three datasets (four-fold
255 degenerate sites [4dTVs], first-codon sites, and second-codon sites) extracted from the
256 single-copy genes as the input file. Fossil records were downloaded from the
257 TIMETREE website [45] and used to calibrate the results. Results from the three
258 different datasets were very similar, showing that the mustache toad diverged from the
259 common ancestor of the marine toad, bullfrog and Tibetan frog about 206.1 million
260 years ago (Figure 6; Additional Figure S2, Additional Figure S3).

261

262 **Gene family expansion and contraction**

263 We performed gene family expansion and contraction analysis using CAFÉ software
264 (CAFÉ, RRID:SCR_005983, v4.0) [46], and found 201 and 326 expanded and

265 contracted gene families in the mustache toad ($P < 0.05$), respectively. Using the Gene
266 Ontology (GO) and KEGG databases, functional enrichment analysis of expanded
267 gene families revealed 210 GO terms (adjusted $P < 0.05$) and 9 KEGG pathways
268 ($q < 0.05$) to be significantly enriched (Additional Table S11, Additional Table S12).
269 The expanded gene families were mainly related to metabolic processes, intermediate
270 filament terms, enzyme activities, and immune terms. For example, cellular metabolic
271 process (adjusted $P = 6.06E-14$), intermediate filament (adjusted $P = 3.42E-15$),
272 keratin filament (adjusted $P = 2.94E-13$), endoribonuclease activity (adjusted
273 $P = 9.19E-08$), and immune response ($q = 8.36E-03$) were enriched (Additional Table
274 S11, Additional Table S12). In addition, for the contracted gene families, 220 GO
275 terms (adjusted $P < 0.05$) and 9 KEGG pathways ($q < 0.05$) were enriched,
276 respectively (Additional Table S13, Additional Table S14). These enriched terms were
277 mainly involved in ion binding and transporter activity, including neurotransmitter
278 transporter activity (adjusted $P = 1.89E-11$), sodium ion transmembrane transporter
279 activity (adjusted $P = 3.33E-06$), and secondary active transmembrane transporter
280 activity (adjusted $P = 1.86E-08$) (Additional Table S13, Additional Table S14). Thus,
281 these biological processes may be related to the special characteristics of the mustache
282 toad.

283

284 **Relative evolutionary rate of species**

285 The evolutionary rate of species can reflect its evolutionary history and status. The
286 relative evolutionary rate of the mustache toad to other closely related species was

287 analyzed using LINTRE [47] and MEGA (MEGA, RRID:SCR_000667, v7.0.26)
288 software. Two-cluster analysis was applied to test the molecular evolution of multiple
289 sequences in a phylogenetic context, based on concatenated supergenes (protein
290 sequences) using *tpcv* (a module in LINTRE software). Concatenated supergenes
291 were also used for Tajima's relative rate test. We used zebrafish as the outgroup in
292 both methods, and found that, except for the axolotl, the mustache toad had a
293 relatively faster evolutionary rate than its closely related species (e.g., *X. tropicalis*, *R.*
294 *marina*, *R. catesbeiana*, and *N. parkeri*) (Additional Table S15, Additional Table S16).
295 The crocodile had a slower evolutionary rate, relative to its closely related species,
296 which is consistent with previous work [48] (Additional Table S15, Additional Table
297 S16).

298

299 **Discussion**

300 Using Illumina, PacBio, and Hi-C sequencing technologies, we report the first
301 chromosome-level genome assembly of the mustache toad. We successfully annotated
302 the high-quality protein-coding genes by integrating results from three different
303 methods. Phylogenetic analysis indicated that the mustache toad is closely related to
304 the marine toad, bullfrog, and Tibetan frog. Analysis showed that the mustache toad
305 had a faster evolutionary rate, relative to most other closely related species studied.
306 Analysis of the expansion and contraction of gene families identified several
307 biological processes and pathways, such as metabolism and intermediate filaments,
308 suggesting that these terms may relate to the special adaptations of the mustache toad

309 to its habitat. This work not only offers a valuable chromosome-level genomic data
310 for comparative genomics analysis, but also provides important genomic data for
311 studying the mustache toad traits.

312

313 **Availability of supporting data and materials**

314 Raw sequencing data were deposited in the NCBI database under accession number
315 PRJNA523649. Genome assembly and annotation results are available via the
316 *GigaScience* repository, GigaDB [49].

317

318 **Declarations**

319 **List of Abbreviations**

320 BLAST: Basic Local Alignment Search Tool; BUSCO: Benchmarking Universal
321 Single-Copy Orthologs; GO: Gene Ontology; Hi-C: High-throughput chromosome
322 conformation capture; KEGG: Kyoto Encyclopedia of Genes and Genomes.

323

324 **Competing interests**

325 The authors declare that they have no competing interests.

326

327 **Funding**

328 This work was supported by the National Key Research and Development Program of
329 China (grant number 2017YFC0505202) and the National Natural Science
330 Foundation of China (grant numbers NSFC-30270175, NSFC-30870278, and

331 NSFC-31372165).

332

333 **Authors' contributions**

334 D.R. designed the project; D.R. and D.Z. collected the samples; Y.L. and Y.R.
335 estimated the genome size and assembled the genome; Y.L. polished the assembled
336 genome and analyzed Hi-C data; H.J. performed the genome annotation; Y.L. and Z.W.
337 assessed the quality of the genome assembly; Y.L. and Y.R. constructed the
338 phylogenetic tree and determined divergence time, relative evolutionary rate of
339 species, and expansion and contraction of gene families. Y.L., D.R., Y.R., and X.L.
340 wrote the manuscript. All authors read and approved the final version of the
341 manuscript.

342

343 **Acknowledgements**

344 Not applicable.

345

346 **References**

- 347 1. Liang F and Changyuan Y. Amphibians of China. Science Press: Beijing;
348 2016. ISBN-10: 9401795630
- 349 2. Matsui M, Hamidy A, Murphy RW, Khonsue W, Yambun P, Shimada T, et al.
350 Phylogenetic relationships of megophryid frogs of the genus *Leptobrachium*
351 (Amphibia, Anura) as revealed by mtDNA gene sequences. *Molecular*
352 *Phylogenetics & Evolution*. 2010;56 1:259-72.

- 353 3. Matsui M. A New *Leptobrachium* (Vibrissaphora) from Laos (Anura:
354 Megophryidae). *Current Herpetology*. 2013;32 2:182-9.
- 355 4. Liu C HS, Zhao EJAHS, Chengdu, Old Ser. Preliminary study of genus
356 *Vibrissaphora* (Amphibia: Salientia) and discussion on problems of amphibian
357 classification. *Copeia*. 1980;3:1-9.
- 358 5. Rao DQ and Wilkinson JA. Phylogenetic relationships of the mustache toads
359 inferred from mtDNA sequences. *Molecular Phylogenetics & Evolution*.
360 2008;46 1:61-73.
- 361 6. Zheng Y, Li S and Fu J. A phylogenetic analysis of the frog genera
362 *Vibrissaphora* and *Leptobrachium*, and the correlated evolution of nuptial
363 spine and reversed sexual size dimorphism. *Molecular Phylogenetics &*
364 *Evolution*. 2008;46 2:695-707.
- 365 7. Zheng Y, Deng D, Li S and Fu J. Aspects of the breeding biology of the Emei
366 mustache toad (*Leptobrachium boringii*): Polygamy and paternal care.
367 *Amphibia-Reptilia*. 2010;31 2:183-94.
- 368 8. Zhang W, Guo Y, Li J, Huang L, Kazitsa EG and Wu H. Transcriptome
369 analysis reveals the genetic basis underlying the seasonal development of
370 keratinized nuptial spines in *Leptobrachium boringii*. *Bmc Genomics*. 2016;17
371 1:978.
- 372 9. Zheng Y, Rao D, Murphy RW and Zeng X. Reproductive Behavior and
373 Underwater Calls in the Emei Mustache Toad, *Leptobrachium boringii*. *Asian*
374 *Herpetological Research*. 2011;02 4:199-215.

- 375 10. Hudson CM, Xianjin HE and Jinzhong FU. Keratinized Nuptial Spines Are
376 Used for Male Combat in the Emei Moustache Toad (*Leptobrachium*
377 *boringii*). Asian Herpetological Research. 2011;02 3:142-8.
- 378 11. Genome 10K: a proposal to obtain whole-genome sequence for 10,000
379 vertebrate species. The Journal of heredity. 2009;100 6:659-74.
380 doi:10.1093/jhered/esp086.
- 381 12. WTDBG. <https://github.com/ruanjue/wtdbg-1.2.8>.
- 382 13. Prysycz LP and Gabaldón T. Redundans: an assembly pipeline for highly
383 heterozygous genomes. Nucleic Acids Research. 2016;44 12:e113-e.
- 384 14. Simão FA, Waterhouse RM, Panagiotis I, Kriventseva EV and Zdobnov EM.
385 BUSCO: assessing genome assembly and annotation completeness with
386 single-copy orthologs. Bioinformatics. 2015;31 19:3210-2.
- 387 15. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al.
388 Pilon: an integrated tool for comprehensive microbial variant detection and
389 genome assembly improvement. PLoS One. 2014;9 11:e112963.
390 doi:10.1371/journal.pone.0112963.
- 391 16. Servant N, Varoquaux N, Lajoie BR, Viara E, Chen CJ, Vert JP, et al. HiC-Pro:
392 an optimized and flexible pipeline for Hi-C data processing. Genome Biology.
393 2015;16 1:259.
- 394 17. Durand N, Shamim M, Machol I, Rao SP, Huntley M, Lander E, et al. Juicer
395 Provides a One-Click System for Analyzing Loop-Resolution Hi-C
396 Experiments. Cell Systems. 2016;3 1:95-8.

- 397 18. Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, et al.
398 De novo assembly of the *Aedes aegypti* genome using Hi-C yields
399 chromosome-length scaffolds. *Science*. 2017;356 6333:92.
- 400 19. Wilkinson JA. A New Species of the Genus *Vibrissaphora* (Anura:
401 Megophryidae) from Yunnan Province, China. *Herpetologica*. 2006;62 1:90-5.
- 402 20. Li H and Durbin R. Fast and accurate short read alignment with
403 Burrows-Wheeler transform. 2009.
- 404 21. Chang Z, Li G, Liu J, Zhang Y, Ashby C, Liu D, et al. Bridger: a new
405 framework for de novo transcriptome assembly using RNA-seq data. *Genome*
406 *Biology*. 2015;16 1:30.
- 407 22. Perteza G, Huang X, Liang F, Antonescu V, Sultana R, Karamycheva S, et al.
408 TIGR Gene Indices clustering tools (TGICL): a software system for fast
409 clustering of large EST datasets. *Bioinformatics*. 2003;19 5:651-2.
- 410 23. Benson G. Tandem repeats finder: a program to analyze DNA sequences.
411 *Nucleic Acids Res*. 1999;27 2:573-80.
- 412 24. Bedell JA, Korf I, . and Gish W, . MaskerAid: a performance enhancement to
413 RepeatMasker. *Bioinformatics*. 2000;16 11:1040-1.
- 414 25. Stanke M and Waack S. Gene prediction with a hidden Markov model and a
415 new intron submodel. *Bioinformatics*. 2003;19 suppl_2:215--25.
- 416 26. Kerstin H, Clark MD, Torroja CF, James T, Camille B, Matthieu M, et al. The
417 zebrafish reference genome sequence and its relationship to the human
418 genome. *Nature*. 2013. 496(7446):498-503. doi: 10.1038/nature12111.

- 419 27. Yan-Bo S, Zi-Jun X, Xue-Yan X, Shi-Ping L, Wei-Wei Z, Xiao-Long T, et al.
420 Whole-genome sequence of the Tibetan frog *Nanorana parkeri* and the
421 comparative evolution of tetrapod genomes. Proc Natl Acad Sci U S A.
422 2015;11 112:E1257-E62.
- 423 28. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al.
424 Initial sequencing and analysis of the human genome. Nature. 2001;409
425 6822:860-921. doi:10.1038/35057062.
- 426 29. Sequence and comparative analysis of the chicken genome provide unique
427 perspectives on vertebrate evolution. Nature. 2004;432 7018:695-716.
428 doi:10.1038/nature03154.
- 429 30. Wang Z, Pascual-Anaya J, Zadissa A, Li W, Niimura Y, Huang Z, et al. The
430 draft genomes of soft-shell turtle and green sea turtle yield insights into the
431 development and evolution of the turtle-specific body plan. Nat Genet.
432 2013;45 6:701-6. doi:10.1038/ng.2615.
- 433 31. Session AM, Uno Y, Kwon T, Chapman JA, Toyoda A, Takahashi S, et al.
434 Genome evolution in the allotetraploid frog *Xenopus laevis*. Nature. 2016;538
435 7625:336-43. doi:10.1038/nature19840.
- 436 32. Smith JJ and Timoshevskaya N. The sea lamprey germline genome provides
437 insights into programmed genome rearrangement and vertebrate evolution.
438 2018;50 2:270-7. doi:10.1038/s41588-017-0036-1.
- 439 33. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated
440 eukaryotic gene structure annotation using EVIDENCEModeler and the Program

- 441 to Assemble Spliced Alignments. *Genome Biology*. 2008;9 1:R7.
- 442 34. Jessica AL, Federica DP, Manfred G, Christina W, Lesheng K, Evan M, et al.
443 The genome of the green anole lizard and a comparative analysis with birds
444 and mammals. *Nature*. 2011.
- 445 35. Edwards RJ, Tuipulotu DE, Amos TG, O'Meally D, Richardson MF, Russell
446 TL, et al. Draft genome assembly of the invasive cane toad, *Rhinella marina*.
447 *GigaScience*. 2018;7 9 doi:10.1093/gigascience/giy095.
- 448 36. Hammond SA, Warren RL, Vandervalk BP, Kucuk E, Khan H, Gibb EA, et al.
449 The North American bullfrog draft genome provides insight into hormonal
450 regulation of long noncoding RNA. 2017;8 1:1433.
451 doi:10.1038/s41467-017-01316-7.
- 452 37. Smith JJ, Timoshevskaya N, Timoshevskiy VA, Keinath MC, Hardy D and
453 Voss SR. A chromosome-scale assembly of the axolotl genome. *Genome Res*.
454 2019;29 2:317-24. doi:10.1101/gr.241901.118.
- 455 38. Nowoshilow S, Schloissnig S, Fei JF, Dahl A, Pang AWC, Pippel M, et al. The
456 axolotl genome and the evolution of key tissue formation regulators. *Nature*.
457 2018;554 7690.
- 458 39. Wan QH, Pan SK, Hu L, Zhu Y, Xu PW, Xia JQ, et al. Genome analysis and
459 signature discovery for diving and sensory properties of the endangered
460 Chinese alligator. *Cell research*. 2013;23 9:1091-105.
461 doi:10.1038/cr.2013.104.
- 462 40. Li L, Jr SC and Roos DS. OrthoMCL: identification of ortholog groups for

- 463 eukaryotic genomes. *Genome Research*. 2003;13 9:2178-89.
- 464 41. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and
465 high throughput. *Nucleic Acids Res*. 2004;32 5:1792-7.
466 doi:10.1093/nar/gkh340.
- 467 42. Edgar RC. MUSCLE: a multiple sequence alignment method with reduced
468 time and space complexity. *BMC Bioinformatics*. 2004; 5 (1)
469 DOI:10.1186/1471-2105-5-113
- 470 43. Alexandros S. RAxML version 8: a tool for phylogenetic analysis and
471 post-analysis of large phylogenies. *Bioinformatics*. 2014;30 9:1312-3.
- 472 44. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular*
473 *biology and evolution*. 2007;24 8:1586-91. doi:10.1093/molbev/msm088.
- 474 45. TIMETREE website. www.timetree.org.
- 475 46. Tjil DB, Nello C, Demuth JP and Hahn MW. CAFE: a computational tool for
476 the study of gene family evolution. *Bioinformatics*. 2006;22 10:1269-71.
- 477 47. Takezaki N, Rzhetsky A, and Nei M. Phylogenetic test of the molecular clock
478 and linearized trees. *Molecular Biology & Evolution*. 1995;12 5:823-33.
- 479 48. Green RE, Braun EL, Joel A, Dent E, Ngan N, Glenn H, et al. Three
480 crocodylian genomes reveal ancestral patterns of evolution among archosaurs.
481 *Science*. 2014;346 6215:1254449.
- 482 49. Li Y; Ren Y; Zhang D; Jiang H; Wang Z; Li X; Rao D (2019): Supporting data
483 for "Chromosomal-level assembly of the mustache toad genome using
484 third-generation DNA sequencing and Hi-C analysis" GigaScience Database.

485 <http://dx.doi.org/10.5524/100624>

486

487 **Tables**

488 **Table 1: Sequencing data used for mustache toad genome assembly and**

489 **annotation**

Sequencing type	Platform	Library size (bp)	Clean data (Gb)	Application
Genome long reads	PacBio Sequel	20,000	277.15	Contig assembly
Genome short reads	Illumina HiSeq 2500	250	225.03	Genome survey, genome base correction, and genome assessment
Genome Hi-C reads	Illumina HiSeq X-Ten	250	378.78	Chromosome construction
Transcriptome short reads	Illumina HiSeq 4000	250	14.18	Genome annotation and assessment

490

491 **Table 2: Assembly data for the mustache toad genome**

Term	Wtdbg contig		Hi-C scaffold	
	Size (bp)	Number	Size (bp)	Number
N90	153,029	4,866	134,864,763	11

N80	301,658	3,285	181,461,513	8
N70	456,829	2,334	220,042,448	6
N60	624,716	1,671	359,321,214	5
N50	821,125	1,180	412,424,790	4
Max length (bp)	9,978,207		592,710,058	
Total size (bp)	3,530,531,046		3,535,795,546	
Total number (>100bp)	15,899		5,370	

492 Note: These data pertain to genome assembly. Wtdbg contig was the genome
 493 assembled by wtdbg and 2-round pilon error-correction. Hi-C scaffold was the
 494 genome finished by Hi-C assembly.

495

496 **Table 3: Assessment of genome assembly and annotation completeness of the**
 497 **mustache toad and other amphibian genomes, using Benchmarking Universal**
 498 **Single-Copy Orthologs (BUSCO)**

Library	<i>V. ailaonica</i> (eukaryota)	<i>V. ailaonica</i> (metazoa)	<i>Nanorana</i> <i>parkeri</i> (eukaryota)	<i>Xenopus</i> <i>tropicalis</i> (eukaryota)	<i>Rhinella</i> <i>marina</i> (eukaryota)	<i>Rana</i> <i>catesbeiana</i> (eukaryota)	<i>Ambystoma</i> <i>mexicanum</i> (eukaryota)
Complete genes (%)	80.8	85.1	90.1	90.1	90.4	58.0	24.4
Complete and single-copy	78.2	83.6	87.8	88.1	86.1	55.4	23.4

genes (%)							
Complete and duplicated genes (%)	2.6	1.5	2.3	2.0	4.3	2.6	1.0
Fragmented genes (%)	8.9	4.9	3.6	2.0	3.3	20.8	24.4
Missing genes (%)	10.3	10.0	6.3	7.9	6.3	21.2	51.2

499 Note: Both 'eukaryote' and 'metazoan' are two core gene sets in the BUSCO database

500

501 **Table 4: Quality data for several published amphibian genomes**

Species	Contig N50 (bp)	Scaffold N50 (bp)	Genome size (bp)	Genome BUSCO (eukaryota) (%)
<i>Nanorana parkeri</i>	32,798	1,069,101	2,053,867,363	90.1
<i>Xenopus tropicalis</i>	71,041	135,134,832	1,440,398,454	90.1
<i>Rhinella marina</i>	166,489	167,498	2,551,759,918	90.4
<i>Rana catesbeiana</i>	5,415	39,363	6,250,353,185	58.0
<i>Ambystoma mexicanum</i>	216,366	3,052,786	32,393,605,577	24.4

502

503 **Table 5: De novo-annotated repeat sequences in the mustache toad genome**

Type	Length (bp)	Percentage in genome (%)
------	-------------	--------------------------

DNA	350,793,270	9.94
LINE	297,954,803	8.45
SINE	11,009,363	0.31
LTR	307,317,539	8.71
Other	43,867,330	1.24
Satellite	9,696,790	0.27
Simple repeat	125,397,072	3.55
Unknown	1,114,326,962	31.59
Total	2,147,505,764	60.87

504

505 **Table 6: Functional annotation for protein-coding genes in the mustache toad**506 **genome**

Database	Annotated gene number	Percent (%)
Interpro	12,997	49.56
KEGG	10,035	38.26
SwissProt	12,410	47.32
Trembl	17,916	68.31

507

508 **Figure legends**509 **Figure 1: The mustache toad, *Vibrissaphora ailaonica*.**510 **(A)** The adult male individual with spines in the upper jaw. **(B)** The adult female511 individual. **(C)** The adult male individual during the process of spines shedding from

512 the upper jaw. **(D)** The adult male individual without spines (after spine have been
513 shed) in the upper jaw. **(E)** The body size of the mustache toad, side view: male (left)
514 and female (right). **(F)** The body size of mustache toad, top view: male (left) and
515 female (right).

516

517 **Figure 2: 17-mer analysis of *Vibrissaphora ailaonica* genome characteristics.**

518

519 **Figure 3: Circos graph showing characteristics of the mustache toad genome.**

520 From outer circle to inner ring: gene distribution, tandem repeats (TR), long tandem
521 repeats (LTR), long interspersed nuclear elements (LINE), short interspersed nuclear
522 elements (SINE), and GC content.

523

524 **Figure 4: Length distributions of annotated protein-coding genes in these species.**

525 The species including *Vibrissaphora ailaonica*, *Homo sapiens*, *Danio rerio*, *Gallus*
526 *gallus*, *Anolis carolinensis*, *Alligator sinensis*, and *Nanorana parkeri*.

527

528 **Figure 5: The statistics of gene family among these 11 species.** The species

529 including *Danio rerio*, *Rana catesbeiana*, *Rhinella marina*, *Vibrissaphora ailaonica*,
530 *Nanorana parkeri*, *Homo sapiens*, *Gallus gallus*, *Anolis carolinensis*, *Xenopus*
531 *tropicalis*, *Ambystoma mexicanum*, and *Alligator sinensis*.

532

533 **Figure 6: The phylogenetic relationships among these species.** The species

534 including *Danio rerio*, *Rana catesbeiana*, *Rhinella marina*, *Vibrissaphora ailaonica*,
535 *Nanorana parkeri*, *Homo sapiens*, *Gallus gallus*, *Anolis carolinensis*, *Xenopus*
536 *tropicalis*, *Ambystoma mexicanum*, and *Alligator sinensis*. Blue numbers represent
537 divergence time. The red dot represents the fossil record used in the node.

538

539 **Additional files**

540 Additional Figure S1: The GC content in these genomes. The species including *Gallus*
541 *gallus*, *Vibrissaphora ailaonica*, *Alligator sinensis*, *Nanorana parkeri*, *Homo sapiens*,
542 *Anolis carolinensis*, *Xenopus tropicalis*, and *Danio rerio*.

543 Additional Figure S2: The divergence time of these species (using first-codon sites).

544 The species including *Danio rerio*, *Rana catesbeiana*, *Rhinella marina*, *Vibrissaphora*
545 *ailaonica*, *Nanorana parkeri*, *Homo sapiens*, *Gallus gallus*, *Anolis carolinensis*,
546 *Xenopus tropicalis*, *Ambystoma mexicanum*, and *Alligator sinensis*.

547 Additional Figure S3: The divergence time of these species (using second-codon sites).

548 The species including *Danio rerio*, *Rana catesbeiana*, *Rhinella marina*, *Vibrissaphora*
549 *ailaonica*, *Nanorana parkeri*, *Homo sapiens*, *Gallus gallus*, *Anolis carolinensis*,
550 *Xenopus tropicalis*, *Ambystoma mexicanum*, and *Alligator sinensis*.

551 Additional Table S1: Illumina sequencing clean data.

552 Additional Table S2: PacBio Sequel sequencing data.

553 Additional Table S3: RNA-sequencing clean data.

554 Additional Table S4: Hi-C sequencing clean data.

555 Additional Table S5: Comparison of the BUSCO and Illumina read mapping results

556 between the raw genome and redundancy-filtered genome.

557 Additional Table S6: Illumina read mapping ratio to the assembled genome.

558 Additional Table S7: The statistics of assembled transcripts by Bridger software. The
559 redundant transcripts were removed by TGICL software.

560 Additional Table S8: Transcript mapping ratio to the assembled genome.

561 Additional Table S9: Annotated repeat sequences in our assembled genome.

562 Additional Table S10: Gene families among these species. The species including
563 *Danio rerio*, *Rana catesbeiana*, *Rhinella marina*, *Vibrissaphora ailaonica*, *Nanorana*
564 *parkeri*, *Homo sapiens*, *Gallus gallus*, *Anolis carolinensis*, *Xenopus tropicalis*,
565 *Ambystoma mexicanum*, and *Alligator sinensis*.

566 Additional Table S11: Gene Ontology (GO) enrichment analysis of expanded gene
567 families.

568 Additional Table S12: Kyoto Encyclopedia of Genes and Genomes (KEGG)
569 enrichment analysis of expanded gene families.

570 Additional Table S13: Gene Ontology (GO) enrichment analysis of contracted gene
571 families.

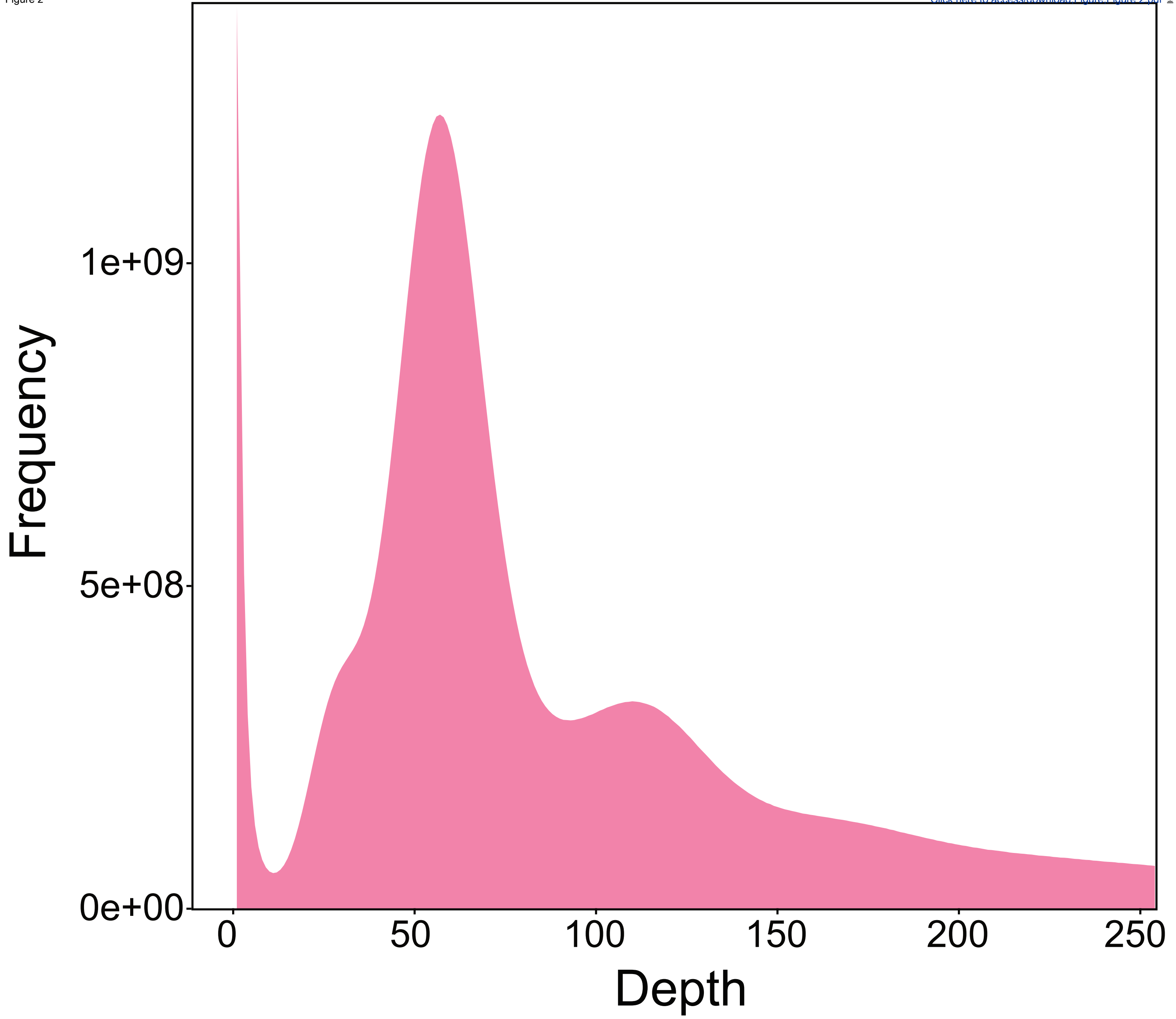
572 Additional Table S14: Kyoto Encyclopedia of Genes and Genomes (KEGG)
573 enrichment analysis of contracted gene families.

574 Additional Table S15: Two-cluster analysis of mustache toad and other species.

575 Additional Table S16: The relative evolutionary rate of mustache toad and other
576 species analyzed by Tajima's test.

577





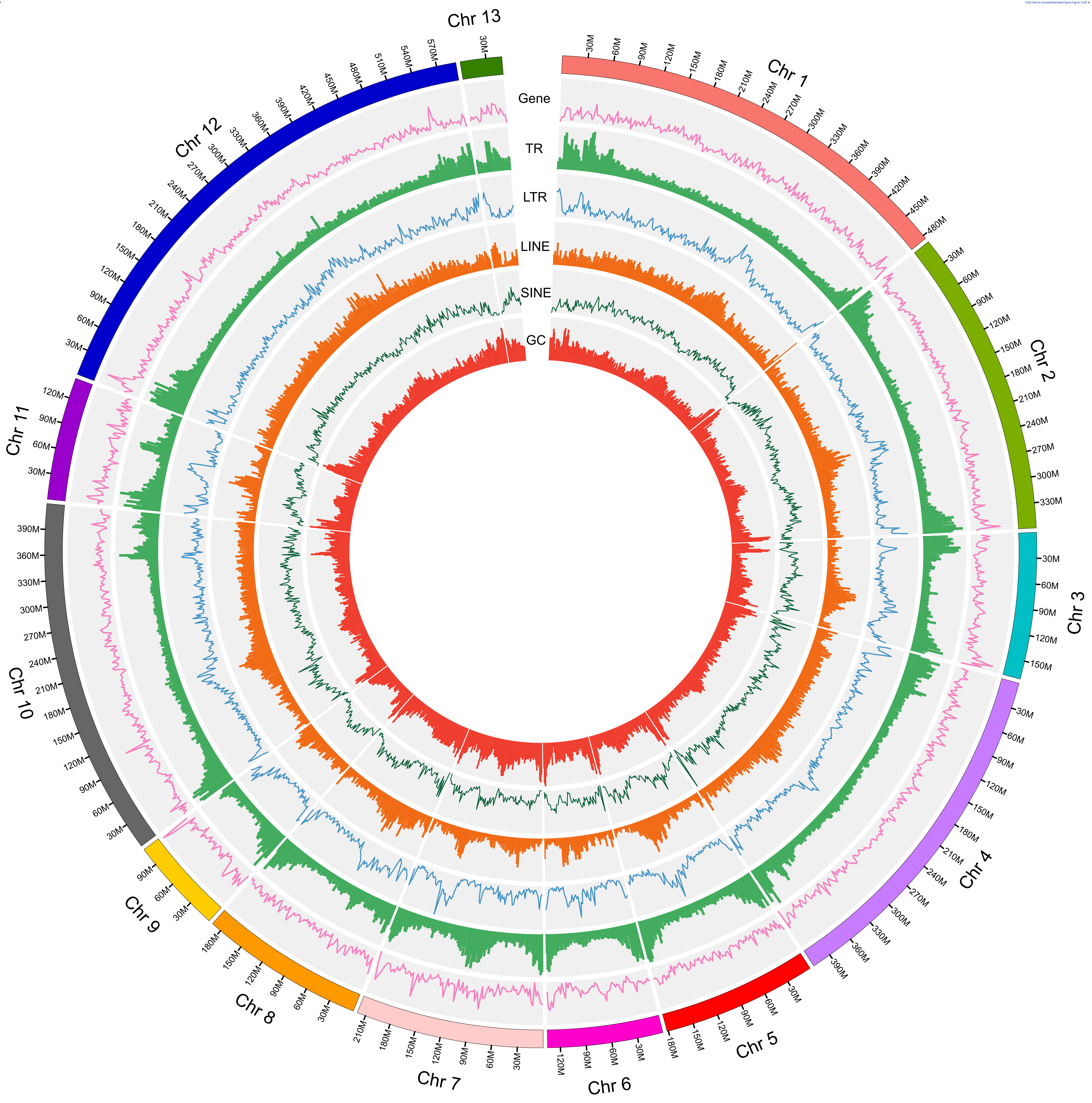
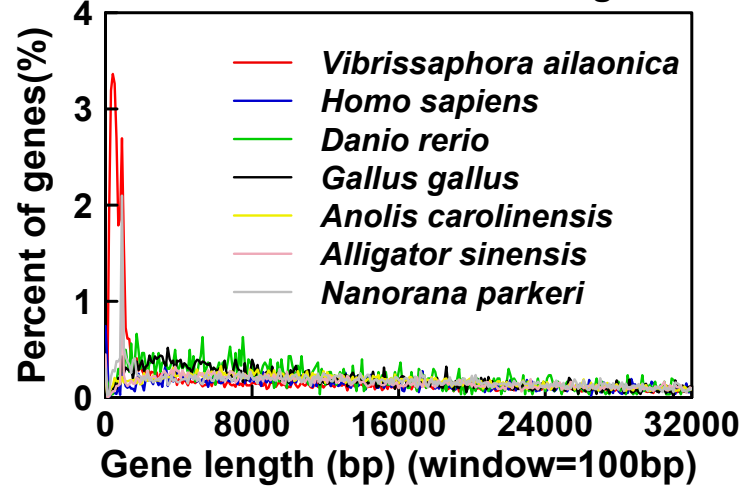


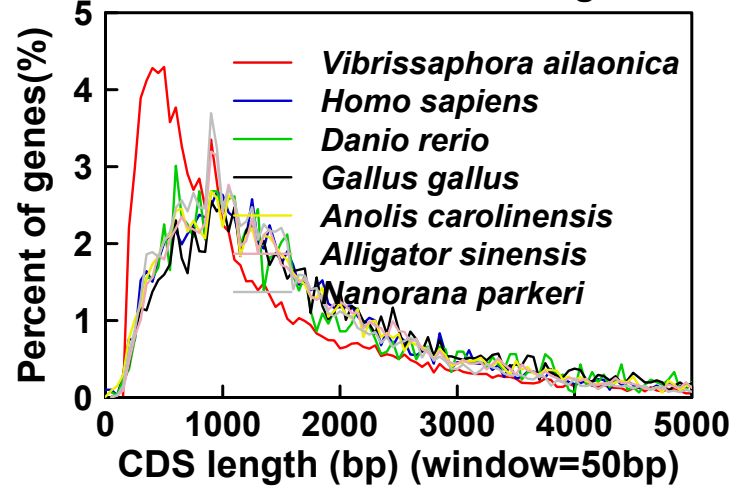
Figure 4

[Click here to access/download/Figure/Figure 4.pdf](#)

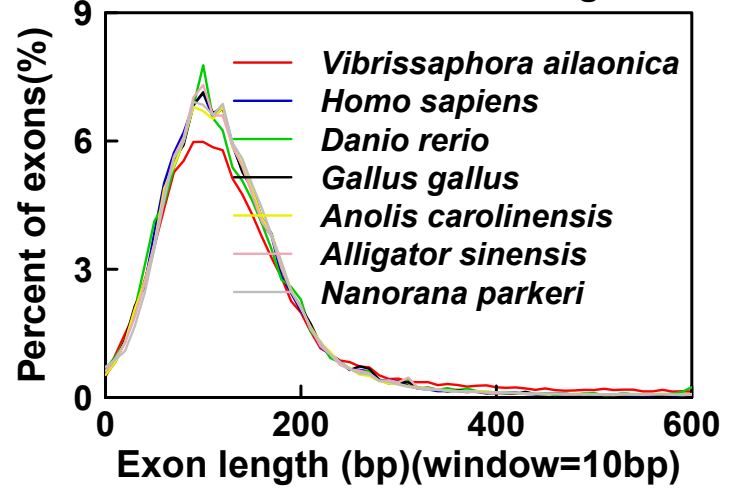
Distribution of mRNA length



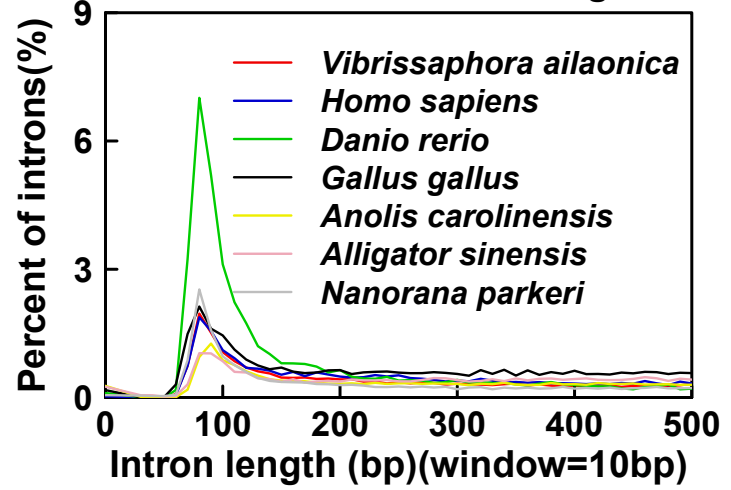
Distribution of CDS length



Distribution of exon length



Distribution of intron length



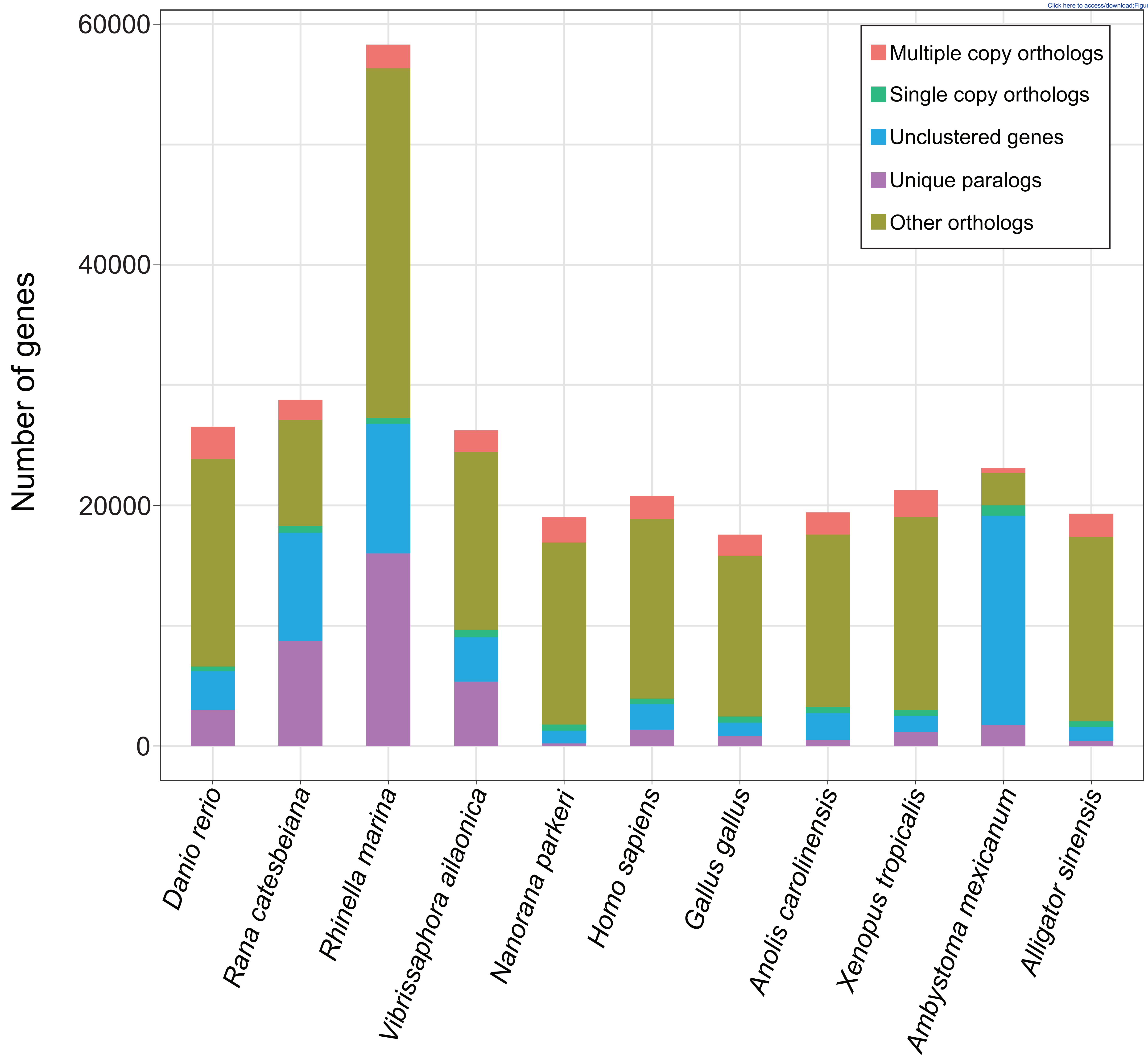
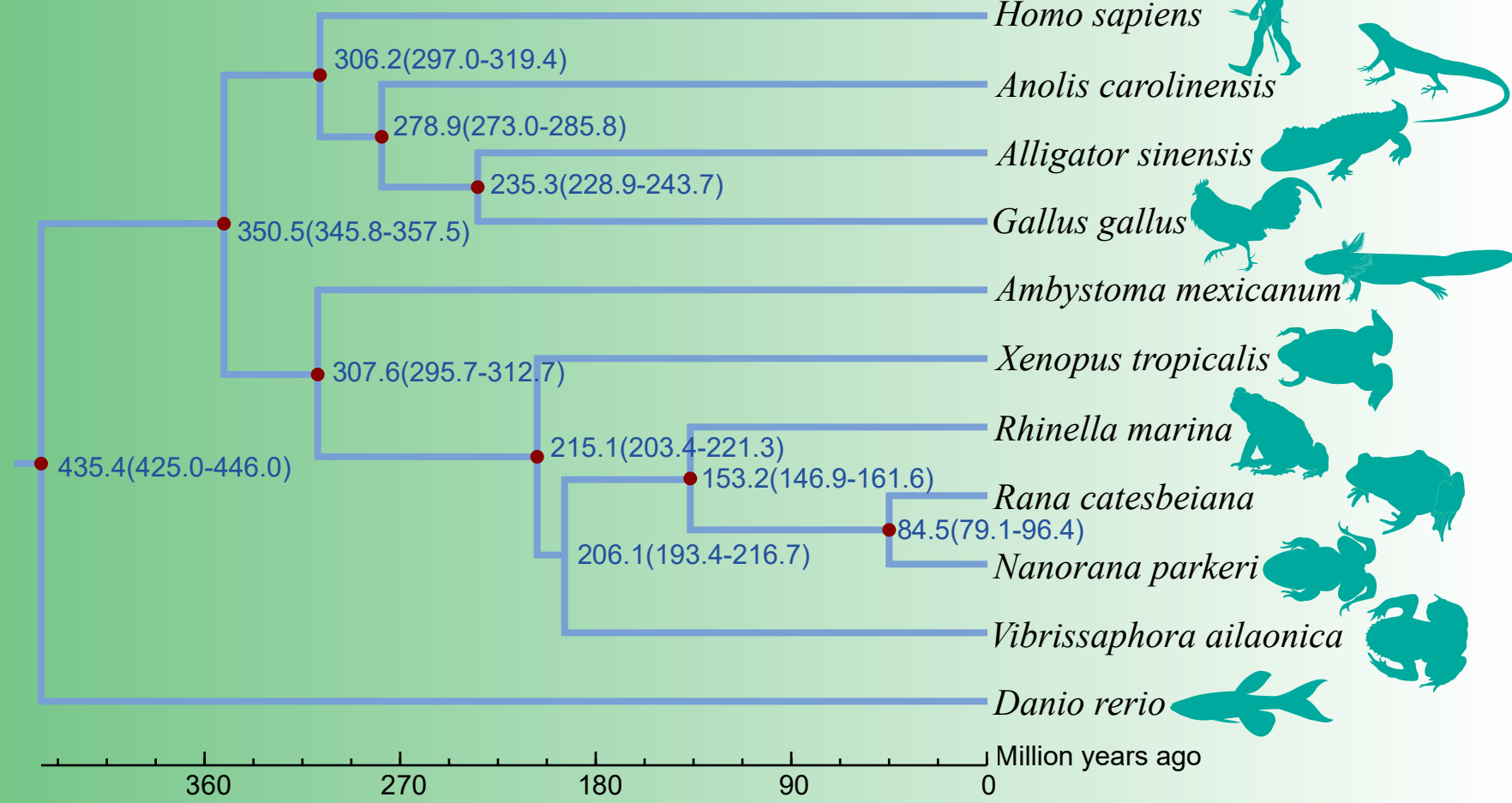


Figure 6

[Click here to access/download;Figure;Figure 6.pdf](#)





Click here to access/download
Supplementary Material
Supplementary Materials.doc

