

# GigaScience

## A high-quality genome assembly of the endangered golden snub-nosed monkey (*Rhinopithecus roxellana*) --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-19-00030	
<b>Full Title:</b>	A high-quality genome assembly of the endangered golden snub-nosed monkey ( <i>Rhinopithecus roxellana</i> )	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	National Natural Science Foundation of China (31622053)	Dr. Xiao-Guang Qi
<b>Abstract:</b>	<p><b>Background:</b> The golden snub-nosed monkey (<i>Rhinopithecus roxellana</i>), is an endangered colobine monkey species endemic to China. This species has several distinctive traits, and it is an ideal model for analysing evolutionary development of the social structure due to its unique social organization. Although there has been reported a genome assembly of the subspecies <i>R. roxellana hubeiensis</i>, the assembly is incomplete and fragmented due to employing short reads sequencing technology. This drawback may lose information, such as structural variation and repeat sequences which are important for understanding this endangered species. Therefore, to have a better understanding of evolutionary history and genetic-specific signatures, a high-quality reference genome of the taxon is need.</p> <p><b>Findings:</b> To obtain a high-quality chromosome assembly of <i>R. roxellana qinlingensis</i>, we combined a total of five techniques including Pacific Bioscience's single-molecule real-time sequencing, Illumina's paired-end sequencing, BioNano optical maps, 10X Genomics link-reads and high-throughput chromosome conformation capture. The results indicate the assembled genome is about 3.04 Gb with a contig N50 of 5.72 Mbp and a scaffold N50 of 144.56 Mbp, which have made a 10-fold improvement compared to past published. It is shown that a total of 22497 protein coding genes were predicted, of which 22053 were functionally annotated. Moreover, gene family analysis shows that 993 and 2745 gene families are expanded and contracted in the <i>R. roxellana</i> genome, respectively.</p> <p><b>Conclusion:</b> We present the updated high-quality genome assembly of <i>R. roxellana</i> with superior continuity and accuracy. The assembled genome can be used as reference for future genetic studies of the species. Also, the updated genome assembly may contribute to our comprehensive understanding of the species, which is particularly helpful in the conservation of this endangered species. Furthermore, such genome with superior continuity and accuracy can provide a new standard reference for Colobine primates.</p>	
<b>Corresponding Author:</b>	Xiao-Guang Qi  CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>		
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Xiao-Guang Qi	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Xiao-Guang Qi Lu Wang Jinwei Wu Xiaomei Liu	

	Dandan Di
	Yuhong Liang
	Yifei Feng
	Baoguo Li
<b>Order of Authors Secondary Information:</b>	
<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
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or</p>	Yes

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               **A high-quality genome assembly of the endangered golden snub-nosed monkey**

2                                    **(*Rhinopithecus roxellana*)**

3  
4 3  
5  
6 4 Lu Wang<sup>1,†</sup>, Jinwei Wu<sup>1,†</sup>, Xiaomei Liu<sup>1</sup>, Dandan Di<sup>1</sup>, Yuhong Liang<sup>1</sup>, Yifei Feng<sup>1</sup>, Baoguo

7  
8  
9 5 Li<sup>1,2</sup>, Xiao-Guang Qi<sup>1,\*</sup>

10  
11  
12  
13 6  
14  
15  
16  
17 7 <sup>1</sup> Shaanxi Key Laboratory for Animal Conservation, College of Life Sciences, Northwest  
18  
19  
20 8 University, Xi'an, 710069, China.

21  
22  
23  
24 9 <sup>2</sup> Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences,  
25  
26  
27 10 Kunming, 650223, China.

28  
29  
30  
31 11 **\*Correspondence address.** Xiao-Guang Qi, E-mail: qixg@nwu.edu.cn

32  
33  
34  
35 12 <sup>†</sup>These authors contributed equally to this work.

1 13 **ABSTRACT**

2  
3  
4 14  
5  
6  
7 15 **Background:** The golden snub-nosed monkey (*Rhinopithecus roxellana*), is an endangered  
8  
9  
10 16 colobine monkey species endemic to China. This species has several distinctive traits, and it is  
11  
12  
13  
14 17 an ideal model for analysing evolutionary development of the social structure due to its unique  
15  
16  
17 18 social organization. Although there has been reported a genome assembly of the subspecies *R.*  
18  
19  
20 19 *roxellana hubeiensis*, the assembly is incomplete and fragmented due to employing short reads  
21  
22  
23 20 sequencing technology. This drawback may lose information, such as structural variation and  
24  
25  
26 21 repeat sequences which are important for understanding this endangered species. Therefore, to  
27  
28  
29 22 have a better understanding of evolutionary history and genetic-specific signatures, a high-  
30  
31  
32  
33 23 quality reference genome of the taxon is need.

34  
35  
36 24 **Findings:** To obtain a high-quality chromosome assembly of *R. roxellana qinlingensis*, we  
37  
38  
39 25 combined a total of five techniques including Pacific Bioscience's single-molecule real-time  
40  
41  
42 26 sequencing, Illumina's paired-end sequencing, BioNano optical maps, 10X Genomics link-  
43  
44  
45 27 reads and high-throughput chromosome conformation capture. The results indicate the  
46  
47  
48 28 assembled genome is about 3.04 Gb with a contig N50 of 5.72 Mbp and a scaffold N50 of  
49  
50  
51 29 144.56 Mbp, which have made a 10-fold improvement compared to past published. It is shown  
52  
53  
54  
55 30 that a total of 22497 protein coding genes were predicted, of which 22053 were functionally  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 31 annotated. Moreover, gene family analysis shows that 993 and 2745 gene families are expanded  
2  
3  
4 32 and contracted in the *R. roxellana* genome, respectively.  
5  
6

7 33 **Conclusion:** We present the updated high-quality genome assembly of *R. roxellana* with  
8  
9  
10 34 superior continuity and accuracy. The assembled genome can be used as reference for future  
11  
12  
13 35 genetic studies of the species. Also, the updated genome assembly may contribute to our  
14  
15  
16 36 comprehensive understanding of the species, which is particularly helpful in the conservation  
17  
18  
19 37 of this endangered species. Furthermore, such genome with superior continuity and accuracy  
20  
21  
22 38 can provide a new standard reference for Colobine primates.  
23  
24  
25

26 39  
27  
28  
29 40 **Keywords:** high-quality; *Rhinopithecus roxellana*; genome assembly; annotation; BioNano  
30  
31  
32 41 optical maps  
33  
34  
35

36 42  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 43 **Data Description**

### 45 **Background information**

46 Snub-nosed monkeys (*Rhinopithecus*) consist of five endangered species narrowly restricted  
47 to China and Vietnam [1]. Among those, the golden snub-nosed monkey (*Rhinopithecus*  
48 *roxellana*) is also referred to as the Sichuan snub-nosed monkey, with the northernmost  
49 distribution of all Asian colobine species, found only in three isolated regions (Sichuan and  
50 Gansu, Shaanxi and Hubei provinces) in central and northwest China [2, 3]. This species is  
51 characterized by several distinctive traits, such as golden fur, blue facial colour, odd-shaped  
52 nose, more folivorous, most striking unique social system with multilevel societies, a rare and  
53 complex system that is found only in a few mammal species, including human beings [4].  
54 Therefore, *R. roxellana* is an ideal model for the studies analysing evolutionary development  
55 of the social structure in primates and understanding the behaviour patterns of human society  
56 in social-anthropology.

57 As a research hotspot, studies on *R. roxellana* have investigated various aspects [4-6].

58 Recently, genomic analysis offered a powerful tool and has successfully been employed to  
59 underlie the molecular evolution of several groups [7-9]. According to the morphological  
60 variation and distribution difference, *R. roxellana* can differentiate into three subspecies:  
61 *Rhinopithecus roxellana roxellana* from Minshan mountains of Sichuan and Gansu province,

1 62 *R. r. Qinlingensis* from Qinling mountains of Shaanxi province, *R. r. hubeiensis* from  
2  
3  
4 63 shennongjia mountains [3]. Up to now, the best genome assembly of *R. roxellana* was  
5  
6  
7 64 published in 2014 [10], which was derived from short reads sequencing on Illumina HiSeq  
8  
9  
10 65 2000 platform. Based on this achievement, studies on its folivorous dietary adaptations and the  
11  
12  
13 66 evolutionary history of *R. roxellana* have been conducted [10-12]. Despite such progress, the  
14  
15  
16 67 information including structural variation and repeat sequences was largely absent or unreliable  
17  
18  
19  
20 68 due to the incomplete and fragmented genome assembly [13, 14].  
21  
22

23 69 Owing to the advances in sequencing technology, it is possible to obtain high-quality  
24  
25  
26 70 genome assembly that can provide new insights into the understanding of the organisms.  
27  
28  
29 71 Indeed, transposable elements and lineage-specific genes have never been reported to be  
30  
31  
32 72 identified by using the new improved maize reference genome [15]. By combining new  
33  
34  
35 73 sequencing approaches, Seo et al. [13] discovered clinically relevant structural variants and  
36  
37  
38  
39 74 genes never reported before in updated human genome. New sequencing technologies have  
40  
41  
42 75 also been widely used in Gorillas [16] and Sumatran orangutan [17] that have the closet genetic  
43  
44  
45 76 relationship with humans, and domestic goat [18]. More importantly, a lot of new findings have  
46  
47  
48 77 been reported based on these updated genome assembly. However, a high-quality genome  
49  
50  
51 78 assembly of *R. roxellana* has not been reported yet, lagging the progress of understanding this  
52  
53  
54  
55 79 endangered species.  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



1 80 Here, we report a greatly improved assembly and annotation of the reference genome for  
2  
3  
4 81 *R. roxellana* from through combined five technologies: Pacific Bioscience's single-molecule  
5  
6  
7 82 real-time sequencing (SMRT), Illumina's HiSeq paired-end sequencing (HiSeq), BioNano  
8  
9  
10 83 optical maps (BioNano), 10X Genomics link-reads (10X Genomics) and high-throughput  
11  
12  
13 84 chromosome conformation capture (Hi-C). Also, this is the first Colobine genome sequenced  
14  
15  
16 85 and assembled with both long reads and short reads. The updated genome assembly may allow  
17  
18  
19  
20 86 us to comprehensively understand *R. roxellana*, offering new opportunities in analysing  
21  
22  
23 87 evolutionary history and genetic-specific signatures for this species, which may provide new  
24  
25  
26 88 insights in the conservation of this endangered primate. In addition, this genome with superior  
27  
28  
29  
30 89 continuity and accuracy will provide a new standard reference for Colobine primates.

31  
32  
33 90

34  
35  
36 91

## 37 38 39 92 **Sample collection and sequencing**

40  
41  
42 93 An adult dead male *R. roxellana qinlingensis* used for sequencing was from Louguantai  
43  
44  
45 94 Breeding Centre, Xi'an, Shaanxi province, China, originally from Qinling Mountains. Total  
46  
47  
48 95 genomic DNA was extracted from heart tissue. To acquire a high-quality genome assembly,  
49  
50  
51  
52 96 we applied a combined five sequencing methods. Initially, PacBio's SMRT sequencing was  
53  
54  
55 97 conducted on the SEQUEL platform according to manufactures, after removing adaptors in  
56  
57  
58 98 polymerase reads, resulting a total of 304.84 Gb clean long reads (95.86X coverage). Different  
59  
60  
61  
62  
63  
64  
65

1 99 from PacBio sequencing, paired-end sequencing was performed using an Illumina NovaSeq  
2  
3  
4 100 6000 platform with an insert size of 350 bp. Short reads derived from this step were filtered by  
5  
6  
7 101 SOAPfilter v. 2.2 [19] (a package from SOAPdenovo2) with the following criteria: filtering  
8  
9  
10 102 those reads with adapters, contaminations, N bases more than 10% and low quality, which  
11  
12  
13 103 generated 423.32 Gb sequencing clean reads (133.12X coverage). In addition, a high-quality  
14  
15  
16 104 optical genome map was constructed with Irys platform (BioNano Genomics), from which we  
17  
18  
19 105 acquired 463.75 Gb clean reads (145.83X coverage). Besides, 10X genomic link-reads  
20  
21  
22 106 sequencing was carried out on Illumina Hiseq Xten platform, and 348.41 Gb clean reads  
23  
24  
25 107 (109.56X coverage) were generated in total. Finally, a Hi-C library was prepared and  
26  
27  
28 108 sequenced with an Illumina NovaSeq 6000 platform for chromosome-scale scaffolding of  
29  
30  
31 109 genome assembly. Adapter sequences and low quality reads were discarded by using  
32  
33  
34 110 Cutadapter v1.0 [20], yielding a total of 310.92 Gb clean data (97.77X coverage). Statistics of  
35  
36  
37 111 the sequencing data was detailed in **Table 1**.

## 42 112 ***De novo* assembly of the *R. roxellana* genome**

43  
44  
45 113 Estimation of genome size is helpful to our understanding of *R. roxellana*. Generally, we  
46  
47  
48 114 estimated the genome size of *R. roxellana* with the formula of  $G = (K_{total} - K_{error})/D$ , in which  
49  
50  
51 115  $G$  represents genome size, while  $K_{total}$ ,  $K_{error}$  and  $D$  indicates the total number of k-mers,  
52  
53  
54 116 number of k-mers which caused by sequencing errors and k-mer depth respectively. Finally,  
55  
56  
57 117 109,210,004,556 k-mers were generated, and the peak k-mer depth was 34. Thus, the genome  
58  
59  
60  
61  
62  
63  
64  
65

1 118 size of *R. roxellana* was estimated to be about 3.18 Gb. The distribution of k-mer frequency  
2  
3  
4 119 was shown in **Supplementary Fig. S1**.

7 120 The *de novo* assembly of newly sequenced *R. roxellana* genome was performed in four  
8  
9  
10 121 progressive steps. Firstly, the assembly was conducted with the FALCON assembler (default  
11  
12  
13 122 parameters) [21] with the long reads obtained from the PacBio platform, which mainly includes  
14  
15  
16  
17 123 three steps: 1) detection of overlap and reads correction; 2) detection of overlap between  
18  
19  
20 124 corrected reads; and 3) construction of string graph. Following FALCON step, the string graph  
21  
22  
23 125 assembly was further polished by Quier with long reads [22] and then corrected by Pilon with  
24  
25  
26 126 Illumina short reads [23]. Based on this initial genome assembly, sspace-longreads [19] with  
27  
28  
29 127 default settings was implemented for getting a longer scaffold genome by using PacBio long  
30  
31  
32 128 reads. Despite attempts have been made, scaffolding gaps were still found, those gaps were  
33  
34  
35  
36 129 further closed with the help of PBjelly software under default settings, which generated a  
37  
38  
39 130 phased genome assembly with scaffold N50 of 8.20 Mbp (**Supplementary Table S1**).

42 131 Secondly, a hybrid assembly with scaffold N50 of 9.22 Mbp was constructed on the basis  
43  
44  
45 132 of Bionano optical map data using Bionano Solve3.1 ([www.bionanogenomics.com](http://www.bionanogenomics.com)) with  
46  
47  
48 133 default parameters (**Supplementary Table S2**). Thirdly, 10X genomic linked reads were  
49  
50  
51  
52 134 employed to connect scaffolds from the second step by fragScaff software [24], which has  
53  
54  
55 135 updated the scaffold N50 of genome assembly to 24.09 Mbp (**Supplementary Table S3**).

1 136 Subsequently , those short-reads derived from Illumina were applied to correcting errors due  
2  
3  
4 137 to Burrows-Wheeler Aligner (BWA) [25] and pilon-1.18 [23].  
5  
6

7 138 Finally, to build chromosome-level assembly scaffolds, we mapped the Hi-C reads to the  
8  
9  
10 139 assembled scaffolds with BWA [25]. Then Hi-C data was subsequently applied to cluster, order,  
11  
12  
13 140 and orient scaffolds by Lachesis software [26]. The chromosome-level scaffolds for *R.*  
14  
15  
16 141 *roxellana* allowed us to estimate the interaction frequency between chromosome loci, the  
17  
18  
19  
20 142 interaction heatmap shown in **Fig. 2**.  
21  
22

23 143 These processes together yielded a updated genome assembly of *R. roxellana* with its  
24  
25  
26 144 genome size of 3.04 Gb, contig N50 of 5.72 Mbp and scaffold N50 of 144.56 Mbp (**Table 2**).  
27  
28  
29 145 In comparison, the newly acquired *R. roxellana* reference genome has 100-fold higher  
30  
31  
32 146 contiguity than its previous (contig N50: 5.72 Mb versus 25.5 kb and scaffold N50: 144.56 Mb  
33  
34  
35  
36 147 versus 1.55 Mb) [10]. We suppose that the remarkable improvement in contiguity can be  
37  
38  
39 148 attributed to the longer read length, deeper sequencing depth, properly assembled gaps, and  
40  
41  
42 149 increased sophisticated assembly algorithm.  
43  
44

45 150 To assess the genome assembly accuracy, we aligned the Illumina short reads to the  
46  
47  
48 151 assembly by BWA program [25]. With a ratio number of 99.17%, mapped read covered  
49  
50  
51  
52 152 approximately 99.27% of the assembly (**Supplementary Table S4**). In addition, we estimated  
53  
54  
55 153 the assembly completeness by conducting Benchmarking Universal Single-copy Orthologs  
56  
57  
58 154 (BUSCO) analysis with BUSCO V3.0 [27]. Among the 4,104 mammalian BUSCOs, 94% was  
59  
60  
61  
62  
63  
64  
65

1 155 detected in the genome assembly (**Supplementary Table S5**). The assembly completeness was  
2  
3  
4 156 also checked by core eukaryotic gene-mapping approach (CEGMA) [28]. The results showed  
5  
6  
7 157 that 93.95% (233 of 248) conserved genes were found in our genome assembly  
8  
9  
10 158 (**Supplementary Table S6**). Together, these analyses indicated a high accuracy and  
11  
12  
13  
14 159 completeness of our genome assembly.  
15

16  
17 160

## 18 19 20 161 **Identification of repeat elements**

21  
22  
23 162 Repeat sequences occupy a large proportion of the genome sequences. Thus, it is  
24  
25  
26 163 necessary for us to identified those repeat elements. In our study, we combined homolog based  
27  
28  
29 164 and *de novo* based approach to predict and classify repeat elements. As for the homolog  
30  
31  
32  
33 165 approach, we searched transposable elements from the RepBase database [29] with  
34  
35  
36 166 RepeatMasker v4.0.6 (<http://www.repeatmasker.org/>) and RepeatProteinMask (implemented  
37  
38  
39 167 in RepeatMasker). The *de novo* method was employed with RepeatModeler V1.0.11 [30],  
40  
41  
42 168 RepeatMasker v4.0.6 and Tandem Repeat Finder (TRF) (Version 4.07b) [31]. We merged the  
43  
44  
45 169 findings from both methods. Results showed that 45.43% of the genome was predicted as  
46  
47  
48  
49 170 repeat elements (**Supplementary Table S7**). A closer investigation indicated that the largest  
50  
51  
52 171 category of repeat elements in the species is the short (SINEs) and long (LINEs) interspersed  
53  
54  
55 172 nuclear elements. The detailed categories of repeat elements are summarized in  
56  
57  
58 173 **Supplementary Table S8**.

1 174  
2  
3  
4 175  
5  
6  
7 176  
8  
9  
10 177  
11  
12  
13 178  
14  
15  
16 179  
17  
18  
19  
20 180  
21  
22  
23 181  
24  
25  
26 182  
27  
28  
29 183  
30  
31  
32  
33 184  
34  
35  
36 185  
37  
38  
39 186  
40  
41  
42 187  
43  
44  
45 188  
46  
47  
48 189  
49  
50  
51  
52 190  
53  
54  
55 191  
56  
57  
58 192  
59  
60  
61  
62  
63  
64  
65

## **Non-coding RNA prediction**

Non-coding RNA consists of several RNAs, as such ribosomal RNA (rRNA), transfer RNA (tRNA), microRNAs (miRNA) and small nuclear RNA (snRNA). This RNA group mainly plays a regulation role in biological processes. In our study, we detected rRNA from a Human rRNA database with BLASTN command, and the E-value was set as 1E-10. Similarly, miRNAs and snRNAs were searched against the Rfam database [32] with INFERNAL 1.1rc4 [33]. The tRNAs were predicted by tRNAscan-SE 1.3.1 software [34]. The numbers of rRNA, miRNA, snRNA and tRNA were 608, 17,813, 3,656 and 460, respectively in the genome of the species (**Supplementary Table S9**).

## **Gene prediction and functional annotation**

We combined prediction methods based on *de novo*, homolog prediction and transcriptome data to estimate genes. As for *ab initio* based prediction, a total of five programs, namely Augustus v. 3.2.2 [35], GlimmeHMM v. 3.0.1 [36], GENSCAN [37], GENEID [38] and SNAP V2013-11-29 [39] were employed to predict protein-coding genes. Subsequently, we used the homolog-based prediction approach. Protein sequences from five homolog species (*Homo sapiens*, *Gorilla gorilla*, *Macaca mulatta*, *Rhinopithecus bieti*, *Rhinopithecus roxellana hubeiensis*) were downloaded from Ensemble Release 75

1 193 (<http://www.ensembl.org/info/data/ftp/index.html>), and used to perform TBLASTN blast  
2  
3  
4 194 against the repeat-masked genome sequences [40]. The related homologous genome sequences  
5  
6  
7 195 were then annotated to the matching proteins by GeneWise 2.4.1 [41]. Finally, we estimated  
8  
9  
10 196 genes based on transcriptome data. During this process, high-quality RNAs from heart and skin  
11  
12  
13 197 tissue were sequenced by an Illumina Novaseq 6000 platform. RNA-seq reads were assembled  
14  
15  
16 198 with trinityrnaseq-2.1.1 [42]. The assembled transcript sequences were aligned to the *R.*  
17  
18  
19  
20 199 *roxellana* genome by Assemble Spliced Alignment (PASA) [43] with default parameters. In  
21  
22  
23 200 addition, we estimated the expression levels of transcripts by Tophat 2.0.13 [44] and Cufflinks  
24  
25  
26 201 [45].  
27  
28

29 202 The genes predicted from those three approaches were merged with EvidenceModeler  
30  
31  
32 203 [46]. Furthermore, untranslated regions and alternative splicing of those predicted gene sets  
33  
34  
35 204 were further checked by PASA with the help of transcriptome data [43]. Finally, a total of  
36  
37  
38  
39 205 22497 genes were predicted for the assembly genome of *R. roxellana* (**Table 3**), and each of  
40  
41  
42 206 them consisted of 7.71 exons on average. The detailed results generated during the gene  
43  
44  
45 207 prediction process were shown in **Table 3**. And, the gene prediction evidence based on  
46  
47  
48 208 different methods were shown in **Fig. 3**. In addition, we made a comparison between the *R.*  
49  
50  
51 209 *roxellana qinlingensis* and other mammals, suggesting a comparable pattern of the genome  
52  
53  
54 210 assembly for *R. roxellana qinlingensis* (**Supplementary Fig. S2**).  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 211 To have a better understanding the biological functions of those predicted genes, they were  
2  
3  
4 212 annotated with several databases including NCBI nonredundant protein database (NR),  
5  
6  
7 213 SwissProt [47], Kyoto Encyclopedia of Genes and Genomes (KEGG) [48], InterPro [49], Pfam  
8  
9  
10 214 [50] and GO database [51]. In total, 22053 genes (98.42%) were functionally annotated  
11  
12  
13  
14 215 **(Supplementary Table S10).**

15  
16  
17 216

### 20 217 **Phylogenetic relationship analysis and gene family estimation**

23 218 Coding regions and protein sequences of 11 representative mammals were downloaded  
24  
25  
26 219 from Ensemble (Ensemble Release 75). The longest transcript was chosen if genes possess  
27  
28  
29 220 many transcript isoforms. Treefam [52] approach was adopted to estimate gene families.  
30  
31  
32 221 Following all-to-all blast, a total of 17,560 gene families were identified. We reconstructed the  
33  
34  
35  
36 222 phylogenetic relationship between *R. roxellana* and other mammals based on four-fold  
37  
38  
39 223 degenerate sites extracted from the 5,418 single-copy gene families. Phym1 (version 3.2) [53]  
40  
41  
42 224 was employed to construct a maximum-likelihood tree under the GTR + gamma model that  
43  
44  
45 225 was inferred from JMODELTEST (version 2.1.10) [54]. Furthermore, we estimated the  
46  
47  
48 226 divergence time with MCMCTREE in PAML [55]. MCMCTREE was performed on the basis  
49  
50  
51  
52 227 of bayesian method and the fossil calibration times from timetree were used as input. The  
53  
54  
55 228 reconstructed phylogeny confirmed the close relationship between *R. rollexana* and *M. mulatta*.

56  
57  
58  
59  
60  
61  
62  
63  
64  
65



1 229 Moreover, we estimated that *R. rollexana* and *M. mulatta* diverged approximately 13.4 million  
2  
3  
4 230 years ago (Mya) (**Fig. 4**).

5  
6  
7 231 To have a better understanding the evolutionary history of *R. roxellana*, we estimated the  
8  
9  
10 232 expansion and contraction of gene family in *R. roxellana* by using CAFE 3.0 [56]. A gene  
11  
12  
13 233 family with *p*-value less than 0.05 was considered for further analysis. As a result, 993 and  
14  
15  
16 234 2,745 gene families were expanded and contracted in *R. roxellana* genome, respectively (**Fig.**  
17  
18  
19  
20 235 **4**). Its genome showed substantial expansion of gene families which are mainly related to  
21  
22  
23 236 hemoglobin complex, energy metabolisms and oxygen transport (**Supplementary Table S11**).

## 24 25 26 237 27 28 29 238 **Conclusion**

30  
31  
32 239 In this study, we generated a high-quality genome assembly of the golden snub-nosed  
33  
34  
35 240 monkey (*R. roxellana*) by using five advanced technologies. This will be helpful to investigate  
36  
37  
38  
39 241 the origin and evolutionary history of snub-nosed monkey. In addition, the genome may lay a  
40  
41  
42 242 foundation to survey the mechanisms about the formation of distinct characters and understand  
43  
44  
45 243 the unique multilevel societies in *R. roxellana*. Also, such genome may provide new insights  
46  
47  
48 244 for amending the conservation strategies and management of this endangered species.  
49  
50  
51 245 Furthermore, this genome with superior continuity and accuracy can provide a new standard  
52  
53  
54  
55 246 reference for Colobine primates.

## 56 57 58 247 **Declarations**

1 248 **Availability of supporting data**

2  
3  
4 249 The raw data discussed in this publication have been deposited in NCBI's short read archive  
5  
6  
7 250 under the accession number PRJNA524949. Supporting data are available in the GigaDB  
8  
9  
10 251 database.

13 252 **Competing interests**

16 253 The authors declare that they have no competing interests.

20 254 **Funding**

23 255 This work was financially supported by Strategic Priority Research Program of the Chinese  
24  
25  
26 256 Academy of Sciences (XDB31020302), the National Natural Science Foundation of China  
27  
28  
29 257 (31622053, 31730104), the Promotional project for Innovation team, the Department of  
30  
31  
32 258 Science and Technology of Shaanxi Prov. China (2018TD-017), and the National Key  
33  
34  
35  
36 259 Programme of Research and Development, the Ministry of Science and Technology of China  
37  
38  
39 260 (2016YFC0503200).

42 261 **Abbreviations**

45 262 Gb: gigabase; kb: kilobase; Mb: megabase; PE: paired-end; PacBio: Pacific Biosciences;  
46  
47  
48 263 SMRT: single molecule real-time sequencing; Hi-C: high-throughput chromosome  
49  
50  
51 264 conformation capture; BUSCO: Benchmarking Universal Single-copy Orthologs; GEGMA:  
52  
53  
54  
55 265 core eukaryotic gene-mapping approach; GO: gene ontology; TFS: transposable element;  
56  
57  
58 266 TRF: Tandem Repeat Finder; SINEs: Short interspersed nuclear elements; LINEs: long  
59  
60  
61

1 267 interspersed nuclear elements; PASA: genome by Assemble Spliced Alignment; NR: NCBI  
2  
3  
4 268 nonredundant protein database; KEGG: Kyoto Encyclopedia of Genes and Genomes. Mya:  
5  
6  
7 269 million years ago.  
8  
9

10  
11 270 **Author contributions**  
12  
13

14 271 X.G.Q. conceived and designed the project, L.W., J.W. contributed to the work on genomic  
15  
16  
17 272 sequencing and performing data analyses. B.G.L. helped with sample collection. L.W., J.W.  
18  
19  
20 273 and X.G.Q. wrote the manuscript. All authors provided input for the paper and approved the  
21  
22  
23 274 final version.  
24  
25  
26

27 275  
28

29  
30 276 **Acknowledgements**  
31  
32

33 277 We thank Mr. Yiliang Xu, Mr. Qiqi Liang, Mrs. Yue Xie from Novogene for their technical  
34  
35  
36 278 support. Mr. Xuanmin Guang and Mr. Chi Zhang from BGI for their assistance in data analysis.  
37  
38  
39 279 We thank to Mr. Ruliang Pan for his gracious help polishing the language. We are also grateful  
40  
41  
42 280 to Mr. Yinghu Lei from Louguantai Breeding Center, Dr. Zhipang Huang, and Dr. Pei Zhang  
43  
44  
45 281 from Northwest University for their helping with the sampling collection. We specially  
46  
47  
48 282 appreciate Prof. Zhengbing Wang, and Prof. Jiang Chang from Discipline Development  
49  
50  
51 283 Department of Northwest University for their support. This study was fundamentally supported  
52  
53  
54 284 by Discipline Construction Project of Northwest University.  
55  
56  
57

58 285  
59  
60  
61  
62  
63  
64  
65

1 286 **Figures and tables**

2  
3  
4 287

5  
6  
7 288 **Figure legends:**

8  
9  
10 289  
11  
12  
13 290 **Fig. 1. The photo of *R. roxellana* taken in the Qinling mountains.**

14  
15  
16 291 **Fig. 2. Hi-C heatmap between chromosome loci throughout the genome.** Hi- C interactome  
17  
18  
19  
20 292 within and among chromosomes of *R. roxellana* (Chr1–Chr22).

21  
22  
23 293 **Fig. 3. The gene prediction evidence based on different methods.** (a). Number of the genes  
24  
25  
26 294 estimated by the prediction approaches based on de novo (blue color), homolog prediction  
27  
28  
29 295 (pink color) and RNA\_seq data (green color). The rna\_0.5, denovo\_0.5 and homolog\_0.5  
30  
31  
32 296 indicates those genes predicted with an overlap are larger than 50% in each method; (b)  
33  
34  
35  
36 297 Number of the genes shown in combination with the prediction approaches detailed in fig 2a  
37  
38  
39 298 and the expression level standard (rpkm). The rna\_0.5, denovo\_0.5, homolog\_0.5 indicates that  
40  
41  
42 299 those genes predicted with an overlap are larger than 50% in each method, while rpkm>1  
43  
44  
45 300 indicates those genes with an expression level larger than 1.

46  
47  
48 301 **Fig. 4. The phylogenetic relationships of *R. roxellana* and other mammals and Gene**  
49  
50  
51 302 **family analysis in *R. roxellana* genome.** Phylogenetic relationship was inferred from 5418  
52  
53  
54  
55 303 single-copy gene families. All nodes received 100% support values. The estimated divergence  
56  
57  
58 304 times are indicated near the nodes. The images in the figure are credited as “Illustrations

1 305 copyright 2013 Stephen D. Nash / IUCN SSC Primate Specialist Group. Used with permission”.

2  
3

4 306 MYA: million years ago. The numbers on each branch correspond to the numbers of gene

5  
6

7 307 families that have expanded (red) and contracted (green) in mammalian genome. MRCA: most

8  
9

10 308 recent common ancestor.

11  
12

13 309

14 310

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

311 **Table 1. Reads generated from five different sequencing methods**

Pair-end libraries	Insert size (bp)	Total clean data (Gb)	Read length (bp)	Sequence coverage (X)
Illumina reads	350	423.32	150	133.12
Pacbio reads	20 k	304.84	\	95.86
10X Genomics	500 -700	348.41	150	109.56
Bionano	\	463.75	\	145.83
Hi-C	350	310.92	\	97.77
Total	\	1,851.24	\	582.15

312 Note: The sequence coverage was calculated with an estimated genome size of 3.18 Gb. The  
 313 sign of backslash indicates that the insert size was absent.

314

315 **Table 2. The final genome assembly statistics of *R. roxellana***

Sample ID	length		number	
	Contig <sup>a</sup> (bp)	Scaffold (bp)	Contig <sup>b</sup>	Scaffold
Total	3,038,184,325	3,038,467,325	6,099	3,269
Max	30,757,641	206,558,726	\	\
Number >= 2000	\	\	5,708	2,879
N50	5,723,610	144,559,847	151	9
N60	4,241,389	141,075,955	211	11
N70	3,173,235	135,203,321	292	14
N80	2,063,823	118,350,466	408	16
N90	896,517	83,045,532	622	19

316 Note: <sup>a</sup> Contig after scaffolding. The sign of backslash indicates that the length/number was  
 317 absent.

318

319

320 **Table 3. Summary of predicted protein-coding genes and their characteristics**

Gene set	Number	Average transcript length (bp)	Average CDS length (bp)	Average intron length (bp)	Average exon length (bp)	Average exons per gene
Augustus	32,928	23,441	1,052	196	5,112	5.38
GlimmerHMM	618,957	4,204	404	166	2,654	2.43
<i>De novo</i> SNAP	97,298	49,851	755	144	1,1597	5.23
Geneid	36,863	35,242	1,035	188	7,615	5.49
Genscan	50,419	40,635	1,137	167	6,800	6.81
Ggo	25,281	19,893	1,055	184	3,971	5.74
Hsa	38,444	14,763	826	182	3,942	4.54
Homolog Mmu	21,959	29,709	1,470	187	4,123	7.85
Rbi	25,320	25,685	1,387	196	3,991	7.09
Rro	24,121	28,439	1,420	185	4,043	7.68
RNASeq PASA	66,620	28,449	1,219	164	4,247	7.41
Cufflinks	73,199	31,497	2,737	409	5,052	6.69
EVM	30,102	22,298	1,098	182	4,199	6.05
Pasa-update*	29,403	27,638	1,180	181	4,782	6.53
Final set*	22,497	34,153	1,369	178	4885	7.71

321 Note: Pasa-update\* indicates only the UTRs (Untranslated regions) were considered during  
 322 the filter process, and other regions were not included. Final set\* indicates the results were  
 323 acquired following the Pasa-update process, with the criteria of the longest isoform was chosen  
 324 if there were multiple splicing isoforms, and the redundant single exons were also discarded.



- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 327 **Supplementary files:**
- 328 **Supplementary Fig. S1.** Genome size estimation using k-mer method
- 329 **Supplementary Fig. S2.** The comparison of each element in the genome of homologous  
330 species
- 331 **Supplementary Table S1.** The results of assembly with PacBio long reads and gap filling
- 332 **Supplementary Table S2.** The results of assembly with Bionano optical map data
- 333 **Supplementary Table S3.** The results of assembly with 10X Genomics link reads
- 334 **Supplementary Table S4.** The mapping rate of reads and coverage of assembled genome  
335 with BWA
- 336 **Supplementary Table S5.** Assessment results by using BUSCO annotation
- 337 **Supplementary Table S6.** The completeness test results of assembled genome with CEGMA  
338 software
- 339 **Supplementary Table S7.** Results of repeats elements predictions from the genome  
340 assembly
- 341 **Supplementary Table S8.** The results of TEs elements predicted from the genome assembly
- 342 **Supplementary Table S9.** Summary of predicted RNAs and their characteristics

1 343 **Supplementary Table S10.** The functional annotation of the genes predicted from *R. roxellana*

2

3

4 344 genome

5

6

7 345 **Supplementary Table S11.** The GO annotation results of expansion gene families in *R.*

8

9

10 346 *roxellana* genome

11

12

13

14 347

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

348 **References**

- 1  
2 349 1. Li BG, Pan RL, Oxnard CE. Extinction of snub-nosed monkeys in China during the  
3  
4 350 past 400 years. *Int J Primatol* 2002; **23**:1227-1244.
- 5 351 2. Luo MF, Liu ZJ, Pan HJ, et al. Historical geographic dispersal of the golden snub-  
6  
7 352 nosed monkey (*Rhinopithecus roxellana*) and the influence of climatic oscillations.  
8  
9 353 *Am J Primatol* 2012; **74**:91-101.
- 10 354 3. Fang G, Li M, Liu X-J, et al. Preliminary report on Sichuan golden snub-nosed  
11  
12 355 monkeys (*Rhinopithecus roxellana roxellana*) at Laohegou Nature Reserve, Sichuan,  
13 356 China. *Sci Rep* 2018; **8**:16183.
- 14 357 4. Qi XG, Li BG, Garber PA, et al. Social dynamics of the golden snub-nosed monkey  
15  
16 358 (*Rhinopithecus roxellana*): female transfer and one-male unit succession. *Am J*  
17  
18 359 *Primatol* 2009; **71**:670-679.
- 19 360 5. Li H, Meng S-J, Men Z-M, et al. Genetic diversity and population history of golden  
20  
21 361 monkeys (*Rhinopithecus roxellana*). *Genetics* 2003; **164**:269-275.
- 22 362 6. Qi X-G, Garber PA, Ji W, et al. Satellite telemetry and social modeling offer new  
23  
24 363 insights into the origin of primate multilevel societies. *Nat Commun* 2014; **5**:5296.
- 25 364 7. Schnable PS, Ware D, Fulton RS, et al. The B73 Maize Genome: Complexity,  
26  
27 365 Diversity, and Dynamics. *Science* 2009; **326**:1112-1115.
- 28 366 8. Seim I, Fang X, Xiong Z, et al. Genome analysis reveals insights into physiology and  
29  
30 367 longevity of the Brandt's bat *Myotis brandtii*. *Nature Communications* 2013; **4**:2212.
- 31 368 9. Valenzano DR, Benayoun BA, Singh PP, et al. The African turquoise killifish genome  
32  
33 369 provides insights into evolution and genetic architecture of lifespan. *Cell* 2015;  
34  
35 370 **163**:1539-1554.
- 36 371 10. Zhou X, Wang B, Pan Q, et al. Whole-genome sequencing of the snub-nosed monkey  
37  
38 372 provides insights into folivory and evolutionary history. *Nat Genet* 2014; **46**:1303-  
39  
40 373 1310.
- 41 374 11. Kuang W-M, Ming C, Li H-P, et al. The origin and population history of the  
42  
43 375 endangered golden snub-nosed monkey (*Rhinopithecus roxellana*). *Mol Biol Evol*  
44  
45 376 2018:msy220-msy220.
- 46 377 12. Hong YY, Duo HR, Hong JY, et al. Resequencing and comparison of whole  
47  
48 378 mitochondrial genome to gain insight into the evolutionary status of the Shennongjia  
49  
50 379 golden snub-nosed monkey (SNJ *R-roxellana*). *Ecol Evol* 2017; **7**:4456-4464.
- 51 380 13. Seo JS, Rhie A, Kim J, et al. De novo assembly and phasing of a Korean human  
52  
53 381 genome. *Nature* 2016; **538**:243-247.
- 54 382 14. Chaisson MJ, Wilson RK, Eichler EE. Genetic variation and the de novo assembly of  
55  
56 383 human genomes. *Nat Rev Genet* 2015; **16**:627-640.
- 57 384 15. Jiao YP, Peluso P, Shi JH, et al. Improved maize reference genome with single-  
58  
59 385 molecule technologies. *Nature* 2017; **546**:524-527.

- 1 386 16. Gordon D, Huddleston J, Chaisson MJP, et al. Long-read sequence assembly of the  
2 387 gorilla genome. *Science* 2016; **352**:aae0344.
- 3 388 17. Kronenberg ZN, Fiddes IT, Gordon D, et al. High-resolution comparative analysis of  
4 389 great ape genomes. *Science* 2018; **360**:eaar6343.
- 5 390 18. Bickhart DM, Rosen BD, Koren S, et al. Single-molecule sequencing and chromatin  
6 391 conformation capture enable de novo reference assembly of the domestic goat  
7 392 genome. *Nat Genet* 2017; **49**:643-650.
- 8 393 19. Luo R, Liu B, Xie Y, et al. SOAPdenovo2: an empirically improved memory-efficient  
9 394 short-read de novo assembler. *Gigascience* 2012; **1**:18.
- 10 395 20. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing  
11 396 reads. *EMBnetjournal* 2011; **17**:10-12.
- 12 397 21. Pendleton M, Sebra R, Pang A, et al. Assembly and diploid architecture of an  
13 398 individual human genome via singlemolecule technologies. *Nat Methods* 2015;  
14 399 **12**:780-786.
- 15 400 22. Chin CS, Alexander DH, Marks P, et al. Nonhybrid, finished microbial genome  
16 401 assemblies from long-read SMRT sequencing data. *Nat Methods* 2013; **10**:563-+.
- 17 402 23. Walker BJ, Abeel T, Shea T, et al. Pilon: An Integrated Tool for Comprehensive  
18 403 Microbial Variant Detection and Genome Assembly Improvement. *PLoS One* 2014;  
19 404 **9**.
- 20 405 24. Adey A, Kitzman JO, Burton JN, et al. In vitro, long-range sequence information for  
21 406 de novo genome assembly via transposase contiguity. *Genome Res* 2014; **24**:2041-  
22 407 2049.
- 23 408 25. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler  
24 409 transform. *Bioinformatics* 2009; **25**:1754-1760.
- 25 410 26. Burton JN, Adey A, Patwardhan RP, et al. Chromosome-scale scaffolding of de novo  
26 411 genome assemblies based on chromatin interactions. *Nat Biotechnol* 2013; **31**:1119-  
27 412 1125.
- 28 413 27. Simao FA, Waterhouse RM, Ioannidis P, et al. BUSCO: assessing genome assembly  
29 414 and annotation completeness with single-copy orthologs. *Bioinformatics* 2015;  
30 415 **31**:3210-3212.
- 31 416 28. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately annotate core genes in  
32 417 eukaryotic genomes. *Bioinformatics* 2007; **23**:1061-1067.
- 33 418 29. Jurka J, Kapitonov VV, Pavlicek A, et al. Repbase Update, a database of eukaryotic  
34 419 repetitive elements. *Cytogenet Genome Res* 2005; **110**:462-467.
- 35 420 30. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large  
36 421 genomes. *Bioinformatics* 2005; **21 Suppl 1**:i351-358.
- 37 422 31. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic  
38 423 Acids Res* 1999; **27**:573-580.

- 1 424 32. Griffiths-Jones S, Moxon S, Marshall M, et al. Rfam: annotating non-coding RNAs in  
2 425 complete genomes. *Nucleic Acids Res* 2005; **33**:D121-D124.
- 3 426 33. Nawrocki EP, Kolbe DL, Eddy SR. Infernal 1.0: inference of RNA alignments (vol  
4 427 25, pg 1335, 2009). *Bioinformatics* 2009; **25**:1713-1713.
- 5 428 34. Lowe TM, Eddy SR. tRNAscan-SE: a program for improved detection of transfer  
6 429 RNA genes in genomic sequence. *Nucleic Acids Res* 1997; **25**:955-964.
- 7 430 35. Stanke M, Keller O, Gunduz I, et al. AUGUSTUS: ab initio prediction of alternative  
8 431 transcripts. *Nucleic Acids Res* 2006; **34**:W435-W439.
- 9 432 36. Majoros WH, Pertea M, Salzberg SL. TigrScan and GlimmerHMM: two open source  
10 433 ab initio eukaryotic gene-finders. *Bioinformatics* 2004; **20**:2878-2879.
- 11 434 37. Burge C, Karlin S. Prediction of complete gene structures in human genomic DNA. *J*  
12 435 *Mol Biol* 1997; **268**:78-94.
- 13 436 38. Guigo R. Assembling genes from predicted exons in linear time with dynamic  
14 437 programming. *J Comput Biol* 1998; **5**:681-702.
- 15 438 39. Korf I. Gene finding in novel genomes. *BMC Bioinformatics* 2004; **5**:59.
- 16 439 40. Kent WJ. BLAT - The BLAST-like alignment tool. *Genome Res* 2002; **12**:656-664.
- 17 440 41. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome Res* 2004;  
18 441 **14**:988-995.
- 19 442 42. Grabherr MG, Haas BJ, Yassour M, et al. Full-length transcriptome assembly from  
20 443 RNA-Seq data without a reference genome. *Nat Biotechnol* 2011; **29**:644-U130.
- 21 444 43. Haas BJ, Delcher AL, Mount SM, et al. Improving the Arabidopsis genome  
22 445 annotation using maximal transcript alignment assemblies. *Nucleic Acids Res* 2003;  
23 446 **31**:5654-5666.
- 24 447 44. Kim D, Pertea G, Trapnell C, et al. TopHat2: accurate alignment of transcriptomes in  
25 448 the presence of insertions, deletions and gene fusions. *Genome Biol* 2013; **14**:R36.
- 26 449 45. Trapnell C, Roberts A, Goff L, et al. Differential gene and transcript expression  
27 450 analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* 2012;  
28 451 **7**:562-578.
- 29 452 46. Haas BJ, Salzberg SL, Zhu W, et al. Automated eukaryotic gene structure annotation  
30 453 using EVIDENCEModeler and the program to assemble spliced alignments. *Genome*  
31 454 *Biol* 2008; **9**:R7.
- 32 455 47. Bairoch A, Apweiler R, Wu CH, et al. The Universal Protein Resource (UniProt).  
33 456 *Nucleic Acids Res* 2005; **33**:D154-D159.
- 34 457 48. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic*  
35 458 *Acids Res* 2000; **28**:27-30.
- 36 459 49. Mitchell AL, Attwood TK, Babbitt PC, et al. InterPro in 2019: improving coverage,  
37 460 classification and access to protein sequence annotations. *Nucleic Acids Res* 2019;  
38 461 **47**:D351-D360.

1 462 50. Finn RD, Coghill P, Eberhardt RY, et al. The Pfam protein families database: towards  
2 463 a more sustainable future. *Nucleic Acids Res* 2016; **44**:D279-D285.

3 464 51. Ashburner M, Ball CA, Blake JA, et al. Gene Ontology: tool for the unification of  
4 465 biology. *Nat Genet* 2000; **25**:25-29.

5 466 52. Li H, Coghlan A, Ruan J, et al. TreeFam: a curated database of phylogenetic trees of  
6 467 animal gene families. *Nucleic Acids Res* 2006; **34**:D572-D580.

7 468 53. Guindon S, Delsuc F, Dufayard J-F, et al: Estimating maximum likelihood  
8 469 phylogenies with PhyML. In *Bioinformatics for DNA sequence analysis*. Springer;  
9 470 2009: 113-137

10 471 54. Posada D. jModelTest: phylogenetic model averaging. *Mol Biol Evol* 2008; **25**:1253-  
11 472 1256.

12 473 55. Yang Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Mol Biol Evol*  
13 474 2007; **24**:1586-1591.

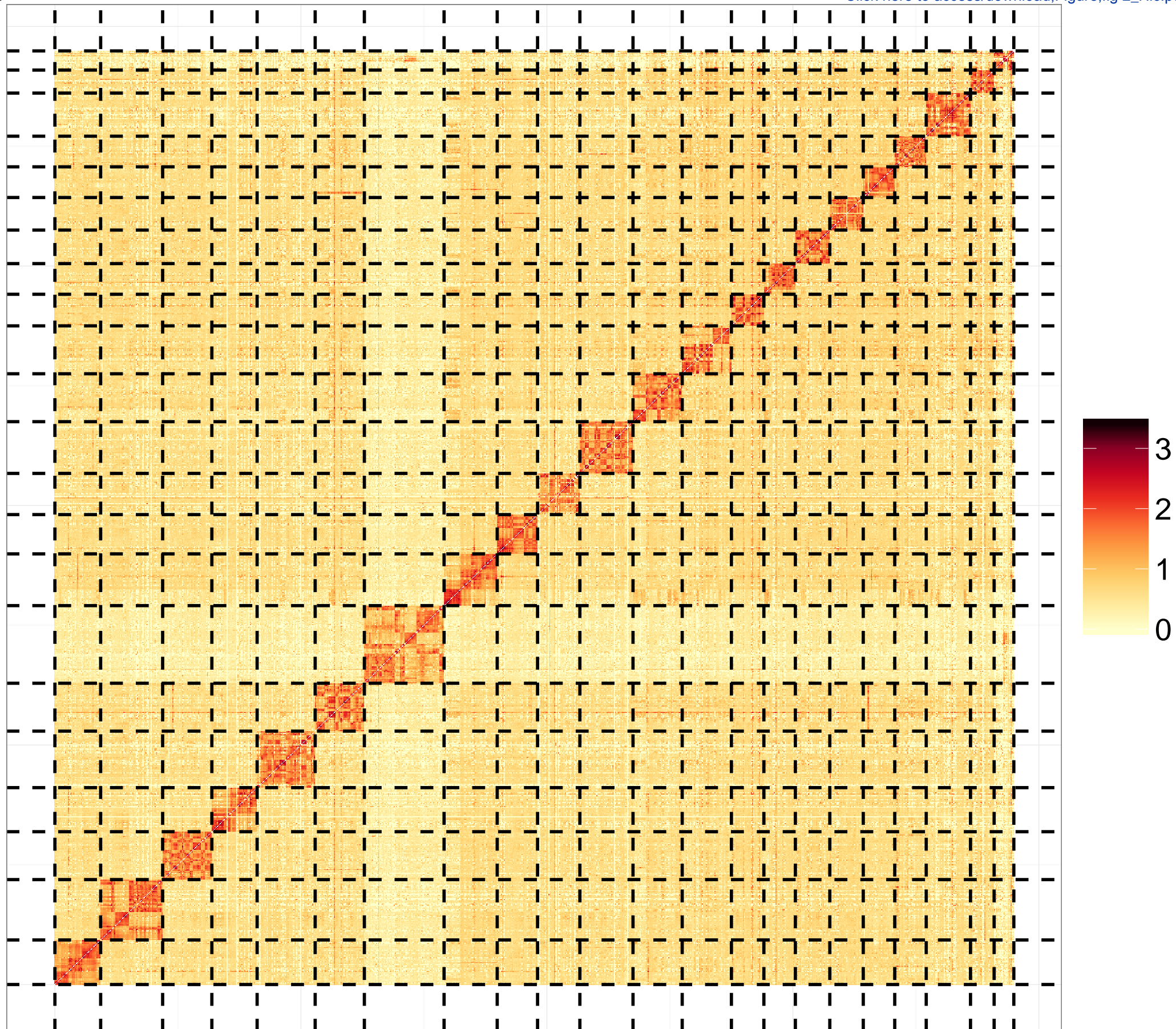
14 475 56. De Bie T, Cristianini N, Demuth JP, et al. CAFE: a computational tool for the study  
15 476 of gene family evolution. *Bioinformatics* 2006; **22**:1269-1271.

16 477

17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

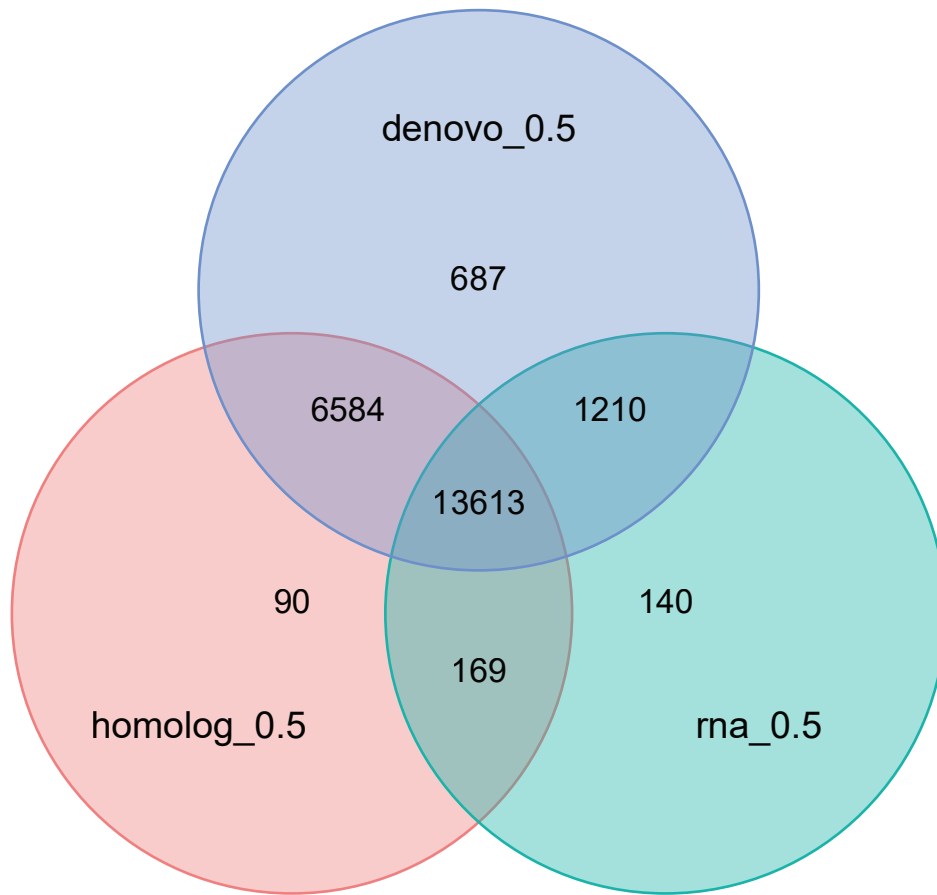


Fig 2





a



b

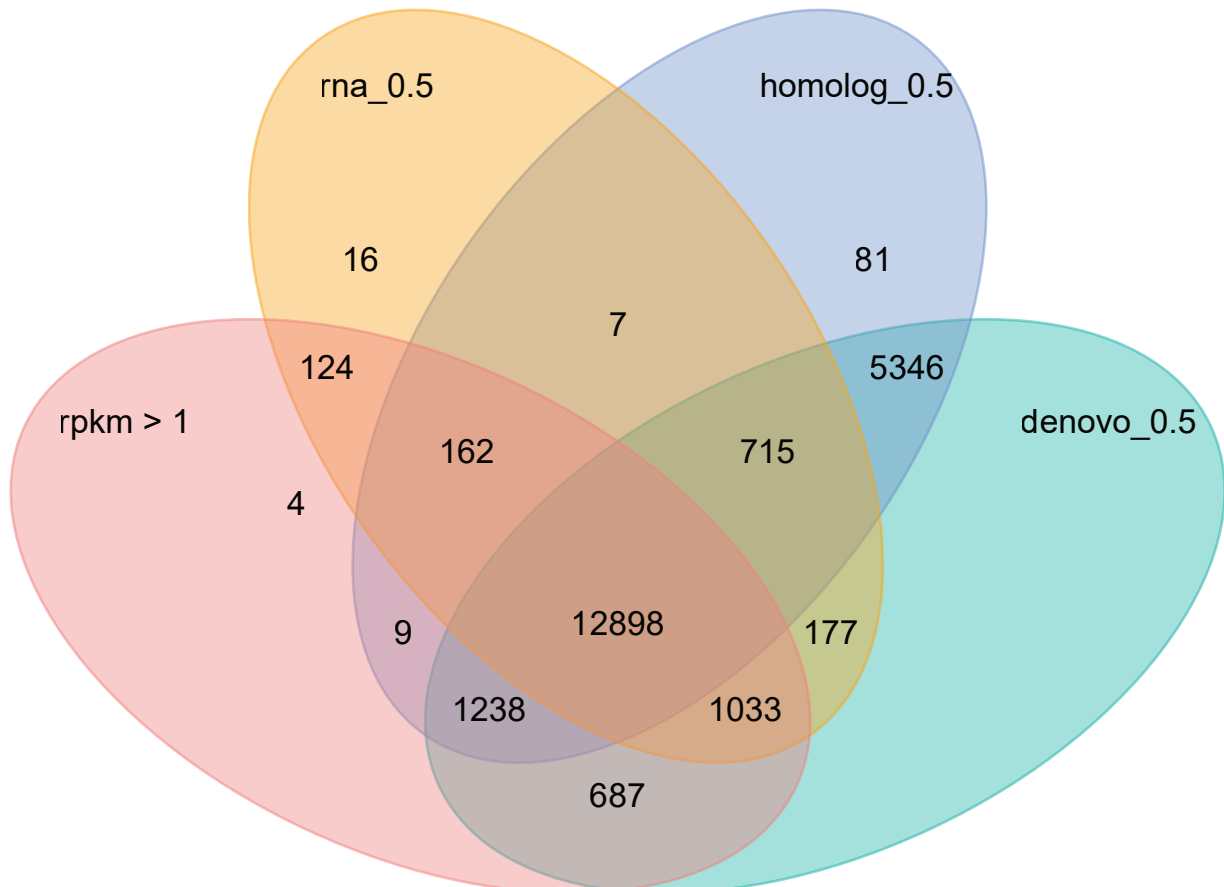
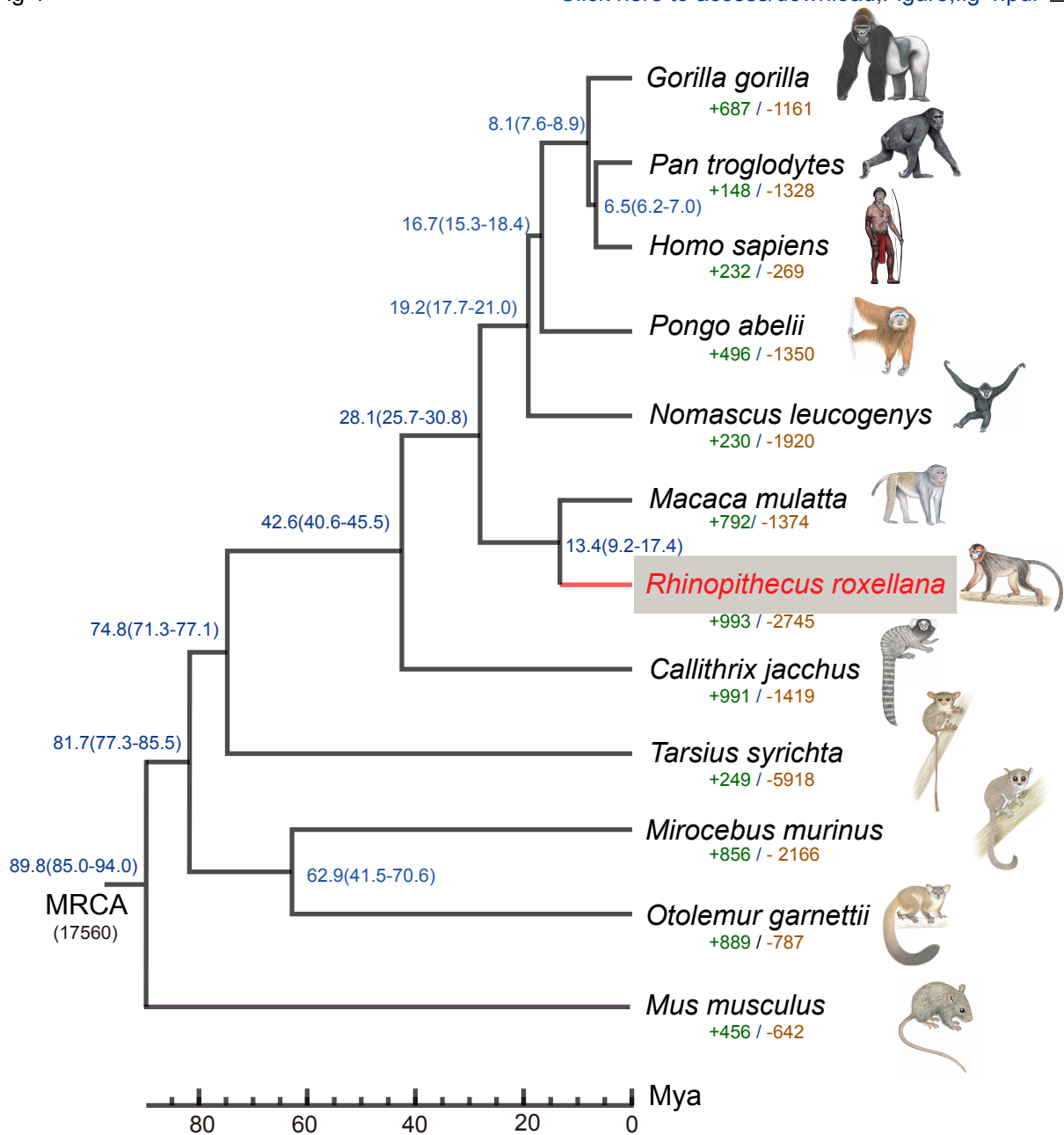


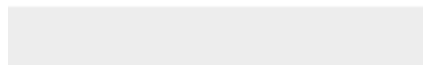
Fig 4

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





Click here to access/download  
**Supplementary Material**  
SI\_V1.docx



March 12, 2019

Dear *GigaScience* editor,

Thanks a lot for having reviewed our manuscript titled “**A high-quality genome assembly of the endangered golden snub-nosed monkey *Rhinopithecus roxellana***” for consideration of publication in *GigaScience*. Now we have revised the manuscript according to the reviewer’s comments. Most explanations regarding the revisions of our manuscript are as follows.

1 The raw data discussed in this publication have been deposited in NCBI’s short read archive under the accession number PRJNA524949. Supporting data are available in the GigaDB database.

2 We have attached a photo of *Rhinopithecus roxellana* taken in the Qinling mountains in the attached file “Figures” (fig 1).

3 We have improved the language by a native speaker.

4 We have updated the abstract, introduction and conclusion.

4 We have improved the introduction to highlight the reason why we sequence the species (On Page 6, Line 84-89).

5 We demonstrated our genome improvement by comparing with previously reported (On Page 9, Line 143-149).

As a result, we believe that our work will be of general interest to a broad scientific audience including evolutionary biologists, phylogeneticists, geneticists, zoologists, ecologists, and popular media. Thank you for your time and consideration.

Yours sincerely,

Xiao-Guang Qi

Shaanxi Key Laboratory for Animal Conservation

College of Life Sciences

Northwest University

Email: qixg@nwu.edu.cn