

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

--Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-19-00099R1	
<b>Full Title:</b>	Chromosomal-level assembly of the mustache toad genome using third-generation DNA sequencing and Hi-C analysis	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	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
<b>Abstract:</b>	<p><b>Background</b></p> <p>The mustache toad, <i>Vibrissaphora ailaonica</i>, is an endemic species to China belonging 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 that fall off when the breeding season ends, which probably reversed the sexual size dimorphism with males being larger than females. To investigate the genetic mechanism of the repeatedly developed keratinized spines, a high-quality reference genome of mustache toad would be a valuable resource.</p> <p><b>Findings</b></p> <p>For genome construction, we generated 225 Gb of short reads and 277 Gb of long reads using Illumina and Pacific Biosciences (PacBio) sequencing, respectively. The sequencing data were assembled into a 3.53 Gb genome assembly with a contig N50 length of 821 Kb. Additionally, we applied Hi-C technology to identify contacts among contigs, then assembled contigs into scaffolds and identified a genome assembly 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 the phylogenetic relationships of the mustache toad with other chordate species. Results showed that the mustache toad has a relatively higher evolutionary rate and separated from the marine toad, bull frog, and Tibetan frog ancestor 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><b>Conclusions</b></p> <p>Using Illumina, PacBio, and Hi-C technologies, we constructed the first high-quality chromosomal-level mustache toad genome. This work not only offers a valuable reference genome for functional studies of mustache toad traits, but also provides important chromosome information for wider genome comparisons.</p>	
<b>Corresponding Author:</b>	Dingqi Rao, Ph.D Kunming Institute of Zoology Chinese Academy of Sciences Kunming, Yunnan CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Kunming Institute of Zoology Chinese Academy of Sciences	

<b>Corresponding Author's Secondary Institution:</b>	
<b>First Author:</b>	Yongxin Li
<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	Yongxin Li
	Yandong Ren
	Dongru Zhang
	Hui Jiang
	Zhongkai Wang
	Xueyan Li
	Dingqi Rao, Ph.D
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>Responses to Reviewers</p> <p>The authors would like to thank the reviewers for the helpful comments, and we have revised our manuscript accordingly. At this time, we hope to have an opportunity to publish this paper in GigaScience.</p> <p>Reviewer 1 Comments for the Author...</p> <p>Reviewer #1: Li et al report the first genome assembly of a mustache toad. They used a combination of PacBio and HiC to generate a highly-contiguous assembly. They used RNA-seq data, ab initio gene prediction and homology to annotate ~26000 genes, analyzed gene family contractions and expansions, and estimated the phylogenetic relationship to other amphibians. Given the sparsity of amphibian genomes, this assembly will be valuable for the community. I recommend accepting the manuscript after a few issues have been addressed, most of which are minor.</p> <p>Major comments:</p> <p>1) Since k-mer based genome size estimation is often not very precise, I find the redundancy reduction of the assembly potentially problematic. The authors removed contigs that overlap with at least 70% another contig using an alignment identity cutoff of 70%. It feels a bit like these parameters were optimized such that the final assembly matches the k-mer predicted size. E.g. the 70% identity cutoff is not compatible with the error rate of PacBio reads, unless the purpose of the Redundans run was to remove single reads that contain much more than the ~15% expected error rate. Also, heterozygosity and alt haplotypes should not result in 30% divergence.</p> <p>I wonder if the authors can check which contigs were removed in this step and ensure that no real sequences were removed. If there is any doubt that some of the contigs may contain functionally important sequences (genes, etc), then I would suggest to provide the removed contigs with the redundancy-filtered assembly as an extra fasta file. Specifically, I wonder if the slightly lower BUSCO scores can be explained by removing real contigs based on 70% similarity.</p> <p>Response: To make sure all the removed contigs not contain real sequences, we checked the BUSCO of the raw genome and the redundancy-filtered genome with eukaryote and metazoan as datasets at the same time. Besides, we also checked the mapping ratio of Illumina reads on the raw genome and the redundancy-filtered genome. The BUSCO results shown that most of the genes were not removed and the mapping ratio results shown that both coding region and non-coding region were remained (Table S5). All these results indicated that the redundancy-filtered step not removed many real sequences. However, as you have said before, we further checked the contigs that were removed in this step, all the removed contigs, the BUSCO results of genome and the redundancy-filtered assembly have been uploaded to GigaDB FTP server. Thank you for your suggestions.</p>

2) The manuscript 'undersells' the contiguity of mustache toad assembly, which has \*substantially higher\* contig and scaffold N50 values than any other amphibian genome.

I therefore recommend to place Table S8 in the main text.

Response: Done. Table S8 has been removed to main text (Table 4). Thank you.

3) Table 3 is hard to understand as absolute numbers are reported. A much better way would be to report '%complete genes, %complete and duplicated genes, %fragmented genes, %missing genes' which sums to 100%. In addition, these 4 BUSCO percentages for the other amphibian genomes should be added to this table to provide a direct comparison of genome assembly completeness.

Response: Table 3 has been corrected and the BUSCO results of the other amphibian genomes were also added.

4) I wonder how the divergence time estimates would change if first or second codon positions instead of four-fold degenerate sites were used. This may be relevant as four-fold degenerate sites are clearly saturated over these phylogenetic distances. Also, the divergence times shown in Figure 6 are quite different to the times from timetree, where e.g. the Rana - Nanorana split was 89 Mya (Figure 6, 44 Mya) and the Rana - Rhinella split was 160 Mya (Figure 6, 137 Mya).

Response: We added the fossil evidence between Rana - Nanorana and Rana - Rhinella from timetree to run mcmctree tree again, and both four-fold degenerate sites and first or second codon positions were used for the divergence time analysis (Figure 6; Figures S2 and S3). The results shown that these three results are much closed and the divergence time between Rana - Nanorana, Rana - Rhinella are similar to timetree results. We have revised the divergence time results in Figure 6. Thank you.

Besides, because the expansion and contraction analysis of gene family are related the divergence time result, so we updated the expansion and contraction results this time, including all the descriptions and Tables (Tables S11-S14). Thank you.

Minor comments:

1) The manuscript should be edited by a native speaker to improve the language. A few examples: "Like other mustache toad species, *V. ailaonica* males develop temporary keratinized nuptial spines on their upper jaw during each breeding season and fall off when the breeding season ends, which probably lead to the reverse of the sexual size dimorphism, namely the size of the male get larger than female." should be improved to "Like other mustache toad species, *V. ailaonica* males temporarily develop keratinized nuptial spines on their upper jaw during each breeding season that fall off when the breeding season ends, which probably reversed the sexual size dimorphism with males being larger than females."

"To investigate the genetic mechanism of the repeatedly develop the keratinized spines"

-->

"To investigate the genetic mechanism of the repeatedly developed keratinized spines"

"Another unique aspect of the mustache toad is that breeding occurs during the cold season, unlike most frogs and toads which breed in the warmer months"

-->

"Another unique aspect of the mustache toad is that breeding occurs during the cold season, whereas most frogs and toads breed in the warmer months"

etc.

Response: Done. We have corrected the above mistakes and other mistakes in our manuscript. Thank you.

2) Please reference Figure 1 in line 55, where the temporary spines are described.

Response: Corrected.

3) Line 75-76: I find this outlook that we will learn from the toad genome (sex

dimorphism) how body size control works in general a bit far-stretched. This could be removed.

Response: Corrected.

4) Line 94/95: Please mention the Illumina read length (paired end 150 bp reads). I find this information more important than library size.

Response: Corrected.

5) Line 164/165: The conclusion that the toad assembly is very complete is justified based on the high percentage of mapping RNA-seq reads and transcripts. However, this sentence should be moved to Line 161 (after "Table S8).", where this analysis is done.

Response: Corrected.

6) Line 180: Please replace 'closely-related' with 'vertebrate' as zebrafish, lamprey and amphibians are not really closely related.

Response: Corrected.

7) Line 179: Would an Augustus model trained from an amphibian (e.g. xenopus) be not more appropriate than a zebrafish model?

Response: The official website of Augustus not included the amphibian species as model. So we selected zebrafish that has high-quality gene set as model. Thank you for your suggestions.

8) Table 4: Please round the percentage to 2 digits (9.94%).

Response: Corrected.

9) Line 204/205: The references don't match: Reference 34 ([www.axolotl-omics.org](http://www.axolotl-omics.org)) and 36 refer to the Ambystoma genome assembly. The Rhinella reference is missing.

Response: Corrected.

Reviewer 2 Comments for the Author...

Reviewer #2: Here, Yongxin Li and colleagues reported the chromosome-level genome with the full annotation of the mustache toad, *Vibrissaphora ailaonica*, using conventional paired-end short read, sufficient amount of PacBio long reads and chromosome conformation capture (Hi-C) data. Although there are several amphibian genomes reported previously, many of them do not have chromosome-level genomes, so I think this is definitely a valuable resource to the community, especially to study the synteny of amphibian genome. So I would like to recommend accepting this manuscript for publication after resolving some issues as mentioned below:

1) On page 5, more details for RNA-Seq library prep construction method should be provided (poly-A capturing or ribosome-depletion? Which library kit do they use?). Also, even the authors mentioned that 9 tissues were dissected from the biospecimen they sequenced the genome (page 4, line 85), it is not clear whether all those tissues were used in this 'mixed RNA-Seq' experiment. Please provide more details for this experiment.

Response: The details for RNA-seq library construction method were added. About the RNA-seq experiment, after equally mixed the DNA of the 9 tissues, the mixed DNA sample was used for library construction and RNA-seq experiments. Thank you.

2) On page 5, the authors mentioned that they used four Hi-C libraries. Are they constructed from the same samples (blood), with the same parameter (four technical replicates)? Or using different samples? If they used a different parameter to construct these four libraries, it should be specified.

	<p>Response: Yes, all these four Hi-C libraries were used the same samples with the same parameter. We also clarify these informations in the main text this time. Thank you for your suggestions.</p> <p>3) Authors claimed that they deposited the data on PRJNA523649, but it looks like they uploaded one single file for each data set. Because they used different libraries, at least for paired-end seq (Table S1) and HiC-seq (Table S4), it would be better to provide those raw data separately.</p> <p>Response: Thank you for your suggestions, we uploaded these sequencing data in the same PROJECT ID, but with different SRA IDs, you could see this by this link (<a href="https://www.ncbi.nlm.nih.gov/sra/?term=PRJNA523649">https://www.ncbi.nlm.nih.gov/sra/?term=PRJNA523649</a>). Thank you.</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. 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>	Yes
<p><b>Resources</b></p> <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 <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes
<b>Availability of data and materials</b>	Yes

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.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

[Click here to view linked References](#)

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

3

4 Yongxin Li<sup>1,2,†</sup>, Yandong Ren<sup>2,†</sup>, Dongru Zhang<sup>1,†</sup>, Hui Jiang<sup>3</sup>, Zhongkai Wang<sup>2</sup>, Xueyan Li<sup>1</sup>,  
5 Dingqi Rao<sup>1,\*</sup>

6

7 1. State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology,  
8 Chinese Academy of Sciences, Kunming 650223, China

9 2. Center for Ecological and Environmental Sciences, Northwestern Polytechnical University,  
10 Xi'an 710072, China

11 3. National Engineering Laboratory of Marine Germplasm Resources Exploration and  
12 Utilization, Zhejiang Ocean University, Zhoushan 316022, China

13

14 †These authors have the equal contribution.

15 \*Corresponding author: Dingqi Rao (kizar@mail.kiz.ac.cn).

16

17

18

19

20

21

22

23 **Abstract**

24 **Background:** The mustache toad, *Vibrissaphora ailaonica*, is an endemic species to China  
25 belonging to the Megophryidae family. Like other mustache toad species, *V. ailaonica* males  
26 temporarily develop keratinized nuptial spines on their upper jaw during each breeding season  
27 that fall off when the breeding season ends, which probably reversed the sexual size  
28 dimorphism with males being larger than females. To investigate the genetic mechanism of  
29 the repeatedly developed keratinized spines~~Like other mustache toad species, *V. ailaonica*~~  
30 ~~males develop temporary keratinized nuptial spines on their upper jaw during each breeding~~  
31 ~~season and fall off when the breeding season ends, which probably lead to the reverse of the~~  
32 ~~sexual size dimorphism, namely the size of the male get larger than female. To investigate the~~  
33 ~~genetic mechanism of the repeatedly develop the keratinized spines~~, a high-quality reference  
34 genome of mustache toad would be a valuable resource. **Findings:** For genome construction,  
35 we generated 225 Gb of short reads and 277 Gb of long reads using Illumina and Pacific  
36 Biosciences (PacBio) sequencing, respectively. The sequencing data were assembled into a  
37 3.53 Gb genome assembly with a contig N50 length of 821 Kb. Additionally, we applied Hi-C  
38 technology to identify contacts among contigs, then assembled contigs into scaffolds and  
39 identified a genome assembly with 13 chromosomes and a scaffold N50 length of 412.42 Mb.  
40 Based on the 26,227 protein-coding genes annotated in the genome, we analyzed the  
41 phylogenetic relationships of the mustache toad with other chordate species. Results showed  
42 that the mustache toad has a relatively higher evolutionary rate and separated from the marine  
43 toad, bull frog, and Tibetan frog ancestor ~~206.1+94.8~~ million years ago. Furthermore, we  
44 identified ~~201349~~ expanded gene families in the mustache toad, which were mainly enriched



45 in immune pathway, keratin filament, and metabolic processes. **Conclusions:** Using Illumina,  
46 PacBio, and Hi-C technologies, we constructed the first high-quality chromosomal-level  
47 mustache toad genome. This work not only offers a valuable reference genome for functional  
48 studies of mustache toad traits, but also provides important chromosome information for  
49 wider genome comparisons.

50

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

52

### 53 **Introduction**

54 The mustache toad, *Vibrissaphora ailaonica* (NCBI Taxonomy ID: 428466), belongs to the  
55 Megophryidae family and is an endemic amphibian species to China (including the  
56 China-Vietnam border) [1-3]. This mustache toad species exhibits many interesting features,  
57 including unique keratinized spines along the upper jaw [1, 4-6]. These spines repeatedly  
58 grow in sexually mature males during the breeding season and fall off at the end of this  
59 process [5-8] ([Figure 1](#)). This morphological difference between males and females is further  
60 highlighted by their sexual dimorphism in body size (males are significantly larger than  
61 females) and may be used as a weapon for sexually mature individuals to compete for nests  
62 and mating opportunities [7, 9, 10]. [Another unique aspect of the mustache toad is that](#)  
63 [breeding occurs during the cold season, whereas most frogs and toads breed in the warmer](#)  
64 [months. Another unique aspect of the mustache toad is that breeding occurs during the cold-](#)  
65 [season, unlike most frogs and toads which breed in the warmer months](#) [1]. However, despite  
66 the importance of the mustache toad in spine dynamic development and sexual dimorphism in

67 body size, the genomic resources for the species remain limited. To date, no next-generation  
68 sequencing (NGS) data have been reported in the *Vibrissaphora* genus. Therefore, this lack of  
69 genome sequence and transcriptome data for this important species (*V. ailaonica*) has  
70 hindered identification of genome-based functional genes related to their attractive dynamic  
71 body appearance (e.g., spine and body size). Besides, as there is such a shortage of amphibian  
72 genomes in Genome 10K project, and it is necessary to analyze other important genomes to  
73 large-scale study the phylogenetic relationships in amphibian [11].

74 In this study, we combined genomic sequencing data from Illumina short reads, PacBio long  
75 reads, and Hi-C data to generate the first chromosomal-level reference genome for the  
76 mustache toad. The completeness and continuity of the genome were comparable with that of  
77 other important amphibian species. The high-quality reference genome generated in this study  
78 will facilitate research on population genetic traits and functional gene identification related  
79 to important characteristics of the mustache toad, ~~which will, in turn, accelerate the~~  
80 ~~development of more efficient body size control techniques and improve the artificial~~  
81 ~~breeding industry for other economically important species.~~

82

### 83 **Sampling and sequencing**

84 A male mustache toad (*V. ailaonica*) with keratinized nuptial spines on its upper jaw was  
85 caught from the Ailao Mountain during the breeding season for sequencing (Figure 1). To  
86 obtain sufficient high-quality DNA for the PacBio Sequel platform (Pacific Biosciences,  
87 USA), the mustache toad was dissected, and fresh liver tissue was used for DNA extraction  
88 using phenol/chloroform extraction. DNA quality was checked by agarose gel electrophoresis,

89 with excellent integrity DNA molecules were obtained. Other tissues, including spines, brain,  
90 stomach, intestine, liver, lung, spleen, blood, and tongue, were snap-frozen in liquid nitrogen  
91 for 10 min and then all these 9 organs/tissues were stored at  $-80^{\circ}\text{C}$  for RNA-seq  
92 experiments subsequent use. Isolated total RNA were used to isolate intact poly(A)+ RNA by  
93 the NEBNext Poly(A) mRNA Magnetic Isolation Module. The mRNA was further  
94 fragmented and randomly primed during the first-strand synthesis by reverse transcription.  
95 This procedure was followed by second-strand synthesis with DNA polymerase I to create  
96 double-stranded cDNA fragments using Transcriptor First Strand cDNA Synthesis Kit  
97 (Roche). In the Hi-C experiments, the collected blood was used for library construction. The  
98 blood sample (150  $\mu\text{l}$ ) was cross-linked for 10 min with formaldehyde (1% final  
99 concentration), after which glycine (0.2 M final concentration) was added for 5 min to stop  
100 the cross-linking process, with the sample then stored until the further analysis.  
101 Extracted DNA was sequenced using the Illumina and PacBio Sequel platforms. The short  
102 reads generated from the Illumina platform were used for estimation of genome size and error  
103 correction of the assembled genome, and the PacBio long reads were used for genome  
104 assembly. To this end, five libraries with insertion lengths of 220 bp or 500 bp were ~~generated~~  
105 ~~from sequenced on the~~ Illumina HiSeq 2500 platform generating 150 bp paired-end reads and  
106 a 20 Kb library was constructed using the PacBio platform according to the manufacturers'  
107 protocols. Finally, we obtained 225.03 Gb of Illumina short reads and 277.15 Gb of PacBio  
108 long reads (Table 1; Additional File: Tables S1 and S2). The average subreads N50 length  
109 reached to 14.78 Kb, providing ultra-long genomic sequences for the following assembly and  
110 analysis (Additional File: Table S2). The RNA-seq samples were obtained by mixing an equal

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

116

#### 117 **Genome characteristics estimation**

118 The Illumina short reads were quality filtered by the following steps: First, the adaptors were  
119 removed from the sequencing reads. Second, read pairs were excluded if any one read had  
120 more than 10% “N”. Third, read pairs with low quality base more than 50% were removed.  
121 Fourth, the PCR duplicates produced during library construction in read pairs were removed.  
122 The filtered reads were used for estimation of genome size and other characteristics. Using  
123 the k-mer method, we calculated the 17-mer depth frequency distribution in the mustache  
124 toad. Genome size was estimated by:  $G = \text{TKN}_{17\text{-mer}} / \text{PKFD}_{17\text{-mer}}$ , where  $\text{TKN}_{17\text{-mer}}$  is the total  
125 kmer number and  $\text{PKFD}_{17\text{-mer}}$  is the peak kmer frequency depth of 17-mer. We estimated a  
126 genome size of 3.52 Gb (peak = 54) and found a heterozygous and repeat sequences peak,  
127 suggesting that the mustache toad genome exhibits complex genome assembly (Figure 2).

128

#### 129 **Genome assembly by PacBio long reads and Hi-C data**

130 Based on 38 single-molecular real-time cells by the PacBio Sequel platform, we generated  
131 277.15 Gb of subreads (Table 1; Additional File: Table S2). The average and N50 length of  
132 subreads were 9.65 Kb and 14.78 Kb, respectively (Additional File: Table S2). All the long

133 reads were used for genome assembly using wtdbg software  
134 (<https://github.com/ruanjue/wtdbg-1.2.8>). As a result, we obtained a 3.95 Gb genome  
135 assembly with a contig N50 length of 739.54 Kb. However, although the size of the genome  
136 assembly was comparable with the estimation k-mer result, it was a slightly larger. This may  
137 be due to the complexity of the mustache toad genome (high heterozygous rate and repeat  
138 sequences). Then the redundant sequences in the genome assembly were removed using  
139 Redundans software (v0.13c) [12] with an identity of 0.7 and overlap of 0.7, resulting in a  
140 genome assembly of 3.58 Gb and contig N50 length of 834.90 Kb. To make sure all the  
141 removed contigs not contain real sequences, we checked the BUSCO result and the mapping  
142 ratio of Illumina reads in the raw genome and the redundancy-filtered genome. Finally, All  
143 these results indicated that the parameters in the redundancy-filtered step are proper in this  
144 study (Additional File: Table S5). To further improve the quality and accuracy of our genome  
145 assembly, the Illumina short reads were used to polish the genome using Pilon software  
146 (RRID:SCR\_014731, v1.21) [13] at the single-base level. The Hi-C data were used to  
147 improve the connection integrity of the contigs (15,899 contigs). We obtained 378.78 Gb of  
148 Hi-C sequencing data, which were first filtered by Hic-Pro (v2.10.0) [14] (Table 1; Additional  
149 File: Table S4), and then mapped to the polished mustache toad genome [15]. The location  
150 and direction of the contigs were determined by 3D *de novo* assembly (3d-DNA) software  
151 (v180419) [16] with default parameters. Most contigs were then successfully clustered and  
152 anchored on 13 groups (Figure 3) [17]. Finally, we obtained the first chromosomal-level high  
153 quality mustache toad assembly (3.53 Gb) with a scaffold N50 length of 412.42 Mb,  
154 providing a solid genomic resource for further study of the mustache toad (Table 2).

155

156 **Genome assembly evaluation**

157 Genome assembly quality is directly related to the accuracy and completeness of  
158 protein-coding gene prediction. Therefore, we evaluated the assembled mustache toad  
159 genome using three methods. First, the assembled genome was compared with the core gene  
160 set in BUSCO software (RRID:SCR\_015008, v2.0) [18]. We found 245 (80.8%) and 833  
161 (85.1%) conserved core genes in the mustache toad genome using the eukaryote and  
162 metazoan databases, respectively (Table 3). When we further considered the fragmented  
163 BUSCO genes found in the genome, there were 272 (89.7%) and 881 (90.1%) conserved core  
164 genes found in the eukaryote and metazoan databases, respectively. ~~(Table 3)~~ This results  
165 indicated that the assembled mustache toad genome is comparable with published amphibian  
166 genomes (Table 3). Second, we aligned all filtered short reads generated from the Illumina  
167 platform to the genome using BWA software (RRID:SCR\_010910, v0.7.12) [19] and found  
168 1,778 million clean reads that could be mapped to the genome, accounting for 97.78% of total  
169 clean reads (Additional File: Table ~~S5S6~~). Third, the RNA-seq reads were *de novo* assembled  
170 using Bridger software (RRID:SCR\_017039, version: r2014-12-01) [20], with redundant  
171 transcripts removed by TGICL [21], resulting in 19,876 transcripts (Additional File: Table  
172 ~~S6S7~~). These transcripts were then aligned to the genome, with 17,878 transcripts (89.95%)  
173 found in the assembled genome and 94.52% of transcripts longer than 1 Kb (Additional File:  
174 Table ~~S7S8~~). Besides, we analyzed the N50 length and BUSCO results and found that the  
175 mustache toad genome was comparable with other published amphibian genomes (Tables ~~2-~~  
176 ~~and Table 3-4; Additional File: Table S8).~~ These results indicate that our assembled mustache

177 ~~toad genome exhibited high completeness and accuracy.~~  
178 The GC distribution of the mustache toad and ~~vertebrate closely-related~~ species was  
179 calculated using the slide window method. Results showed that their GC distributions were  
180 similar to each other, with an average GC content of 43.68% in the mustache toad and 36.60%  
181 to 44.49% in other species (Additional File: Figure S1). ~~These results indicate that our~~  
182 ~~assembled mustache toad genome exhibited high completeness and accuracy.~~

183

#### 184 **Genome annotation**

185 We used Tandem Repeat Finder (TRF, v4.04) [22] to identify repetitive elements and  
186 RepeatModeler software (RRID:SCR\_015027, v1.0.4) to detect transposable elements (TEs)  
187 in the mustache toad genome. Then, the *de novo* library of repeats produced by  
188 RepeatModeler analysis and the rebase (RepBase16.02) database were then used for  
189 RepeatMasker (RRID:SCR\_012954, version: open-4.0) [23] analysis to identify homologous  
190 repeats. RepeatProteinMask was used to query the TE protein database at the protein level.  
191 Lastly, we identified 2.45 Gb of repeat sequences, which accounts for 69.48% of the  
192 estimated genome size (Additional File: Table S9). Among these repeat sequences, 60.87%  
193 (2.15 Gb) was predicted by the *de novo* method (Table 45).

194 After repeat sequence annotation, we masked all repeats, except for the tandem repeat  
195 sequences, for protein-coding gene annotation. Augustus software (RRID:SCR\_008417,  
196 v2.5.5) [24] was used to *de novo* predict coding genes using a zebrafish (*Danio rerio*) dataset  
197 as the train species. For the homology-based method, protein sequences of ~~closely-~~  
198 ~~related chordate~~ species, including *D. rerio* (GCF\_000002035.6) [25], *Nanorana parkeri*

199 (GCF\_000935625.1) [26], *Homo sapiens* (GCF\_000001405.38) [27], *Gallus gallus*  
200 (GCF\_000002315.5) [28], *Pelodiscus sinensis* (GCF\_000230535.1) [29], *Xenopus laevis*  
201 (GCF\_001663975.1) [30], and *Petromyzon marinus* [31]  
202 (<https://genomes.stowers.org/organism/Petromyzon/marinus>), were downloaded and aligned  
203 against the mustache toad genome using the TBLASTN module (BLAST version: 2.3.0). The  
204 transcripts assembled by RNA-seq reads were first translated into amino acids and then  
205 aligned to the genome using TBLASTN software for gene annotation. EVidenceModeler  
206 (RRID:SCR\_014659, version: r2012-06-25) [32] was used to integrate results from the three  
207 methods, and genes with poor transcriptome evidence support were filtered out. Finally,  
208 26,227 high-quality protein-coding genes were predicted in the mustache toad genome.  
209 Moreover, the distributions of mRNA, CDS, exon, and intron lengths were comparable with  
210 closely related species (Figure 4).  
211 Gene functional annotation can help to elucidate gene function. Thus, we aligned all 26,227  
212 protein-coding genes to protein databases, including InterProScan, KEGG, SwissProt, and  
213 TrEMBL. Results showed that most obtained genes could be annotated in these functional  
214 databases (Table 56).

215

#### 216 **Phylogenetic tree and divergence time analysis**

217 To reveal the phylogenetic relationships of the mustache toad with other closely related  
218 species, we identified the single-copy genes among the species. First, protein sequences,  
219 including those of *D. rerio* (GCF\_000002035.6) [25], *N. parkeri* (GCF\_000935625.1) [26], *H.*  
220 *sapiens* (GCF\_000001405.38) [27], *G. gallus* (GCF\_000002315.5) [28], *Anolis carolinensis*



221 (GCF\_000090745.1) [33], *Xenopus tropicalis* (GCF\_000004195.3) [30], *Rhinella marina*  
222 (GigaDB) [34], *Rana catesbeiana* (GCA\_002284835.2) [35], *Ambystoma mexicanum*  
223 (www.axolotl-omics.org) [36,37], and *Alligator sinensis* (GCF\_000455745.1) [3738], were  
224 downloaded from NCBI and the longest transcript of each gene in each species was selected.  
225 The BLASTP program (BLAST version: 2.2.24) was then used to align these protein  
226 sequences among these 11 species (including the mustache toad) with an e-value of 1e-5. The  
227 homolog relationships (including ortholog and paralog) were then determined using  
228 OrthoMCL software (v1.4) [3839]. Genes with only one copy in the species were identified as  
229 single-copy genes. In total, 238 genes were identified (Figure 5), with the detailed results of  
230 gene family statistics shown in the supplementary information (Additional File: Table S10).  
231 The 238 single-copy genes were aligned using MUSCLE software (RRID:SCR\_011812,  
232 v3.8.31) [3940, 4041] and concatenated to supergenes for maximum-likelihood-based  
233 phylogenetic analyses. We performed phylogenetic analysis, with zebrafish as the outgroup,  
234 using RAxML software (RRID:SCR\_006086, v8.2.3) [4442] with the parameter -m for  
235 PROTGAMEAUTO. Results indicated that the mustache toad exhibited a close  
236 relationship with the ancestor of the marine toad (*R. marina*), bull frog (*R. catesbeiana*), and  
237 Tibetan frog (*N. parkeri*), with topological relationships in other clades found to be the same  
238 as reported previously (Figure 6). To further investigate the divergence time of these species,  
239 especially toad and frogs, the MCMCTREE model in PAML software (RRID:SCR\_014932,  
240 v4.8) [4243] was used with three datasets (four-fold degenerate sites (4dTVs); the first-codon  
241 sites; the second-codon sites) extracted from the single-copy genes as the input file. Fossil  
242 records were downloaded from the TIMETREE website (www.timetree.org) and used to

243 calibrate the results. We found that the [results from the three different datasets are very](#)  
244 [similar and the](#) mustache toad diverged with the common ancestor of the marine toad, bull  
245 frog, and Tibetan frog about [194.8206.1](#) million years ago (Figure 6: [Additional File: Figures](#)  
246 [S2 and S3](#)).

247

### 248 **Gene family expansion and contraction**

249 We performed gene family expansion and contraction analysis using CAFÉ software  
250 (RRID:SCR\_005983, v4.0) [[4344](#)], and found [349-201](#) and [2,607326](#) expanded and  
251 contracted gene families in the mustache toad ( $P < 0.05$ ), respectively. Using the GO/KEGG  
252 databases, functional enrichment analysis of the expanded gene families found [174-210](#) GO  
253 terms (adjusted  $P$ -value  $< 0.05$ ) and [29](#) KEGG pathways ( $q$ -value  $< 0.05$ ) to be significantly  
254 enriched (Additional File: Tables S11 and S12). The expanded gene families were mainly  
255 related to metabolic processes, intermediate filament terms, enzyme activities, and immune  
256 terms. For example, the cellular metabolic process (adjusted  $P$ -value = [2.376.06E-14](#)),  
257 intermediate filament (adjusted  $P$ -value = [3.923.42E-153](#)), keratin filament (adjusted  $P$ -value  
258 = [1.662.94E-132](#)), endoribonuclease activity (adjusted  $P$ -value = [1.069.19E-087](#)), and  
259 immune response ( $q$ -value = [4.818.36E-036](#)) were enriched (Additional File: Tables S11 and  
260 S12). In addition, for the contracted gene families, [226-220](#) GO terms (adjusted  $P$ -value  $<$   
261  $0.05$ ) and [11-9](#) KEGG pathways ( $q$ -value  $< 0.05$ ) were enriched, respectively (Additional File:  
262 Tables S13 and S14). These enriched terms were mainly involved in ion binding and  
263 transporter activity, including neurotransmitter transporter activity (adjusted  $P$ -value =  
264 [4.201.89E-1109](#)), sodium ion transmembrane transporter activity (adjusted  $P$ -value =

265 ~~1.553.33E-068~~), and secondary active transmembrane transporter activity (adjusted  $P$ -value =  
266 ~~6.371.86E-08~~) (Additional File: Tables S13 and S14). Thus, these biological processes may be  
267 related to the special characteristics of the mustache toad.

268

### 269 **Relative evolutionary rate of species**

270 The evolutionary rate of species can reflect its evolution history and status. The relative  
271 evolutionary rate of the mustache toad to other closely related species was analyzed using  
272 LINTRE [4445] and MEGA (RRID:SCR\_000667, v7.0.26) softwares. Two-cluster analysis  
273 was applied to test the molecular evolution of multiple sequences in a phylogenetic context  
274 based on the concatenated supergenes (protein sequences) using *tpcv* (a module in LINTRE  
275 software). The concatenated supergenes were also used for Tajima's relative rate test. We used  
276 zebrafish as the outgroup in both methods, and found that, except for the axolotl, the  
277 mustache toad had a relatively faster evolutionary rate than its closely related species (e.g.,  
278 *X.tropicalis*, *R. marina*, *R. catesbeiana*, and *N. parkeri*) (Additional File: Tables S15 and S16).  
279 The crocodile had a slower evolutionary rate relative to its closely related species, and this  
280 result is consistent with previous study [4546] (Additional File: Tables S15 and S16).

281

### 282 **Conclusions**

283 Using Illumina, PacBio, and Hi-C sequencing technologies, we reported on the first  
284 chromosomal-level genome assembly of the mustache toad. We successfully annotated 26,227  
285 protein-coding genes by integrating the results of three different methods. The phylogenetic  
286 analysis results indicated that the mustache toad has a close relationship with the marine toad,

287 bull frog, and Tibetan frog, and diverged at [194.8206.1](#) MYA with their common ancestor.

288 Analysis showed that the mustache toad had a faster evolutionary rate relative to most other  
289 closely related species studied. Expansion and contraction of gene family analysis identified  
290 several biological processes and pathways, such as metabolism and intermediate filaments,  
291 suggesting that these terms may relate to the special adaptations of the mustache toad to its  
292 habitat.

293

#### 294 **Availability of supporting data**

295 The raw sequencing data were deposited in the NCBI database under accession number  
296 PRJNA523649. The genome assembly and annotation results are available via the  
297 GigaScience repository GigaDB.

298

#### 299 **Additional files**

300 Figure S1: The GC content in these genomes.

301 [Figure S2: The divergence time of these species \(using the first-codon sites\).](#)

302 [Figure S3: The divergence time of these species \(using the second-codon sites\).](#)

303 Table S1: The statistics of Illumina sequencing clean data.

304 Table S2: The statistics of PacBio Sequel sequencing data.

305 Table S3: The statistics of RNA-seq clean data.

306 Table S4: The statistics of Hi-C sequencing clean data.

307 [Table S5. The comparison of BUSCO and Illumina reads mapping results in these two](#)

308 [genome versions.](#)

309 Table S65: The statistics of Illumina reads mapping ratio to the assembled genome.  
310 Table S76: The statistics of assembled transcripts by Bridger software.  
311 Table S87: The statistics of transcripts mapping ratio to the assembled genome.  
312 ~~Table S8: The quality statistics of several published amphibian genomes.~~  
313 Table S9: The statistics of the annotated repeat sequences in our assembled genome.  
314 Table S10: The statistics of gene family among these species.  
315 Table S11: The GO enrichment analysis of expanded gene families.  
316 Table S12: The KEGG enrichment analysis of expanded gene families.  
317 Table S13: The GO enrichment analysis of contracted gene families.  
318 Table S14: The KEGG enrichment analysis of contracted gene families.  
319 Table S15: Two cluster analysis of mustache toad and other species.  
320 Table S16: The relative evolutionary rate of mustache toad and other species analyzed by  
321 Tajima's Test.

322

### 323 **Abbreviations**

324 BLAST: Basic Local Alignment Search Tool; BUSCO: Benchmarking Universal Single-Copy  
325 Orthologs; BWA: Burrows-Wheeler Aligner; CDS: Coding DNA Sequence; DNA:  
326 Deoxyribonucleic Acid; GO: Gene Ontology; Hi-C: High-throughput chromosome  
327 conformation capture; KEGG: Kyoto Encyclopedia of Genes and Genomes; MHC: Major  
328 Histocompatibility Complex; NCBI: National Center for Biotechnology Information; NR:  
329 Non-Redundant Protein Sequence Database; PCR: Polymerase Chain Reaction; RNA:

330 Ribonucleic Acid; RNA-seq; RNA sequencing.

331

332 **Conflicts of interest**

333 The authors declare that they have no competing interests.

334

335 **Funding**

336 This work was supported by the National Key Research and Development Program of China

337 (NO. 2017YFC0505202) and National Natural Science Foundation of China (NO.

338 NSFC-30270175; NO. NSFC-30870278; NO. NSFC-31372165).

339

340 **Author contributions**

341 D.R. designed the project; D.R. and D.Z. collected the samples; Y.L. and Y.R. estimated the

342 genome size and assembled the genome; Y.L. polished the assembled genome and employed

343 the Hi-C analysis; H.J. performed the genome annotation; Y.L. and Z.W. assessed the quality

344 of the genome assembly; Y.L. and Y.R. constructed the phylogenetic tree and determined

345 divergence time, relative evolutionary rate of species, and expansion and contraction of gene

346 families. Y.L., D.R., Y.R., and X.L. wrote the manuscript.

347

348 **References**

- 349 1. Liang F and Changyuan Y. Amphibians of China. Science Press: Beijing; 2016.  
350 2. Matsui M, Hamidy A, Murphy RW, Khonsue W, Yambun P, Shimada T, et al.  
351 Phylogenetic relationships of megophryid frogs of the genus *Leptobrachium*  
352 (Amphibia, Anura) as revealed by mtDNA gene sequences. *Molecular Phylogenetics*  
353 & Evolution. 2010;56 1:259-72.  
354 3. Matsui M. A New *Leptobrachium* (*Vibrissaphora*) from Laos (Anura: Megophryidae).  
355 *Current Herpetology*. 2013;32 2:182-9.

- 356 4. Liu C HS, Zhao EJAHS, Chengdu, Old Ser. Preliminary study of genus *Vibrissaphora*  
357 (Amphibia: Salientia) and discussion on problems of amphibian classification.  
358 *Copeia*. 1980;3:1-9.
- 359 5. Rao DQ and Wilkinson JA. Phylogenetic relationships of the mustache toads inferred  
360 from mtDNA sequences. *Molecular Phylogenetics & Evolution*. 2008;46 1:61-73.
- 361 6. Zheng Y, Li S and Fu J. A phylogenetic analysis of the frog genera *Vibrissaphora* and  
362 *Leptobrachium*, and the correlated evolution of nuptial spine and reversed sexual size  
363 dimorphism. *Molecular Phylogenetics & Evolution*. 2008;46 2:695-707.
- 364 7. Zheng Y, Deng D, Li S and Fu J. Aspects of the breeding biology of the Emei  
365 mustache toad (*Leptobrachium boringii*): Polygamy and paternal care.  
366 *Amphibia-Reptilia*. 2010;31 2:183-94.
- 367 8. Zhang W, Guo Y, Li J, Huang L, Kazitsa EG and Wu H. Transcriptome analysis  
368 reveals the genetic basis underlying the seasonal development of keratinized nuptial  
369 spines in *Leptobrachium boringii*. *Bmc Genomics*. 2016;17 1:978.
- 370 9. Zheng Y, Rao D, Murphy RW and Zeng X. Reproductive Behavior and Underwater  
371 Calls in the Emei Mustache Toad, *Leptobrachium boringii*. *Asian Herpetological*  
372 *Research*. 2011;02 4:199-215.
- 373 10. Hudson CM, Xianjin HE and Jinzhong FU. Keratinized Nuptial Spines Are Used for  
374 Male Combat in the Emei Moustache Toad (*Leptobrachium boringii*). *Asian*  
375 *Herpetological Research*. 2011;02 3:142-8.
- 376 11. Genome 10K: a proposal to obtain whole-genome sequence for 10,000 vertebrate  
377 species. *The Journal of heredity*. 2009;100 6:659-74. doi:10.1093/jhered/esp086.
- 378 12. Prysacz LP and Gabaldón T. Redundans: an assembly pipeline for highly  
379 heterozygous genomes. *Nucleic Acids Research*. 2016;44 12:e113-e.
- 380 13. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an  
381 integrated tool for comprehensive microbial variant detection and genome assembly  
382 improvement. *PLoS One*. 2014;9 11:e112963. doi:10.1371/journal.pone.0112963.
- 383 14. Servant N, Varoquaux N, Lajoie BR, Viara E, Chen CJ, Vert JP, et al. HiC-Pro: an  
384 optimized and flexible pipeline for Hi-C data processing. *Genome Biology*. 2015;16  
385 1:259.
- 386 15. Durand N, Shamim M, Machol I, Rao SP, Huntley M, Lander E, et al. Juicer Provides  
387 a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Systems*.  
388 2016;3 1:95-8.
- 389 16. Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, et al. De  
390 novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length  
391 scaffolds. *Science*. 2017;356 6333:92.
- 392 17. Wilkinson JA. A New Species of the Genus *Vibrissaphora* (Anura: Megophryidae)  
393 from Yunnan Province, China. *Herpetologica*. 2006;62 1:90-5.
- 394 18. Simão FA, Waterhouse RM, Panagiotis I, Kriventseva EV and Zdobnov EM.  
395 BUSCO: assessing genome assembly and annotation completeness with single-copy  
396 orthologs. *Bioinformatics*. 2015;31 19:3210-2.
- 397 19. Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler  
398 transform. 2009.
- 399 20. Chang Z, Li G, Liu J, Zhang Y, Ashby C, Liu D, et al. Bridger: a new framework for

- 400 de novo transcriptome assembly using RNA-seq data. *Genome Biology*. 2015;16  
401 1:30.
- 402 21. Pertea G, Huang X, Liang F, Antonescu V, Sultana R, Karamycheva S, et al. TIGR  
403 Gene Indices clustering tools (TGICL): a software system for fast clustering of large  
404 EST datasets. *Bioinformatics*. 2003;19 5:651-2.
- 405 22. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic  
406 Acids Res*. 1999;27 2:573-80.
- 407 23. Bedell JA, Korf I, and Gish W, . MaskerAid: a performance enhancement to  
408 RepeatMasker. *Bioinformatics*. 2000;16 11:1040-1.
- 409 24. Stanke M and Waack S. Gene prediction with a hidden Markov model and a new  
410 intron submodel. *Bioinformatics*. 2003;19 suppl\_2:215--25.
- 411 25. Kerstin H, Clark MD, Torroja CF, James T, Camille B, Matthieu M, et al. The  
412 zebrafish reference genome sequence and its relationship to the human genome.  
413 *Nature*. 2013.
- 414 26. Yan-Bo S, Zi-Jun X, Xue-Yan X, Shi-Ping L, Wei-Wei Z, Xiao-Long T, et al.  
415 Whole-genome sequence of the Tibetan frog *Nanorana parkeri* and the comparative  
416 evolution of tetrapod genomes. *Proc Natl Acad Sci U S A*. 2015;11 112:E1257-E62.
- 417 27. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial  
418 sequencing and analysis of the human genome. *Nature*. 2001;409 6822:860-921.  
419 doi:10.1038/35057062.
- 420 28. Sequence and comparative analysis of the chicken genome provide unique  
421 perspectives on vertebrate evolution. *Nature*. 2004;432 7018:695-716.  
422 doi:10.1038/nature03154.
- 423 29. Wang Z, Pascual-Anaya J, Zadissa A, Li W, Niimura Y, Huang Z, et al. The draft  
424 genomes of soft-shell turtle and green sea turtle yield insights into the development  
425 and evolution of the turtle-specific body plan. *Nat Genet*. 2013;45 6:701-6.  
426 doi:10.1038/ng.2615.
- 427 30. Session AM, Uno Y, Kwon T, Chapman JA, Toyoda A, Takahashi S, et al. Genome  
428 evolution in the allotetraploid frog *Xenopus laevis*. *Nature*. 2016;538 7625:336-43.  
429 doi:10.1038/nature19840.
- 430 31. Smith JJ and Timoshevskaya N. The sea lamprey germline genome provides insights  
431 into programmed genome rearrangement and vertebrate evolution. 2018;50 2:270-7.  
432 doi:10.1038/s41588-017-0036-1.
- 433 32. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated  
434 eukaryotic gene structure annotation using EVIDENCEModeler and the Program to  
435 Assemble Spliced Alignments. *Genome Biology*. 2008;9 1:R7.
- 436 33. Jessica AL, Federica DP, Manfred G, Christina W, Lesheng K, Evan M, et al. The  
437 genome of the green anole lizard and a comparative analysis with birds and  
438 mammals. *Nature*. 2011.
- 439 ~~34. Edwards RJ, Tuipulotu DE, Amos TG, O'Meally D, Richardson MF, Russell TL, et al.  
440 Draft genome assembly of the invasive cane toad, *Rhinella marina*. *GigaScience*.  
441 2018;7 9 doi:10.1093/gigascience/giy095.~~
- 442 ~~34. Nowoshilow S, Schloissnig S, Fei JF, Dahl A, Pang AWC, Pippel M, et al. The  
443 axolotl genome and the evolution of key tissue formation regulators. *Nature*.~~



444 [2018;554 7690.](#)

445 35. Hammond SA, Warren RL, Vandervalk BP, Kucuk E, Khan H, Gibb EA, et al. The  
446 North American bullfrog draft genome provides insight into hormonal regulation of  
447 long noncoding RNA. 2017;8 1:1433. doi:10.1038/s41467-017-01316-7.

448 36. Smith JJ, Timoshevskaya N, Timoshevskiy VA, Keinath MC, Hardy D and Voss SR.  
449 A chromosome-scale assembly of the axolotl genome. *Genome Res.* 2019;29  
450 2:317-24. doi:10.1101/gr.241901.118.

451 [37. Nowoshilow S, Schloissnig S, Fei JF, Dahl A, Pang AWC, Pippel M, et al. The](#)  
452 [axolotl genome and the evolution of key tissue formation regulators. \*Nature.\*](#)  
453 [2018;554 7690.](#)

454 [3738.](#) Wan QH, Pan SK, Hu L, Zhu Y, Xu PW, Xia JQ, et al. Genome analysis and signature  
455 discovery for diving and sensory properties of the endangered Chinese alligator. *Cell*  
456 *research.* 2013;23 9:1091-105. doi:10.1038/cr.2013.104.

457 [3839.](#) Li L, Jr SC and Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic  
458 genomes. *Genome Research.* 2003;13 9:2178-89.

459 [3940.](#) Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high  
460 throughput. *Nucleic Acids Res.* 2004;32 5:1792-7. doi:10.1093/nar/gkh340.

461 [4041.](#) Edgar RC. MUSCLE: a multiple sequence alignment method with reduced time and  
462 space complexity. *Bmc Bioinformatics.* 2004.

463 [4142.](#) Alexandros S. RAxML version 8: a tool for phylogenetic analysis and post-analysis  
464 of large phylogenies. *Bioinformatics.* 2014;30 9:1312-3.

465 [4243.](#) Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular biology*  
466 *and evolution.* 2007;24 8:1586-91. doi:10.1093/molbev/msm088.

467 [4344.](#) Tijn DB, Nello C, Demuth JP and Hahn MW. CAFE: a computational tool for the  
468 study of gene family evolution. *Bioinformatics.* 2006;22 10:1269-71.

469 [4445.](#) Takezaki N, Rzhetsky A, and Nei M. Phylogenetic test of the molecular clock and  
470 linearized trees. *Molecular Biology & Evolution.* 1995;12 5:823-33.

471 [4546.](#) Green RE, Braun EL, Joel A, Dent E, Ngan N, Glenn H, et al. Three crocodylian  
472 genomes reveal ancestral patterns of evolution among archosaurs. *Science.* 2014;346  
473 6215:1254449.

474 **Tables and Figures**

475

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

477

478 **Table 2: Assembly statistics of 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	

479 Note: Statistics of genome assembly. Wtdbg contig was the genome assembled by wtdbg and  
 480 2-round pilon error-correction. Hi-C scaffold was the genome finished by Hi-C assembly.

481

482 **Table 3: The BUSCO results of the mustache toad and other amphibian genomes.**

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

483 Note: Both “eukaryote” and “metazoan” are two core gene sets in BUSCO database.

484 **Table 3: The BUSCO results of the mustache toad genome.**

Library	eukaryota	metazoa
Complete BUSCOs (C)	245	833
Complete and single-copy BUSCOs (S)	237	818
Complete and duplicated BUSCOs (D)	8	15
Fragmented BUSCOs (F)	27	48
Missing BUSCOs (M)	34	97
Total BUSCO groups searched	303	978
Summarize	80.8%	85.1%

485

486

487

488

489

490 **Table 4: The quality statistics of several published amphibian genomes.**

Formatted: Left

Formatted: Left

Formatted: Left

Formatted: Left

Formatted: Left

Formatted: Left

Formatted: Left

Formatted: Left

Formatted: Left

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%

491

492 **Table 45: The statistics of *de novo* annotated repeat sequences in mustache toad genome.**

Type	Length (bp)	Percentage in genome (%)
DNA	350,793,270	9.94 <del>3777</del>
LINE	297,954,803	8.4 <del>5989</del>
SINE	11,009,363	0.31 <del>2077</del>
LTR	307,317,539	8.71 <del>1390</del>
Other	43,867,330	1.24 <del>3487</del>
Satellite	9,696,790	0.27 <del>4870</del>
Simple repeat	125,397,072	3.55 <del>4574</del>
Unknown	1,114,326,962	31.5 <del>8732059</del>
Total	2,147,505,764	60.87 <del>4369</del>

493

494 **Table 56: The functional annotation results of protein-coding genes in mustache toad.**

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

495

496 **Figure 1: The mustache toad, *Vibrissaphora ailaonica*.** (A) The adult male individual with  
497 spines in the upper jaw. (B) The adult female individual. (C) The adult male individual during  
498 the fall off process of spines in the upper jaw. (D) The adult male individual without spines  
499 (after fall off process of spines) in the upper jaw. (E) The body size of mustache toad in side  
500 view, male (left) and female (right). (F) The body size of mustache toad in top view, male (left)  
501 and female (right).

502

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

504

505 **Figure 3: The circos graph showing genome characteristics.** From outer circle to inner ring  
506 are: gene distribution, tandem repeats (TR), long tandem repeats (LTR), long interspersed  
507 nuclear elements (LINE), short interspersed nuclear elements (SINE), and GC content.

508

509 **Figure 4: The length distributions of annotated protein-coding genes in these species.**

510

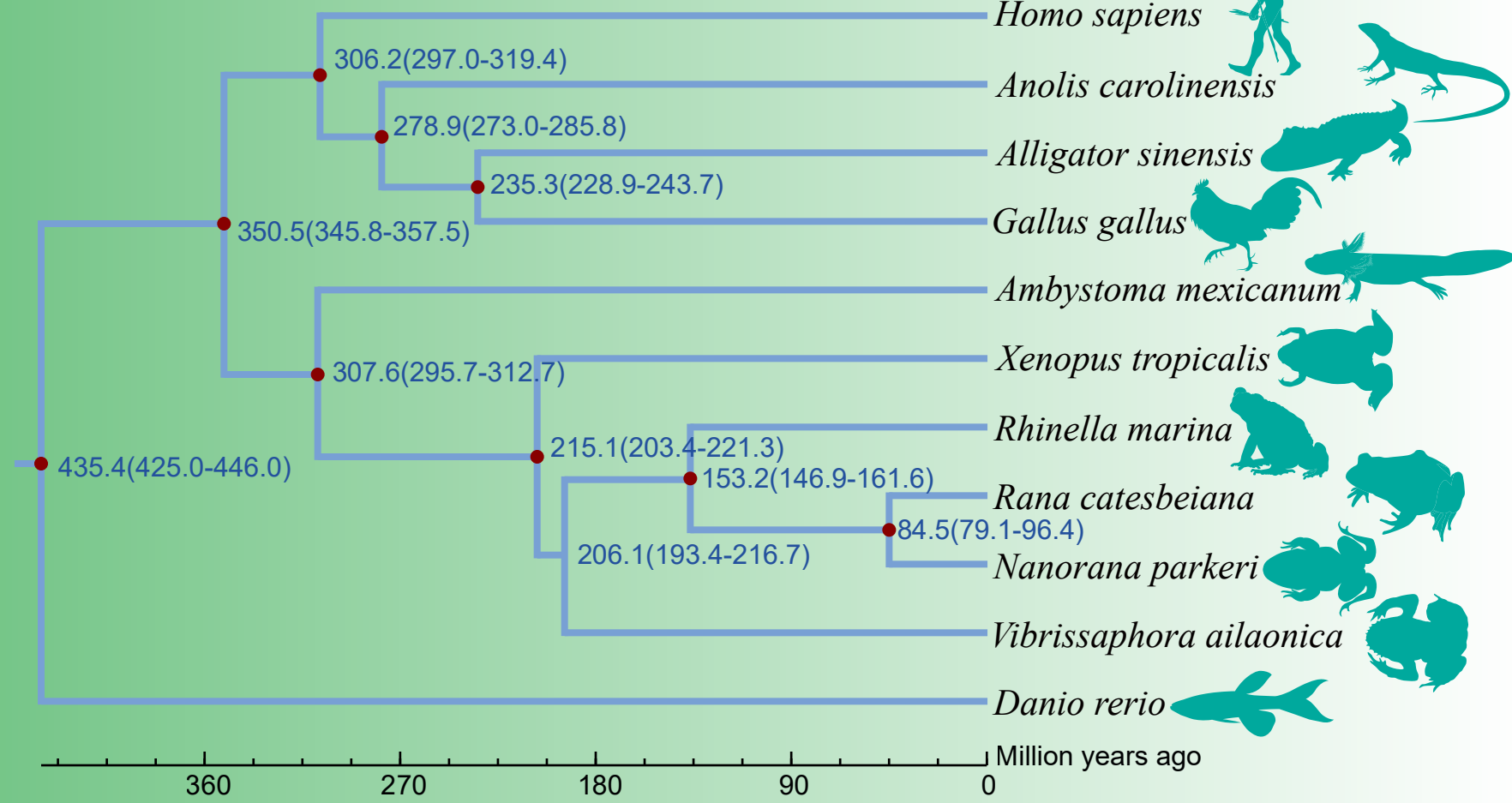
511 **Figure 5: The statistics of gene family among these 11 species.**

512

513 **Figure 6: The phylogenetic relationships among these species.** The blue numbers represent  
514 divergence time. The red dot represents the fossil record used in the node.  
515

Figure 6

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





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

